

# PEOPLE + DATA + COMPUTATION

Working Productively with Data

Jeffrey Heer  
Joe Hellerstein



**LIFE**

The Future of Data Analysis, Tukey 1962



The last few decades have seen the rise of formal theories of statistics, "legitimizing" variation by confining it by assumption to random sampling, often assumed to involve tightly specified distributions, and restoring the appearance of security by emphasizing narrowly optimized techniques and claiming to make statements with "known" probabilities of error.

LIFE



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

LIFE



Exposure, the effective laying open of the data to display the unanticipated, is to us a major portion of data analysis. Formal statistics has given almost no guidance to exposure; indeed, 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



It is too much to ask for close and effective guidance for data analysis from any highly formalized structure, either now or in the near future.

Data analysis can gain much from formal statistics, but only if the connection is kept adequately loose.

LIFE



LIFE

The Future of Data Analysis, Tukey 1962



Visualization

Acquisition

Cleaning

Integration

Visualization

Modeling

Presentation

Dissemination

Acquisition



Cleaning



Integration



Visualization



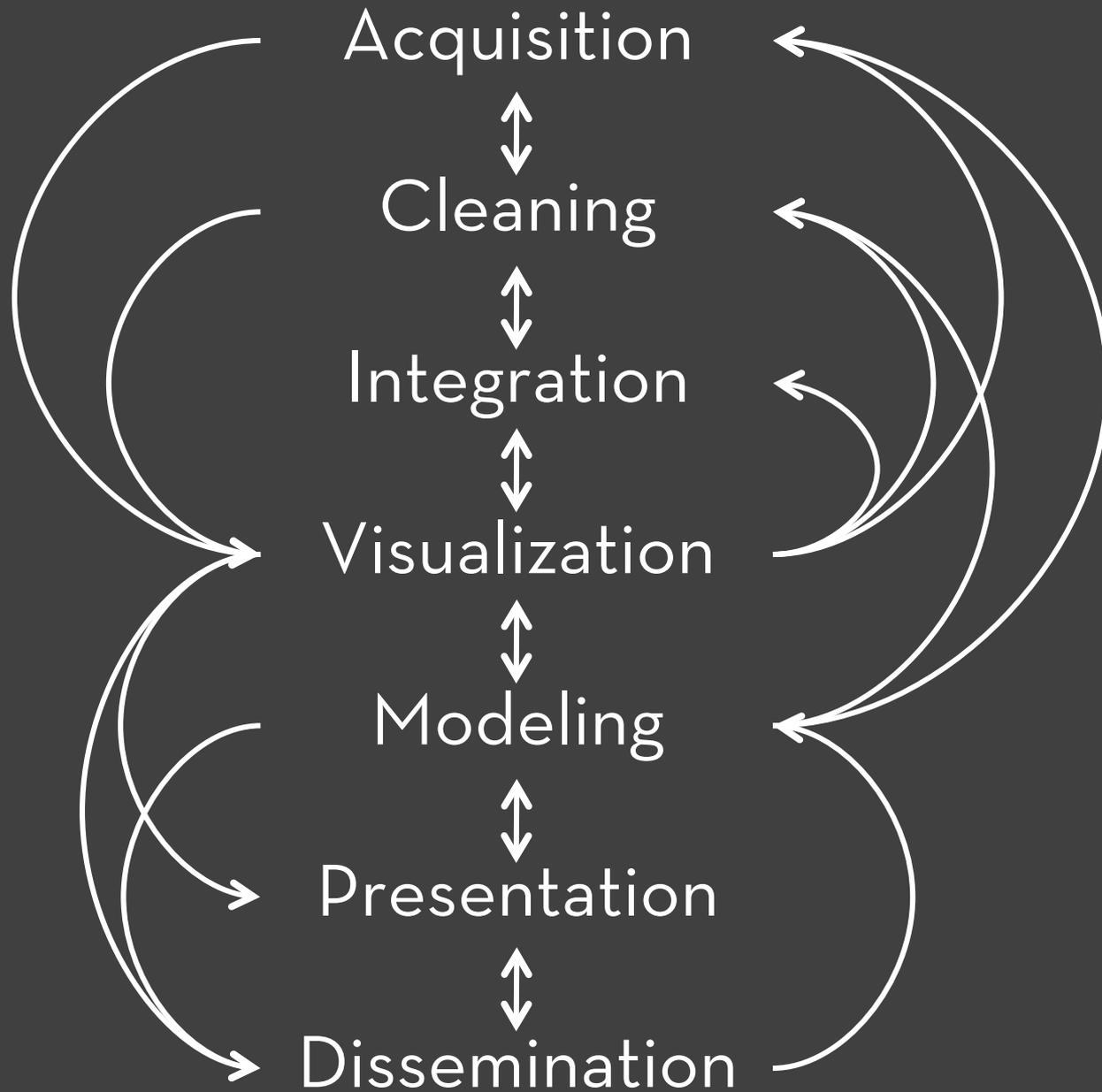
Modeling

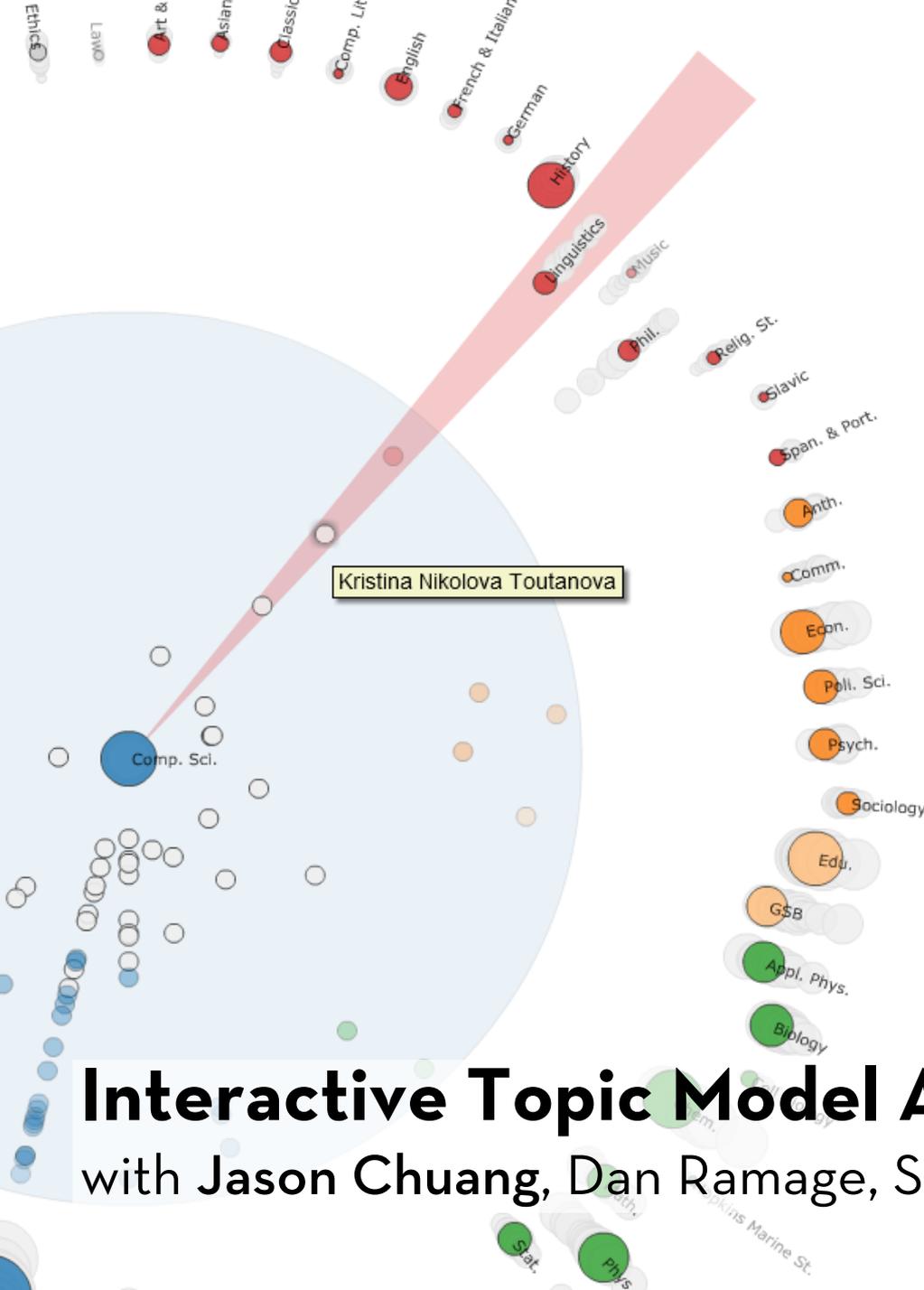


Presentation



Dissemination





## 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.

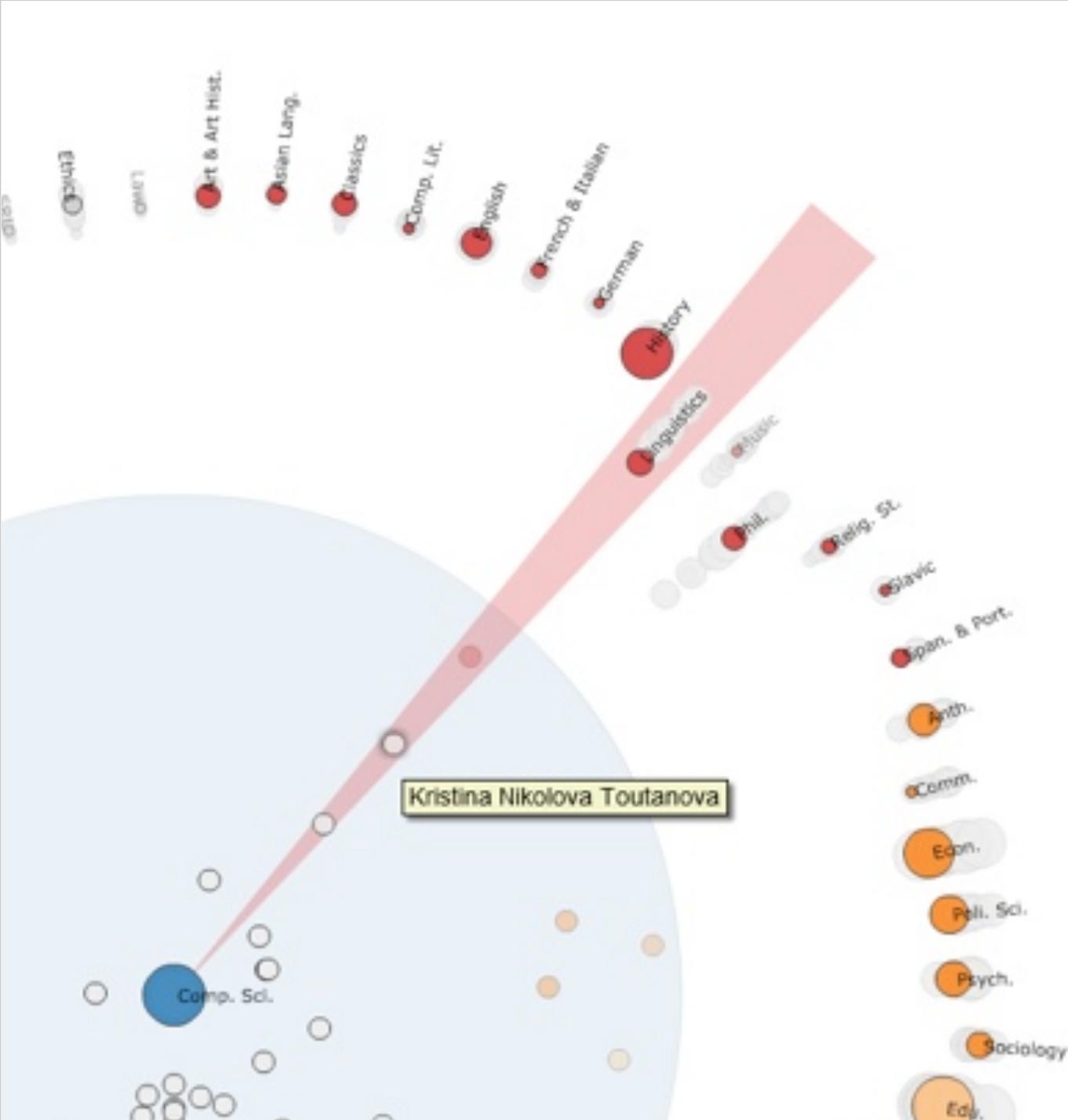
# Interactive Topic Model Assessment

with Jason Chuang, Dan Ramage, Sonal Gupta & Chris Manning





Oh, the humanities!



## Effective statistical models for syntactic disambiguation

Student: Kristina Nikolova Toutanova

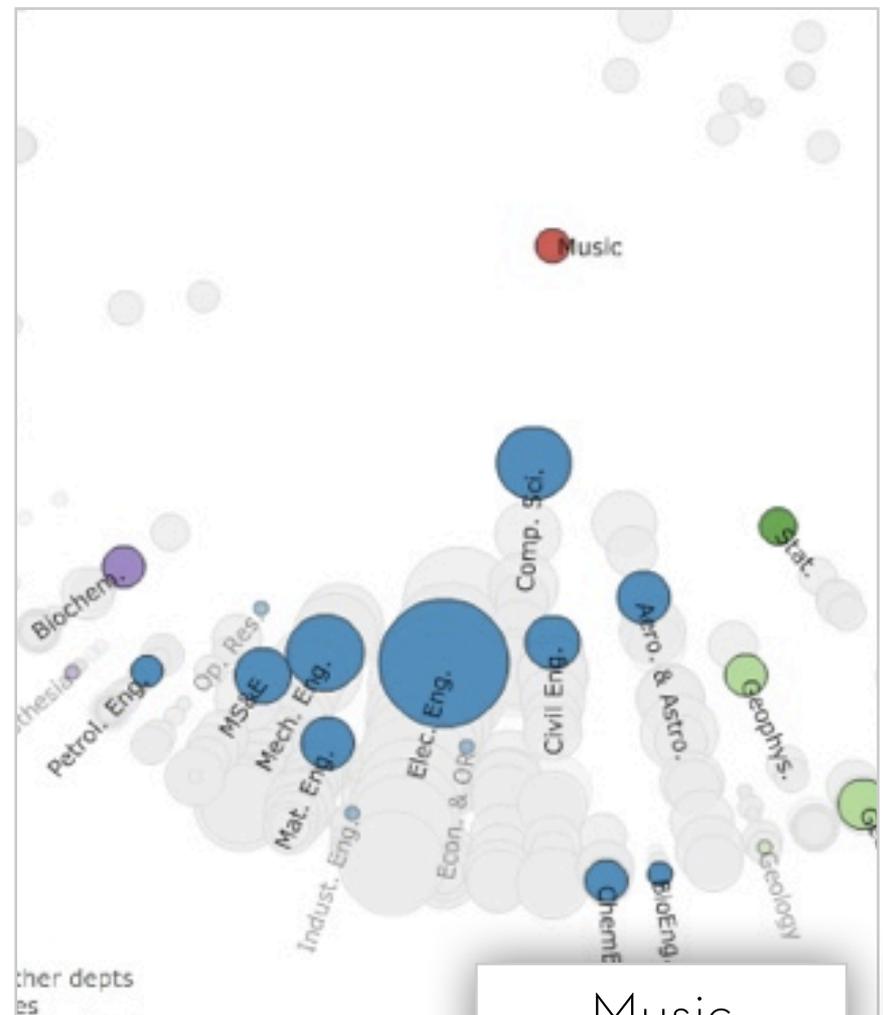
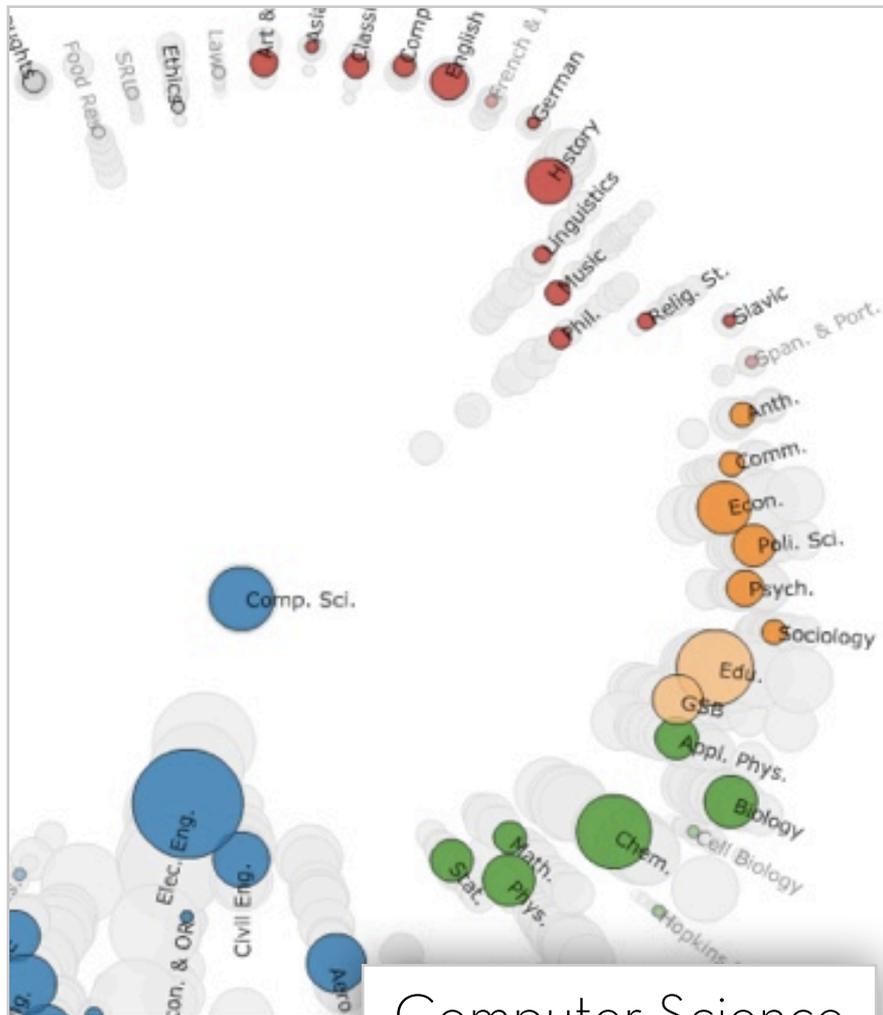
Advisor: Christopher D. Manning

Computer Science (2005)

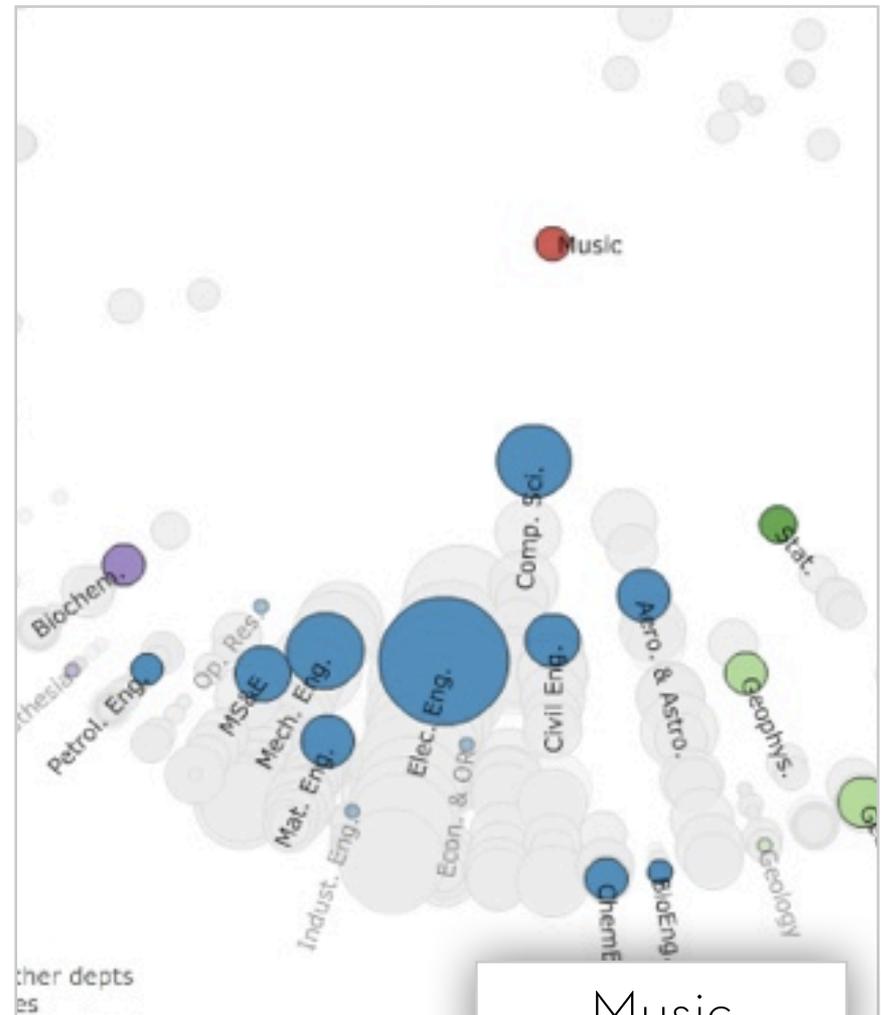
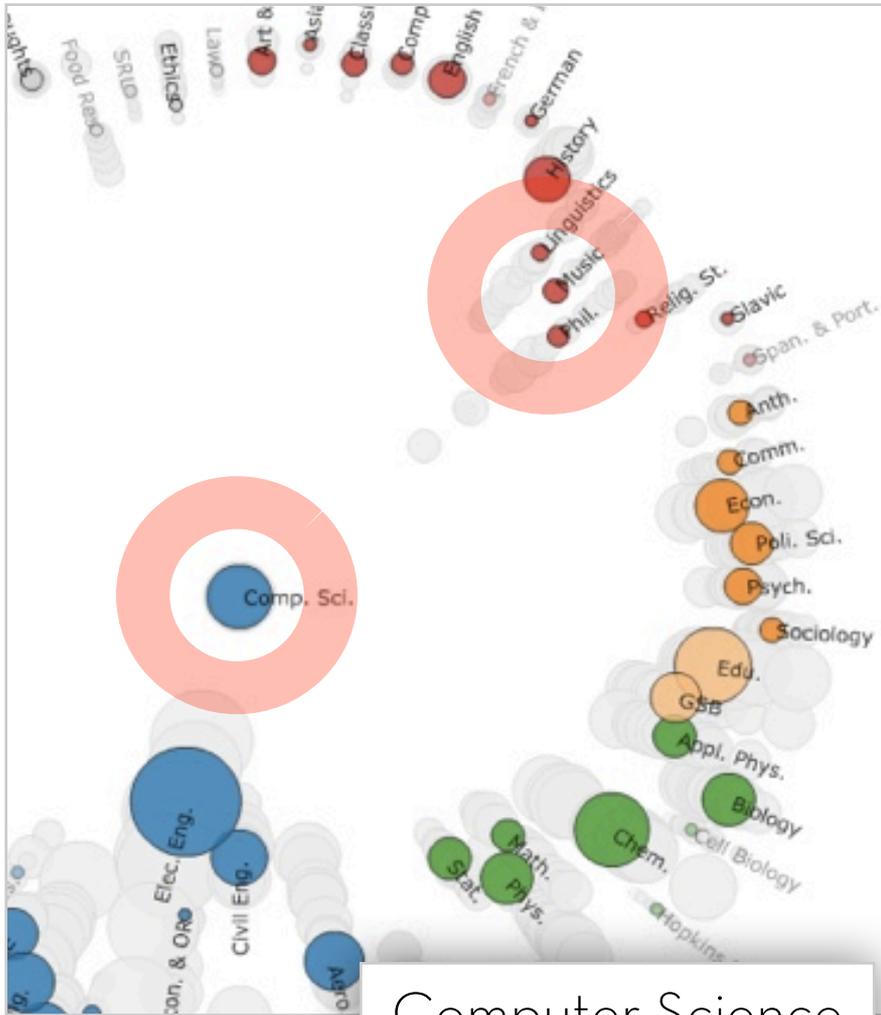
Keywords: Syntactic, Semantic, Tree

Abstract:

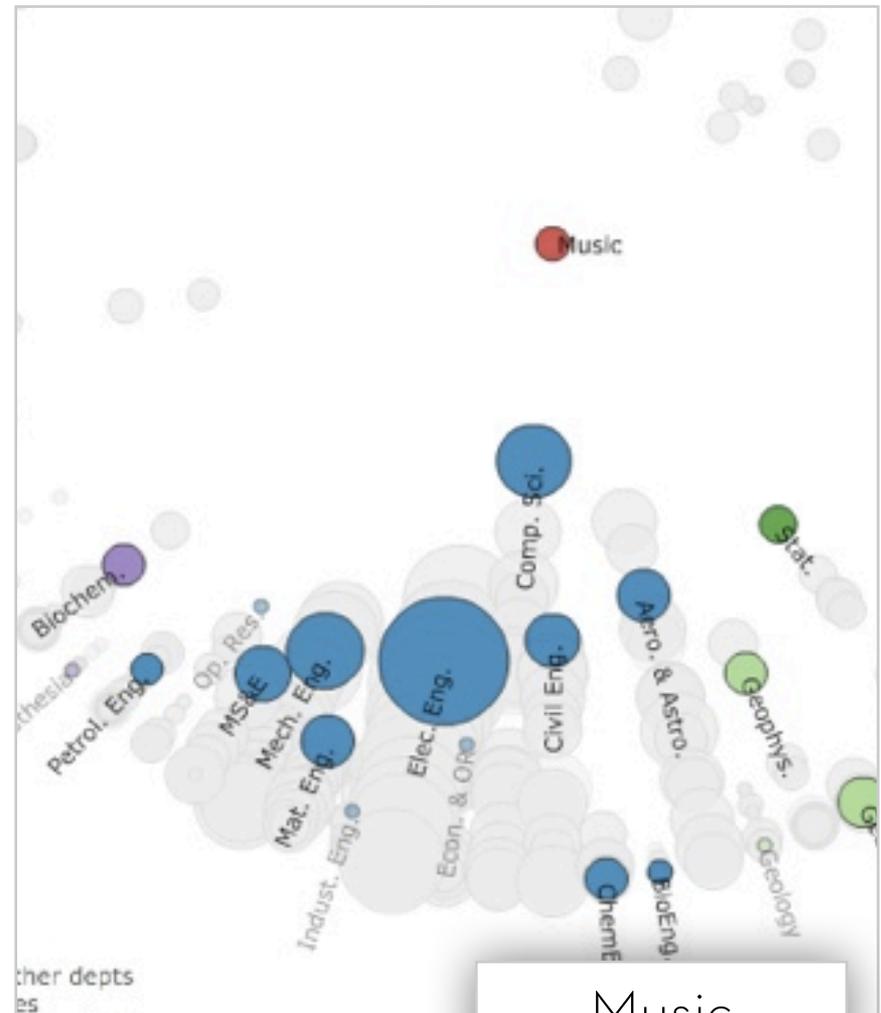
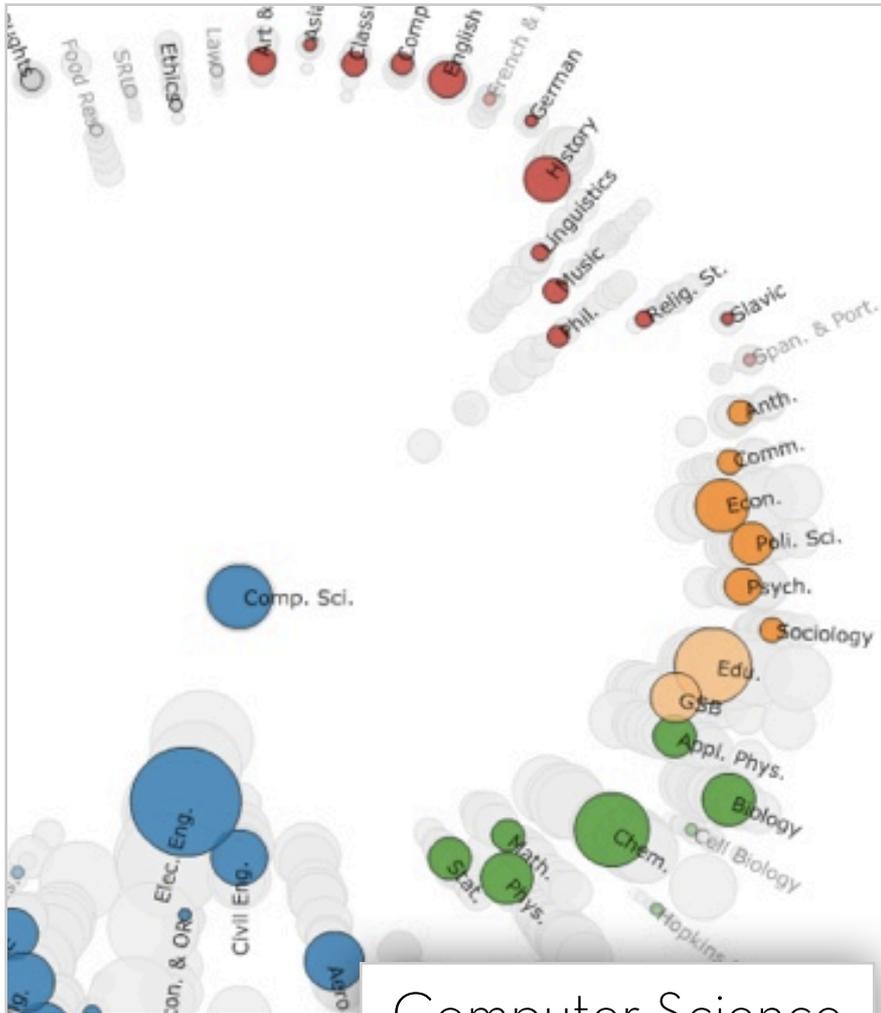
This thesis focuses on building effective models for the syntactic disambiguation of sophisticated syntactic (NL) structures. We advance the state of the art by (i) choosing representations that encode syntactic information effectively and (ii) developing machine learning models that capture the specific properties of NL disambiguation. To address the lack of data and large, structured spaces of syntactic disambiguation, we propose a novel approach to parse trees that connects the words of the structure in a direct way. Experimentally, a Head Driven Phrase Structure Grammar-based model achieves superior performance compared to existing models on the task of disambiguating the semantic structure of a sentence. To build a more accurate model, which captures the semantic structure of a verb frame of a verb is a joint structure with its arguments. We achieve this using a novel approach to Markov independence assumptions and semantic labels. To address the sparsity problem, we develop a method for incorporating semantic information, using Markov chains in a chain framework makes it possible to learn from multiple sources, to learn how much to trust inferences together. It achieves large improvements in disambiguating prepositional phrases.



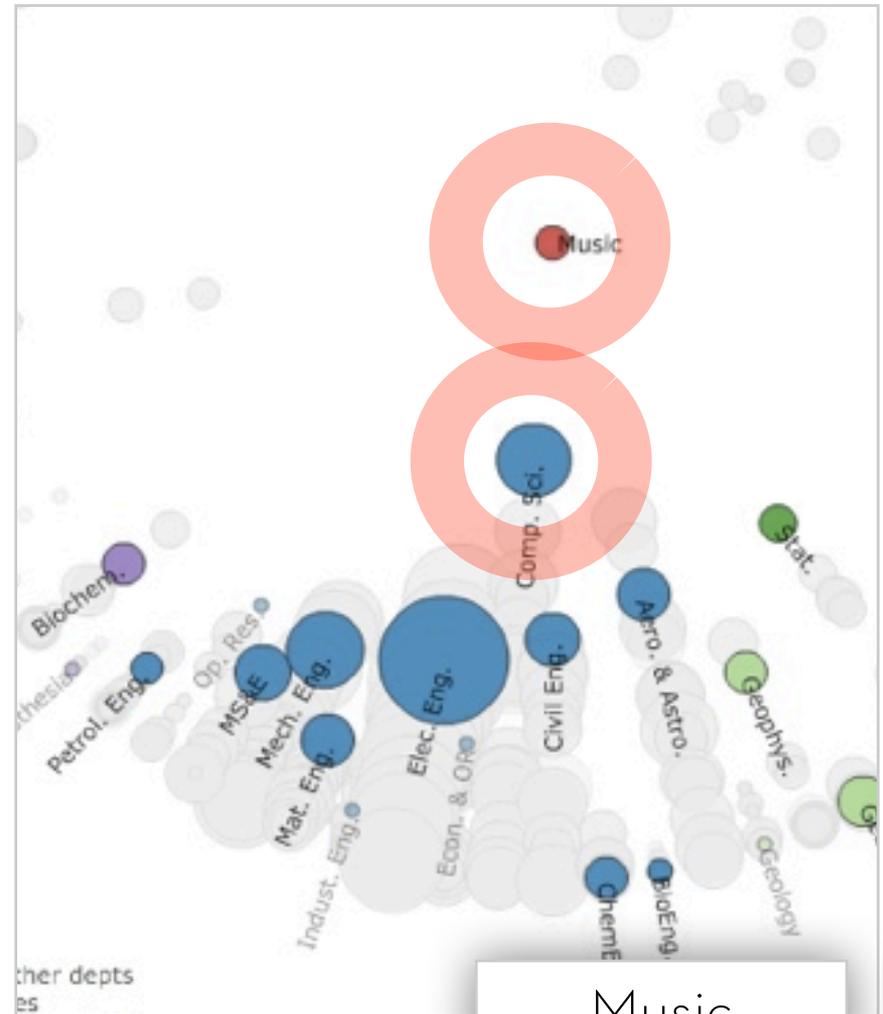
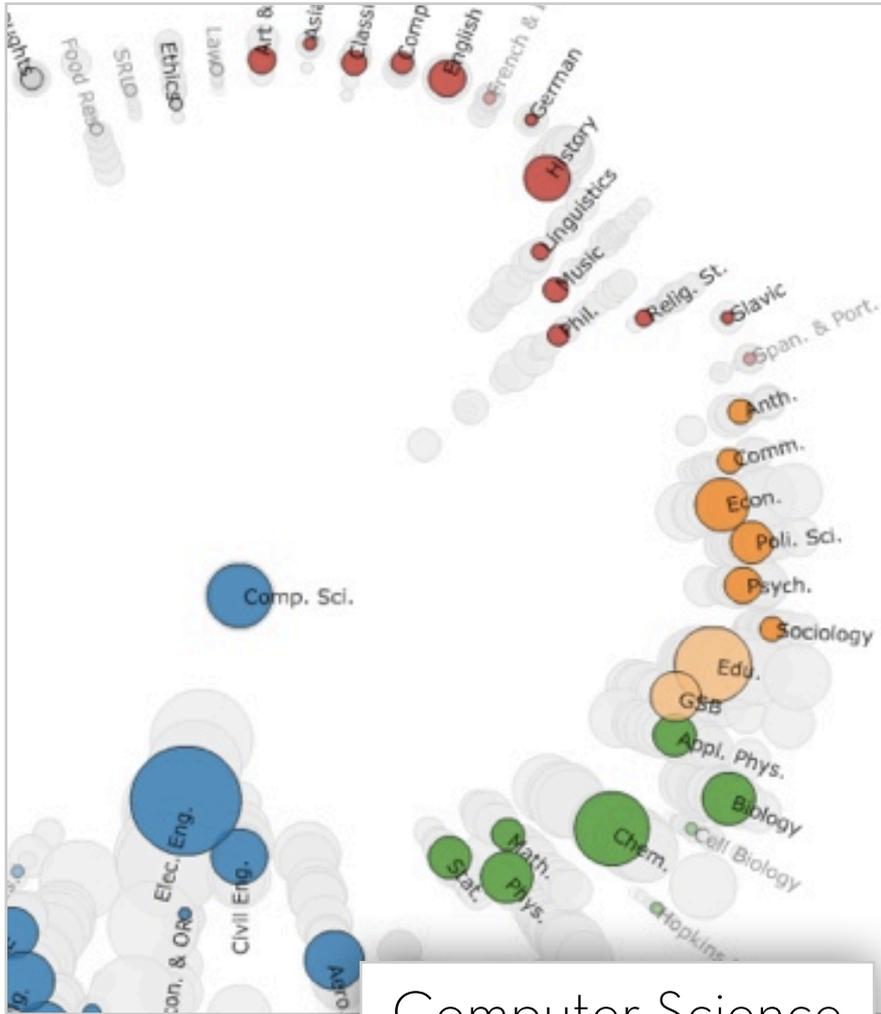
“Word Borrowing” via Labeled LDA



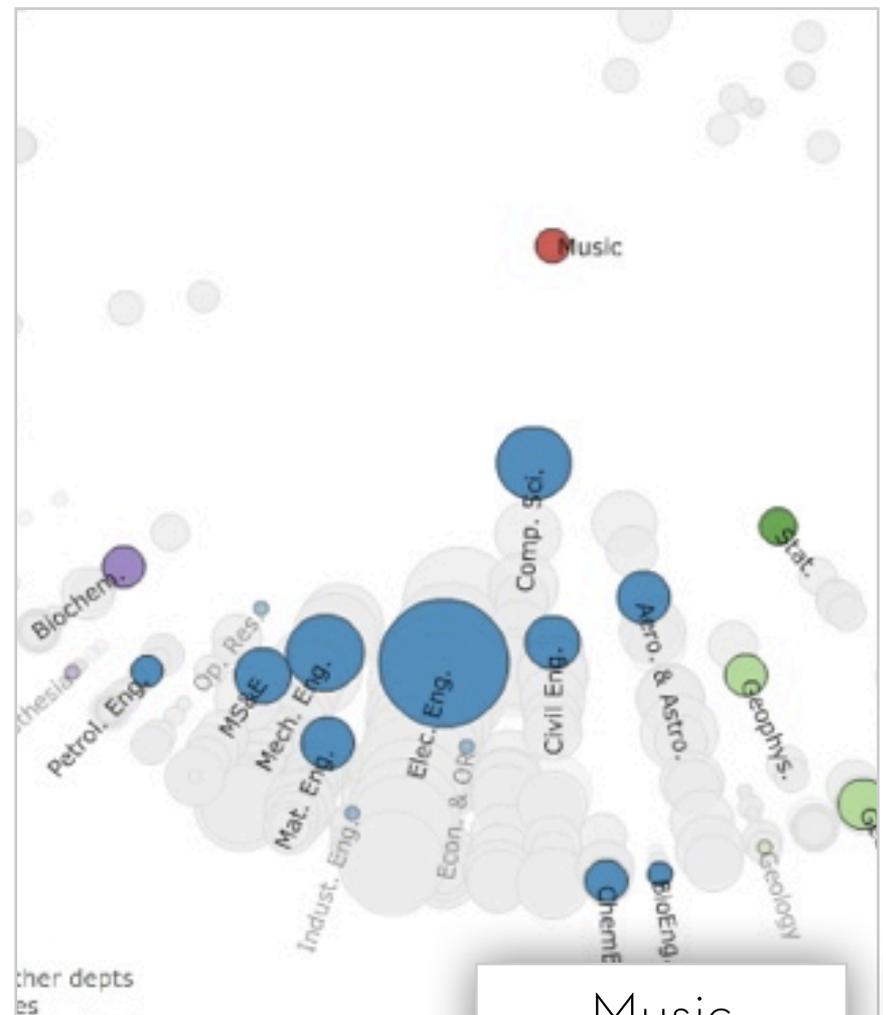
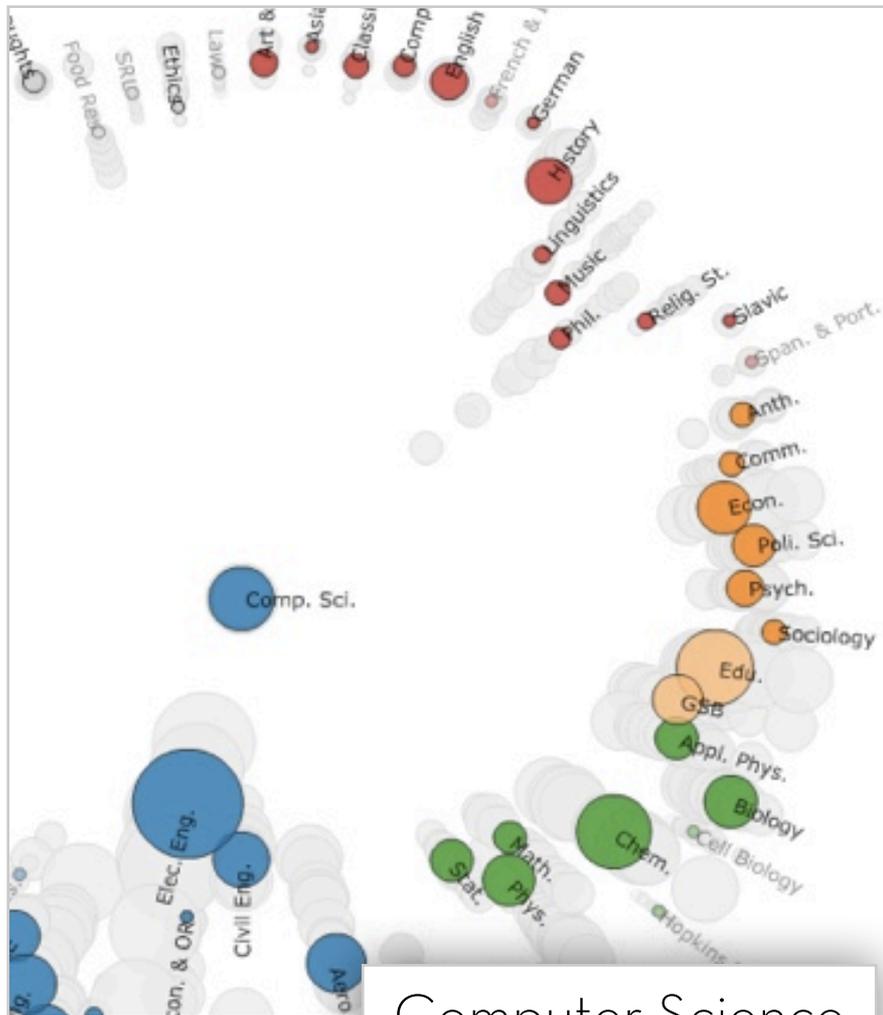
“Word Borrowing” via Labeled LDA



“Word Borrowing” via Labeled LDA



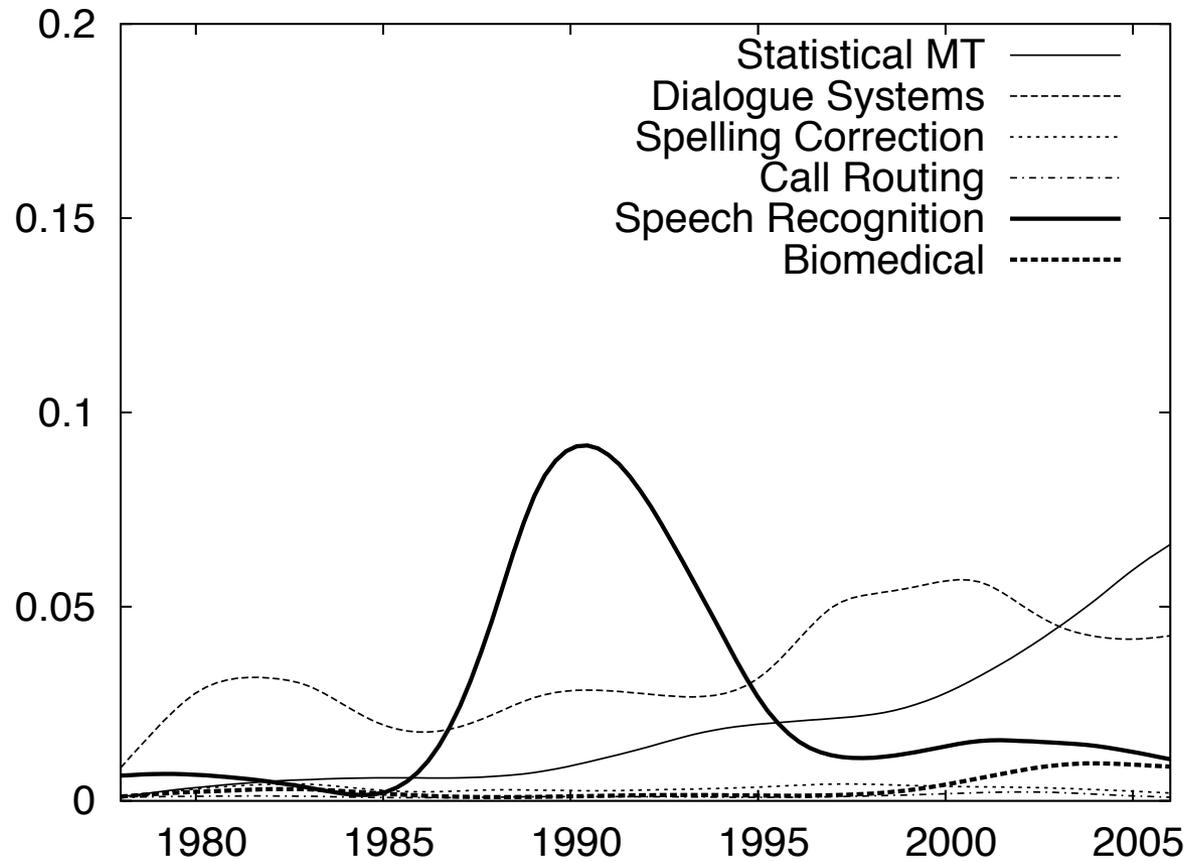
“Word Borrowing” via Labeled LDA



“Word Borrowing” via Labeled LDA

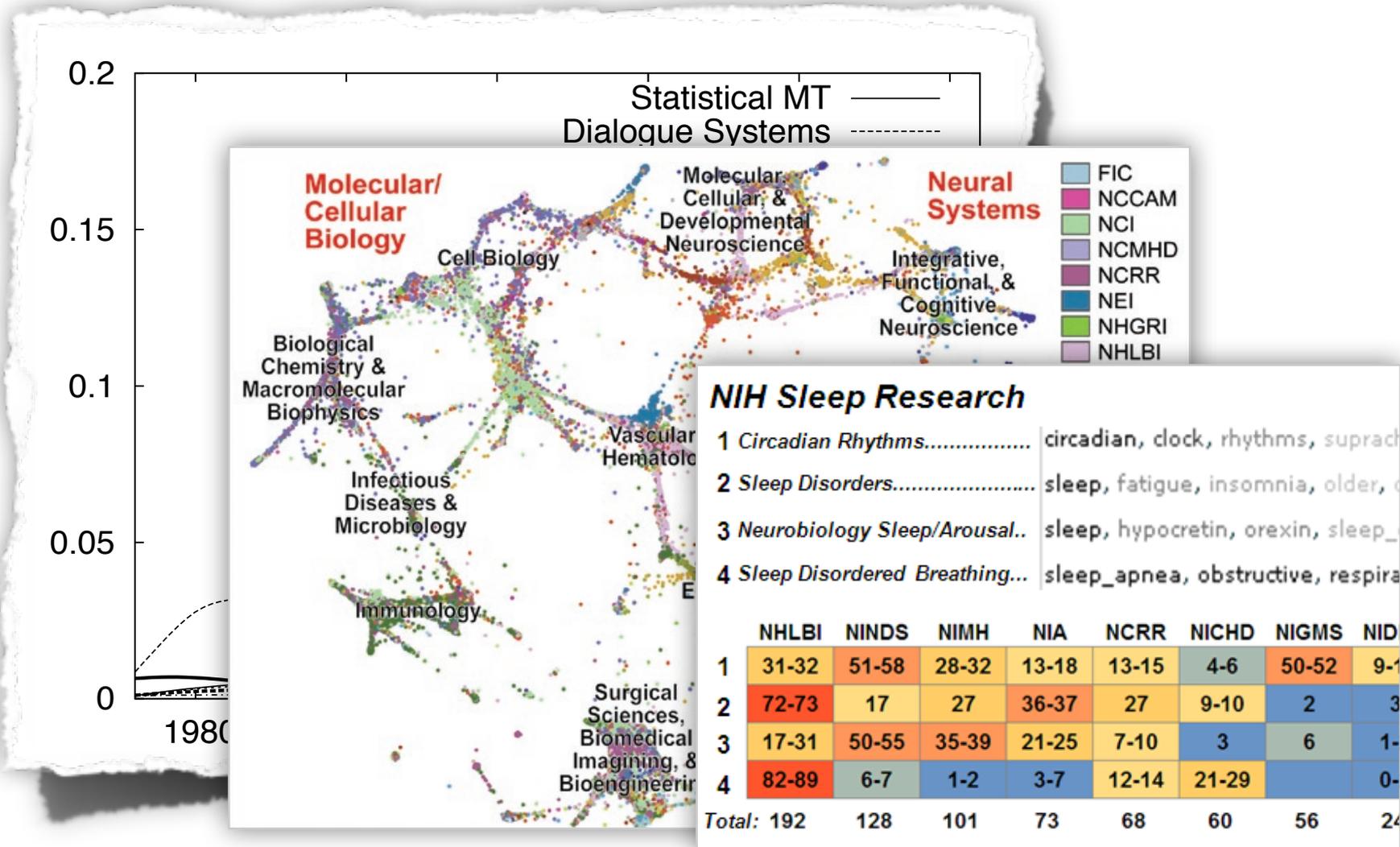
# Real-World Topical Analysis

# Real-World Topical Analysis



[Hall et al. 2008]

# Real-World Topical Analysis



## NIH Sleep Research

- Circadian Rhythms*..... circadian, clock, rhythms, suprach
- Sleep Disorders*..... sleep, fatigue, insomnia, older, c
- Neurobiology Sleep/Arousal*.. sleep, hypocretin, orexin, sleep\_
- Sleep Disordered Breathing*... sleep\_apnea, obstructive, respira

	NHLBI	NINDS	NIMH	NIA	NCRR	NICHD	NIGMS	NID
1	31-32	51-58	28-32	13-18	13-15	4-6	50-52	9-1
2	72-73	17	27	36-37	27	9-10	2	3
3	17-31	50-55	35-39	21-25	7-10	3	6	1-
4	82-89	6-7	1-2	3-7	12-14	21-29		0-
Total:	192	128	101	73	68	60	56	24

[Talley et al. 2011]

# Current Practices

# Current Practices

**Topic Words:** vegf angiogenesis vascular\_endothelial\_growth\_factor angiogenic end  
antiangiogenic anti\_angiogenic vegf\_a tumor\_angiogenesis vegfr2 growth signaling t  
**Title Words:** angiogenesis, vegf, vascular\_endothelial\_growth\_factor, angiogenic, tu  
neovascularization, angiopoietin, signaling, vegfr, vascular, human  
**Phrases:** vascular\_endothelial\_growth\_factor vegf, vegf angiogenesis, vegf receptor,

[Talley et al. 2011]

# Current Practices

**Topic Words:** vegf angiogenesis vascular\_endothelial\_growth\_factor angiogenic end  
antiangiogenic anti\_angiogenic vegf\_a tumor\_angiogenesis vegfr2 growth signaling t

**Title**  
neova  
**Phras**

**Anaphora Resolution**

resolution anaphora pronoun discourse antecedent pronouns core

**Automata**

string state set finite context rule algorithm strings language sym

**Biomedical**

medical protein gene biomedical wkh abstracts medline patient c

**Call Routing**

call caller routing calls destination vietnamese routed router dest

**Categorial Grammar**

proof formula graph logic calculus axioms axiom theorem proofs

**Centering\***

centering cb discourse cf utterance center utterances theory coher

**Classical MT**

japanese method case sentence analysis english dictionary figure

**Classification/Tagging**

features data corpus set feature table word tag al test

**Comp. Phonology**

vowel phonological syllable phoneme stress phonetic phonology

**Comp. Semantics\***

semantic logical semantics john sentence interpretation scope log

# Current Practices

**Topic Words:** vegf angiogenesis vascular\_endothelial\_growth\_factor angiogenic end  
antiangiogenic anti\_angiogenic vegf\_a tumor\_angiogenesis vegfr2 growth signaling t

**Title**

neova

**Phras**

**Anaphora Resolution**

resolution anaphora pronoun discourse antecedent pronouns core

**Automata**

string state set finite context rule algorithm strings language sym

**Biomed**

**Call Ro**

Topic 0: [ireland featureset dialect objlist distance connacht dialects soa  
ulster rss munster murasaki loebner galway ctl prize icons distances]

**Categor**

**Centeri**

MT/ALIGNMENT Topic 1: [english word alignment language source sentence target

**Classica**

bilingual phrase pairs model parallel mt french al system statistical corpus

**Classific**

**Comp. I**

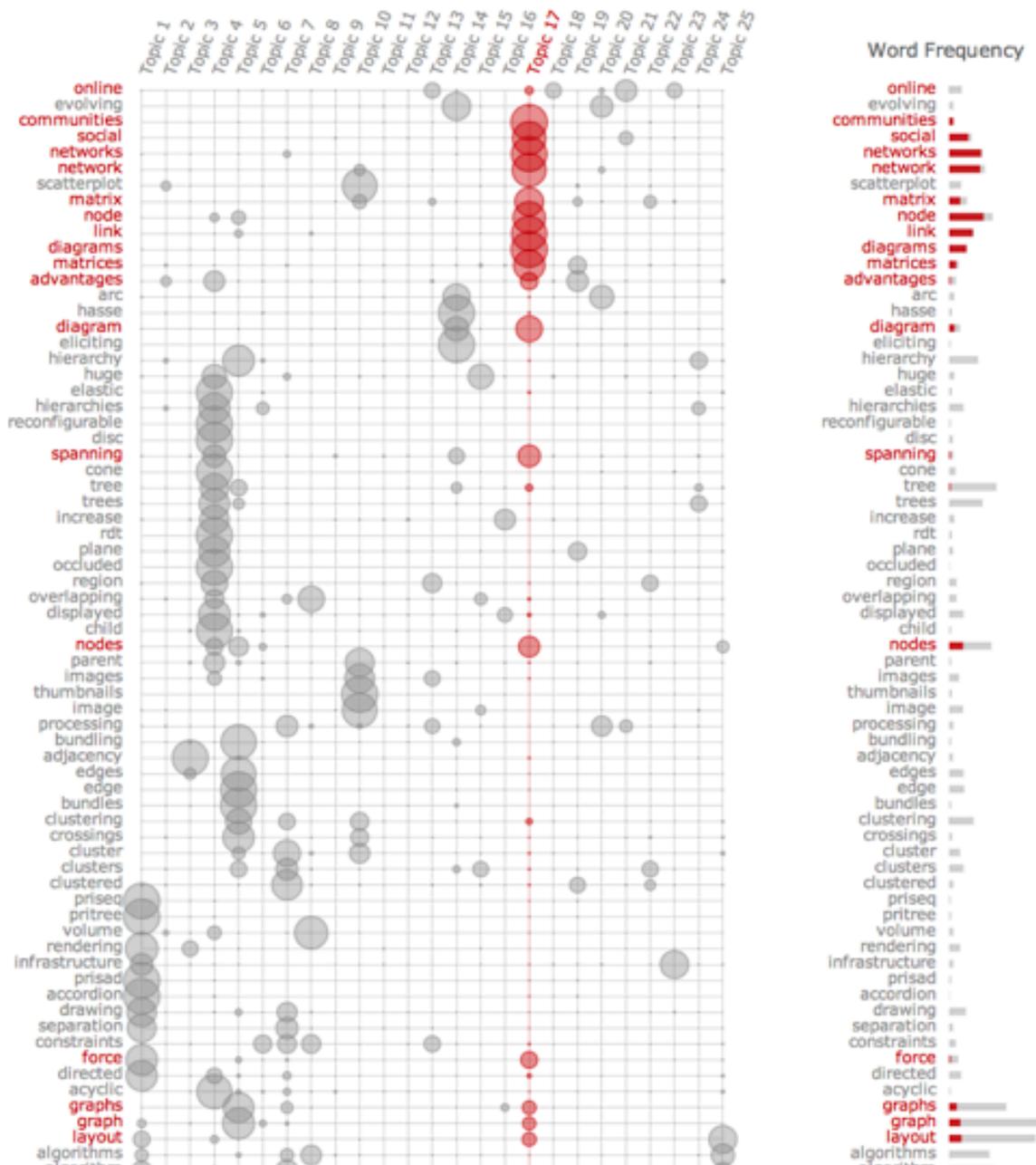
Topic 2: [data features corpus set feature table tag pos al tagging tags acc  
classification test performance part information classifier ing]

**Comp. S**

LEXICAL SEMANTICS Topic 3: [word semantic lexical information senses noun wa  
dictionary relations knowledge nouns text domain disambiguation context simi  
relation number concept]

How might we better support  
**human-in-the-loop verification**  
of topic models?

# Termite | Topic Model Visualization



Representative Document
A Comparison of the Readability of Graphs Using Node-Link and Matrix Representations Mohammad Ghoniem Jean-Daniel Fekete Philippe Castagliola
Using Multilevel Call Matrices in Large Software Projects Frank van Ham
Improving the Readability of Clustered Social Networks using Node-Link Representations Nathalie Henry Anastasia Bezerianos Jean-Daniel Fekete
MatrixExplorer: a Dual-Representation System to Explore Social Networks Nathalie Henry Jean-Daniel Fekete
<b>NodeTriX: a Hybrid Visualization of Social Networks</b> Nathalie Henry Jean-Daniel Fekete Michael J. McGuffin
The need to visualize large social networks is growing as hardware and many new data sets become available. Unfortunately, the visual representation of a network, while arbitrary portions of the network can be shown, does not resolve the basic dilemma of being readable both for the global structure and for the analysis of communities. To address this problem, we present NodeTriX, which combines the advantages of two traditional representations: node-link and matrix. NodeTriX visualization by dragging selections to and from node-link and matrix representations to explore the dataset and create meaningful views. Finally, we present a case study applying NodeTriX to the analysis of a social network to illustrate the capabilities of NodeTriX as both an exploration tool and a visualization.
Visualizing Causal Semantics using Animations Nivedita R. Kadaba Pourang P. Irani Jason Leboe
<b>Balancing Systematic and Flexible Exploration of Social Networks</b> Adam Perer Ben Shneiderman
Social network analysis (SNA) has emerged as a powerful method for analyzing networks. However, interactive exploration of networks is currently limited by the complexity of patterns and comprehend the structure of networks with many nodes and edges. A medley of statistical methods and overwhelming visual output which does not help explore in an orderly manner. This results in exploration that is largely ineffective. SocialAction helps structural analysts understand social networks more effectively by providing attribute ranking and coordinated views to help users systematically explore networks. SocialAction offers analysts a strategy beyond opportunistic exploration. SocialAction offers analysts a strategy beyond opportunistic exploration. SocialAction offers analysts a strategy beyond opportunistic exploration.
Causality Visualization Using Animated Growing Polygons Niklas Elmqvist Philippos Tsigas
SpicyNodes: Radial Layout Authoring for the General Public

## 1 Graph Visualization

Graph, network, node-link diagram, layout, adjacency matrix, reordering

Asymmetric Relations in Longitudinal Social Networks

Multi-Level Graph Layout on the GPU

Balancing Systematic and Flexible Exploration of Social Networks

Parallel Edge Splatting for Scalable Dynamic Graph Visualization

## 3 Multidimensional visualization

Parallel coordinates, small multiples, splom, scatterplot matrix, multidimensional projections, embeddings, MDS, PCA

Rolling the Dice: Multidimensional Visual Exploration using Scatterplots

Scattering Points in Parallel Coordinates

Multidimensional Detective

Improved Similarity Trees and their Application to Visual Data

Steerable, Progressive Multidimensional Scaling

## 5 Software Visualization

algorithm animation, traces, logs

code\_swarm: A Design Study in Organic Software Visualization

The Visual Code Navigator: An Interactive Toolset for Source Code

Using Multilevel Call Matrices in Large Software Projects

## 2 Text Visualization

Text, topics, sentiment analysis

The Shape of Shakespeare: Visualizing Text Using Implicit Surfaces

ThemeRiver: visualizing theme changes over time

From Metaphor to Method: Cartographic Perspectives on Informal

Participatory Visualization with Wordle

FacetAtlas: Multifaceted Visualization for Rich Text Corpora

Mapping Text with Phrase Nets

## 4 Tree Visualization

Treemap, node-link diagram, hierarchies

SpaceTree: supporting exploration in large node link tree, design

Browsing Zoomable Treemaps: Structure-Aware Multi-Scale Navigation

Interactive Visualization of Genealogical Graphs

## 6

**Topic:** Significant and coherent area of research

**Exemplary terms:** techniques, methods, systems, people...  
Separate the terms by commas or semicolons

**Exemplary documents:** 3 or more papers  
Drag and drop from InfoVis proceedings

## 2011 InfoVis Conference

Providence, Rhode Island

### Theory and Foundations

#### Quality Metrics in High-Dimensional Data Visualization: An Overview and Systematization

Enrico BERTINI Andrada TATU Daniel KEIM

#### Benefitting InfoVis with Visual Difficulties

Jessica HULLMAN Eytan ADAR Priti SHAH

#### Product Plots

Hadley WICKHAM Heike HOFMANN

#### Visualization Rhetoric: Framing Effects in Narrative Visualization

Jessica HULLMAN Nick DIAKOPOULOS

#### Adaptive Privacy-Preserving Visualization Using Parallel Coordinates

Aritra DASGUPTA Robert KOSARA

### Techniques

#### Context-Preserving Visual Links

Markus STEINBERGER Manuela WALDNER  
Marc STREIT Alexander LEX  
Dieter SCHMALSTIEG

#### Design Study of LineSets, a Novel Set Visualization Technique

Basak ALPER Nathalie RICHE Gonzalo RAMOS  
Mary CZERWINSKI

#### Developing and Evaluating Quilts for the Depiction of Large Layered Graphs

Juhee BAE Benjamin WATSON

#### Arc Length-based Aspect Ratio Selection

Justin TALBOT John GERTH Pat HANRAHAN

#### Asymmetric Relations in Longitudinal Social Networks

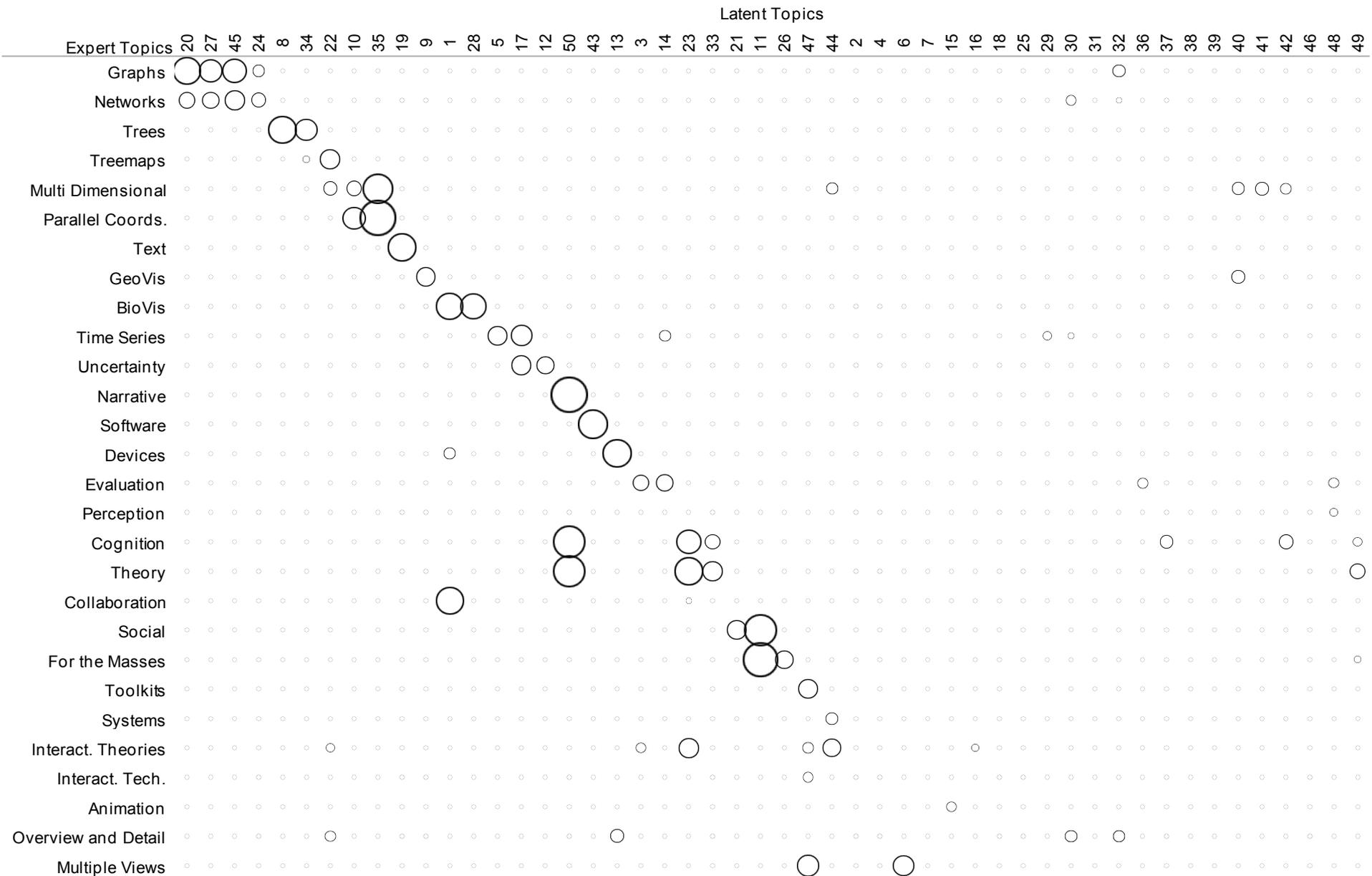
Ulrik BRANDES Bobo NICK

### Systems and Frameworks

#### VisBricks: Multiform Visualization of Large, Inhomogeneous Data

Alexander LEX Hans-Joerg SCHULZ  
Marc STREIT Christian PARTL  
Dieter SCHMALSTIEG

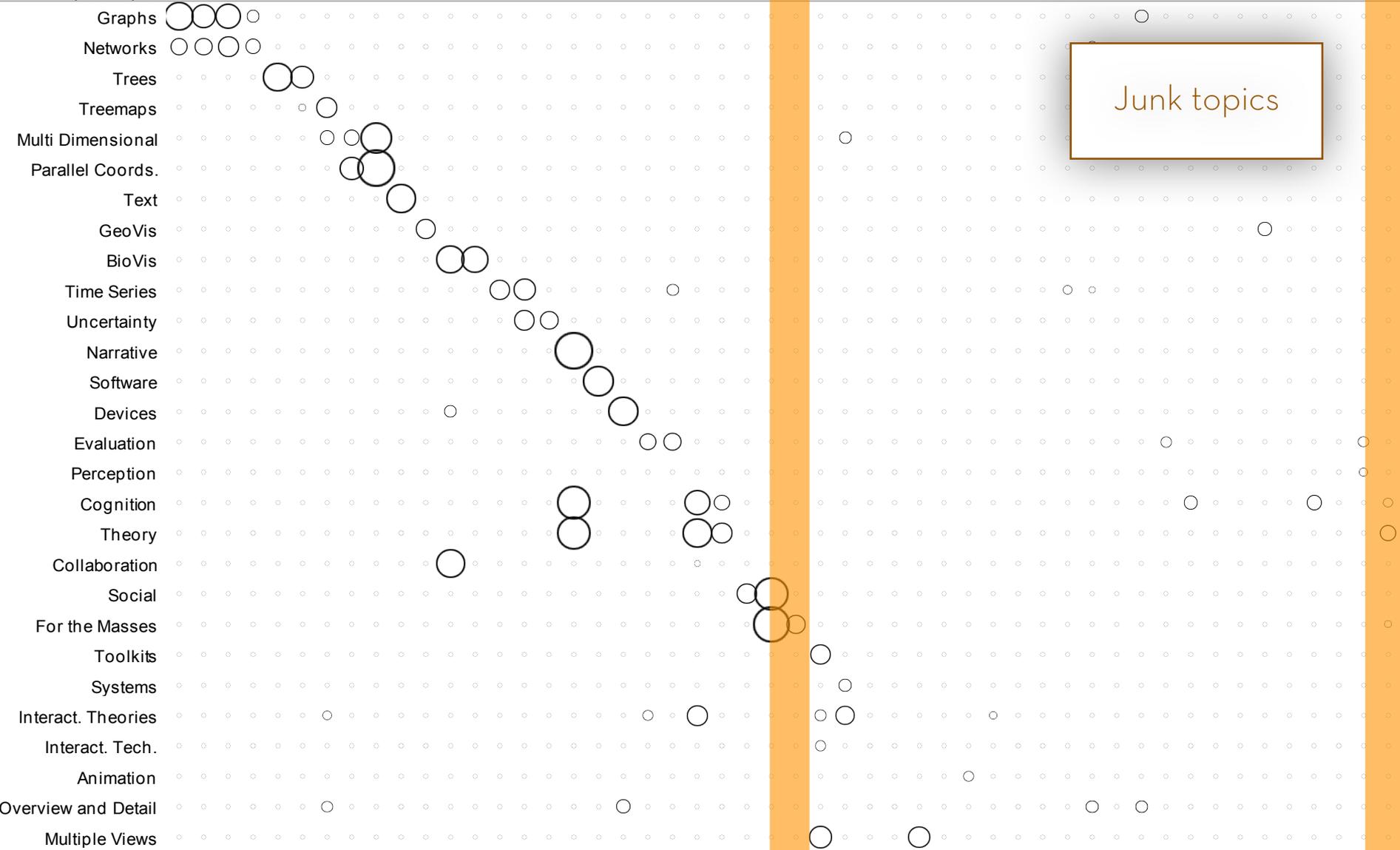
#### Do-Data-Driven Documents



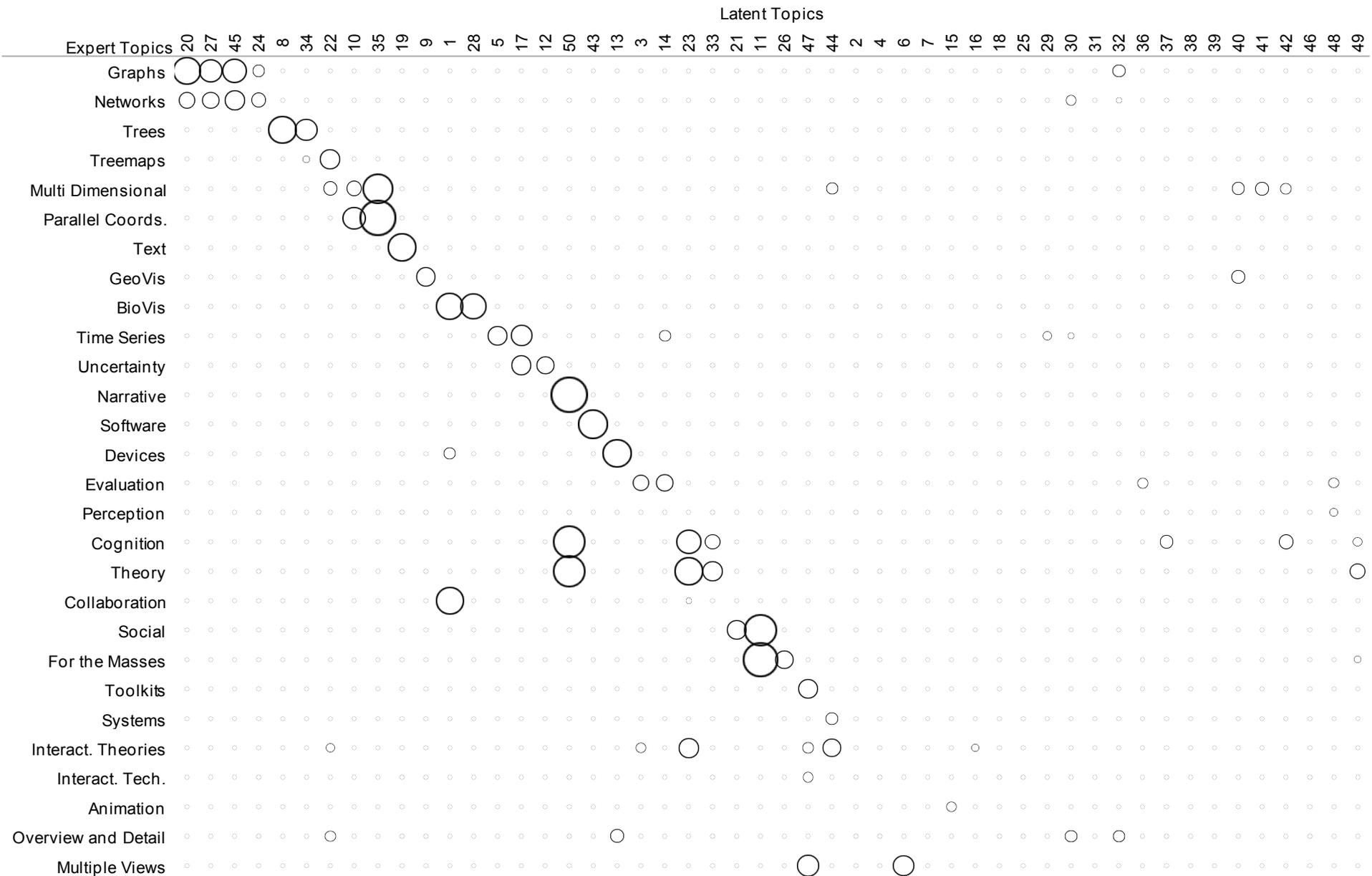
Latent Topics

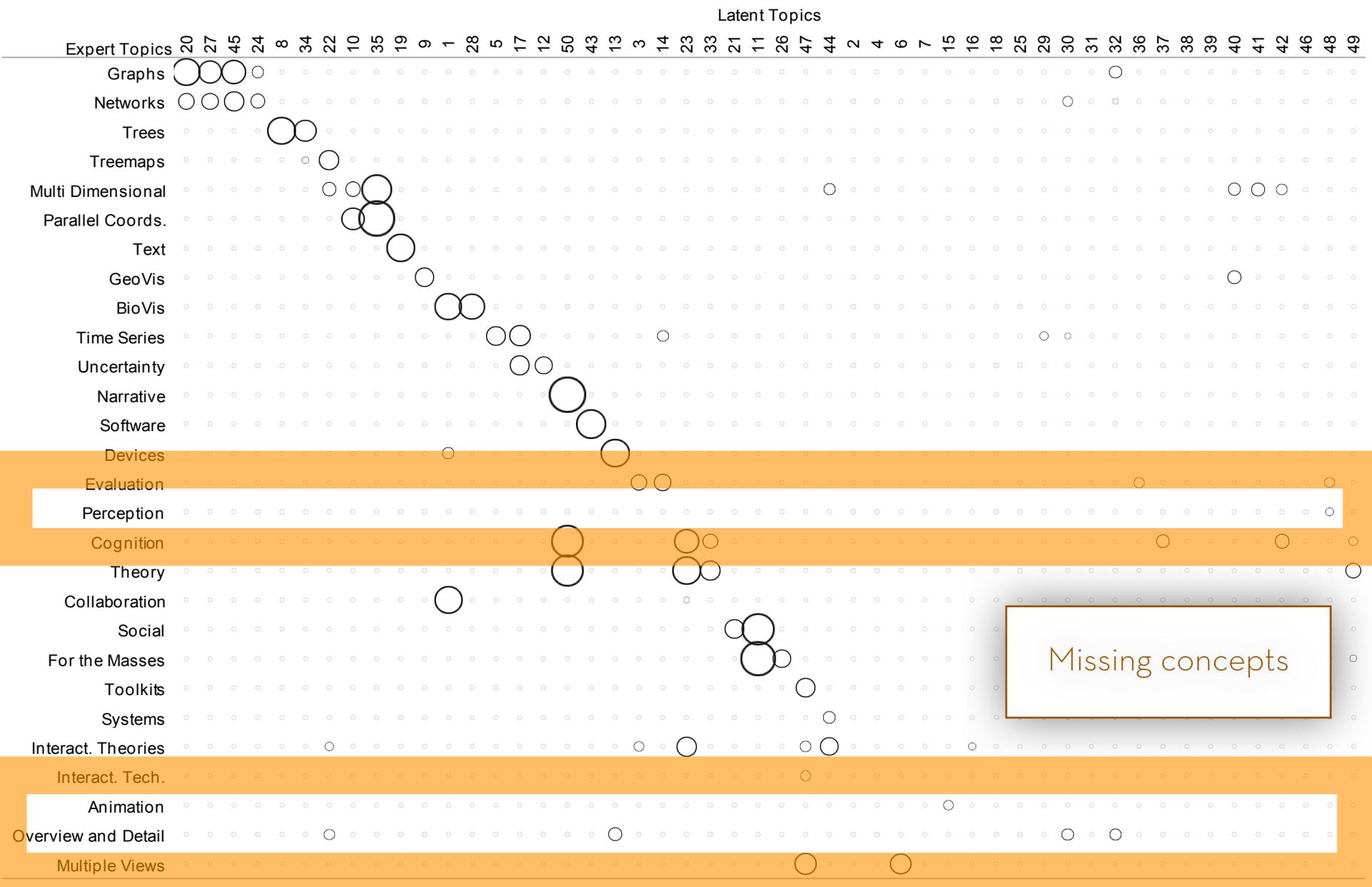
Expert Topics

20 27 45 24 8 34 22 10 35 19 9 1 28 5 17 12 50 43 13 3 14 23 33 21 11 26 47 44 2 4 6 7 15 16 18 25 29 30 31 32 36 37 38 39 40 41 42 46 48 49

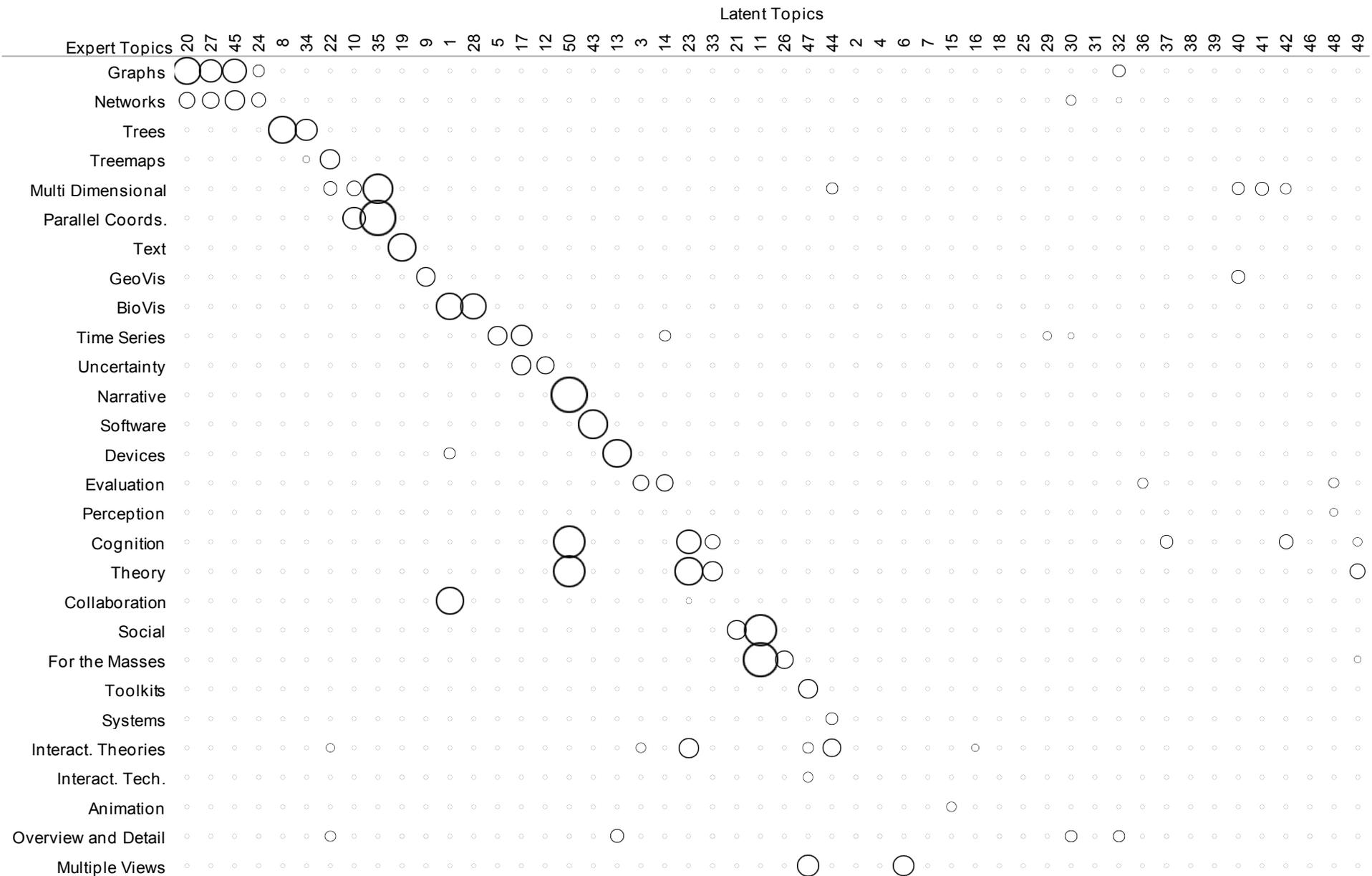


Junk topics





Missing concepts



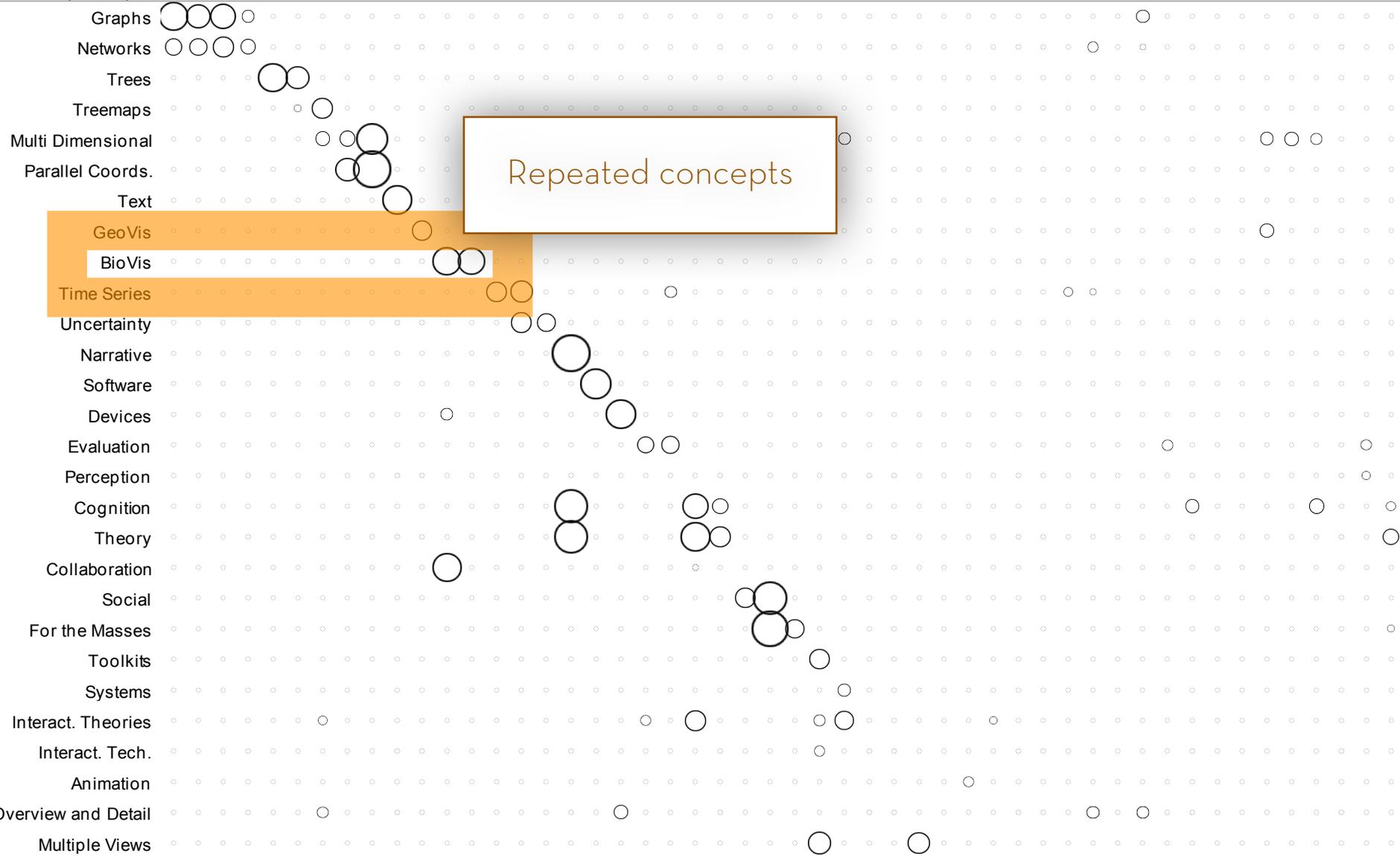




Latent Topics

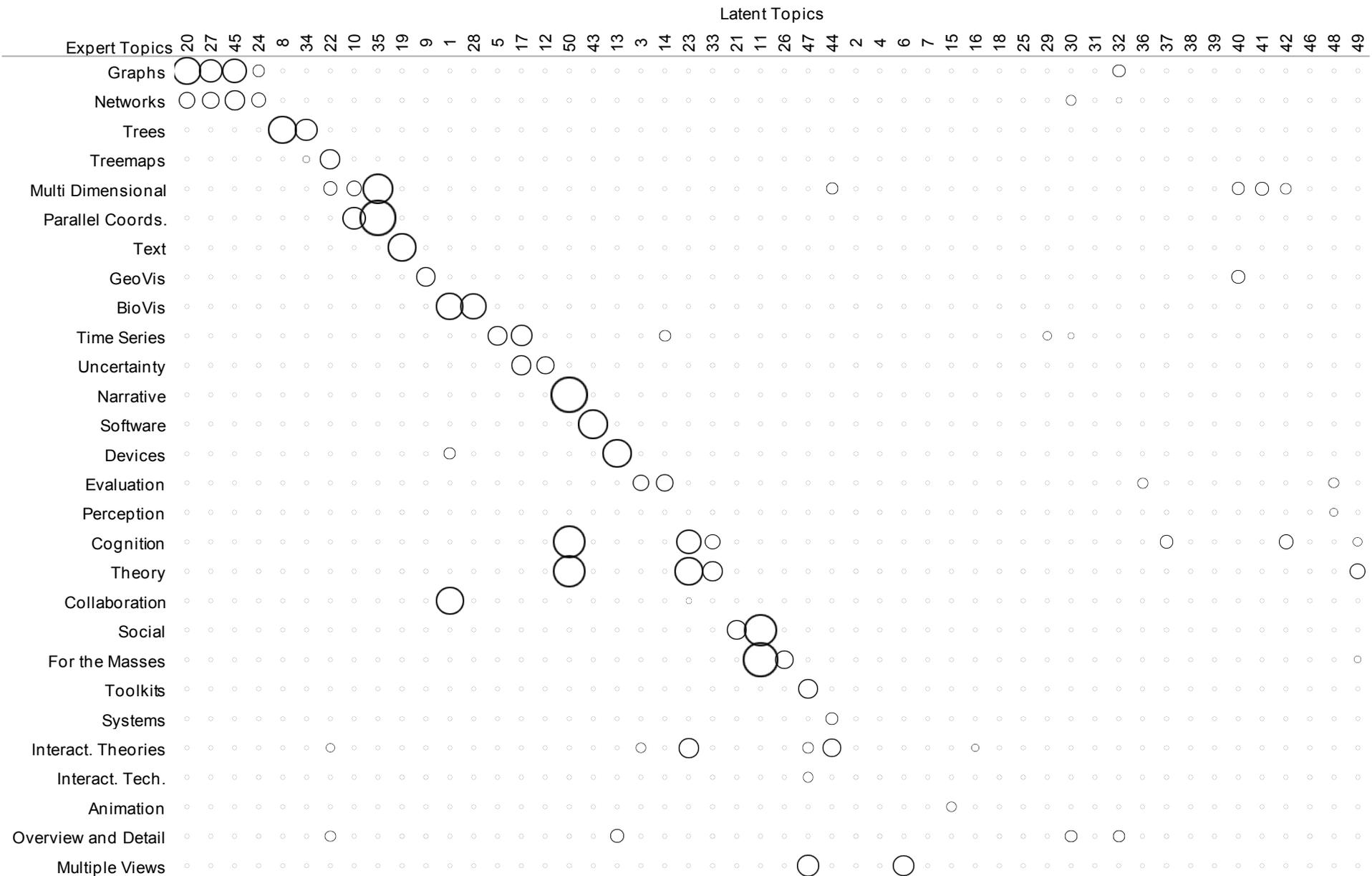
Expert Topics

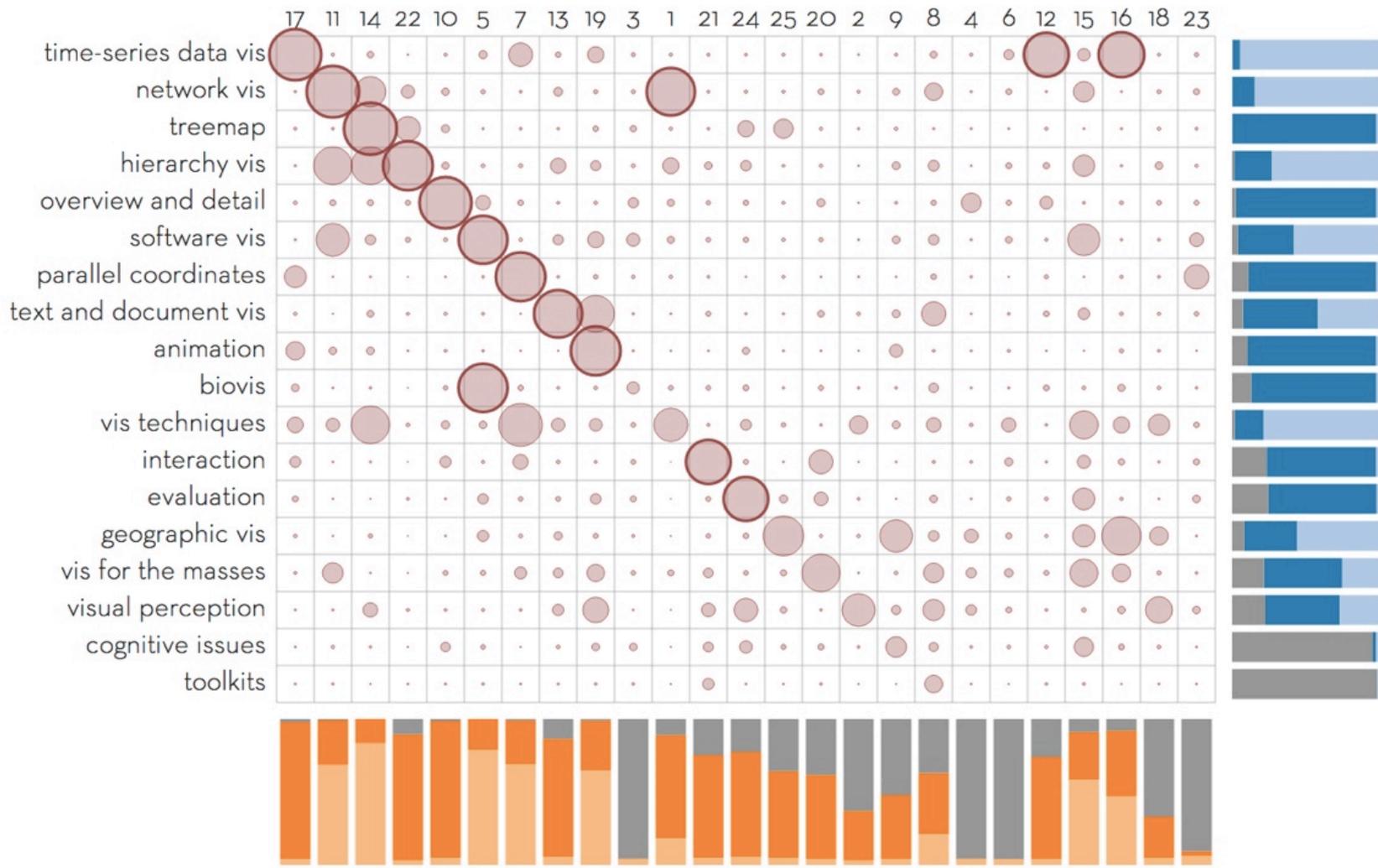
20 27 45 24 8 34 22 10 35 19 9 1 28 5 17 12 50 43 13 3 14 23 33 21 11 26 47 44 2 4 6 7 15 16 18 25 29 30 31 32 36 37 38 39 40 41 42 46 48 49



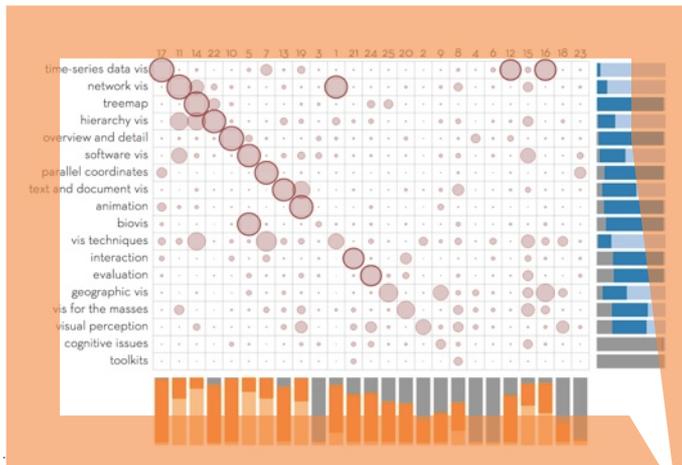
Repeated concepts

GeoVis  
BioVis  
Time Series



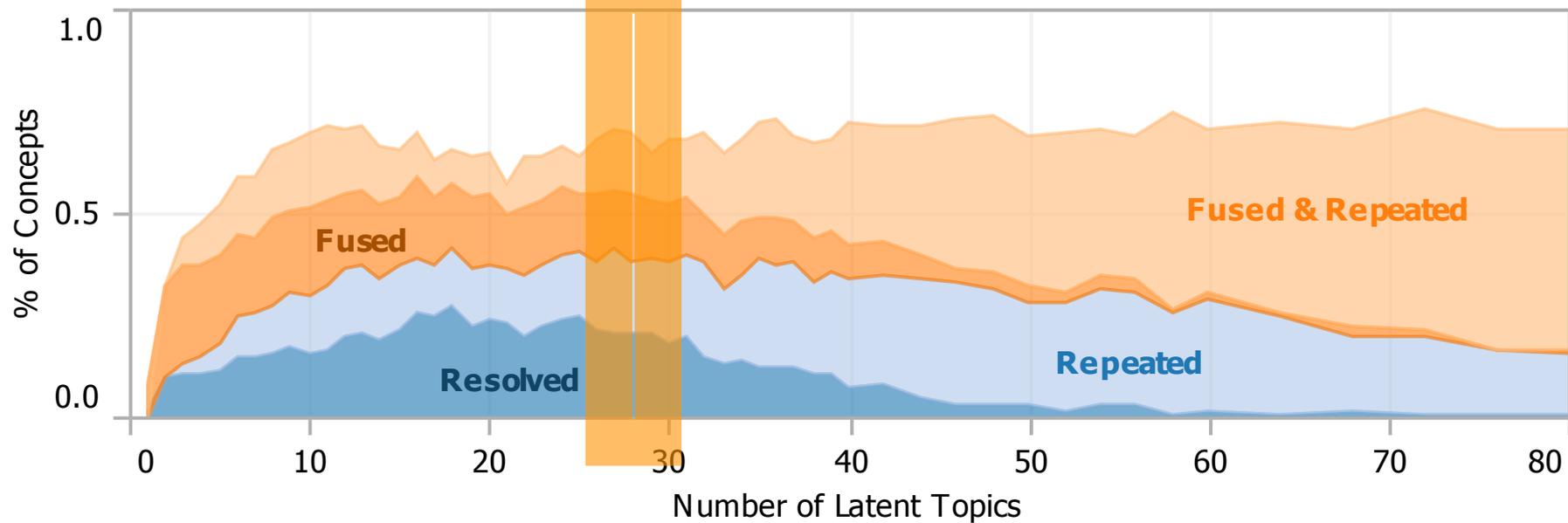






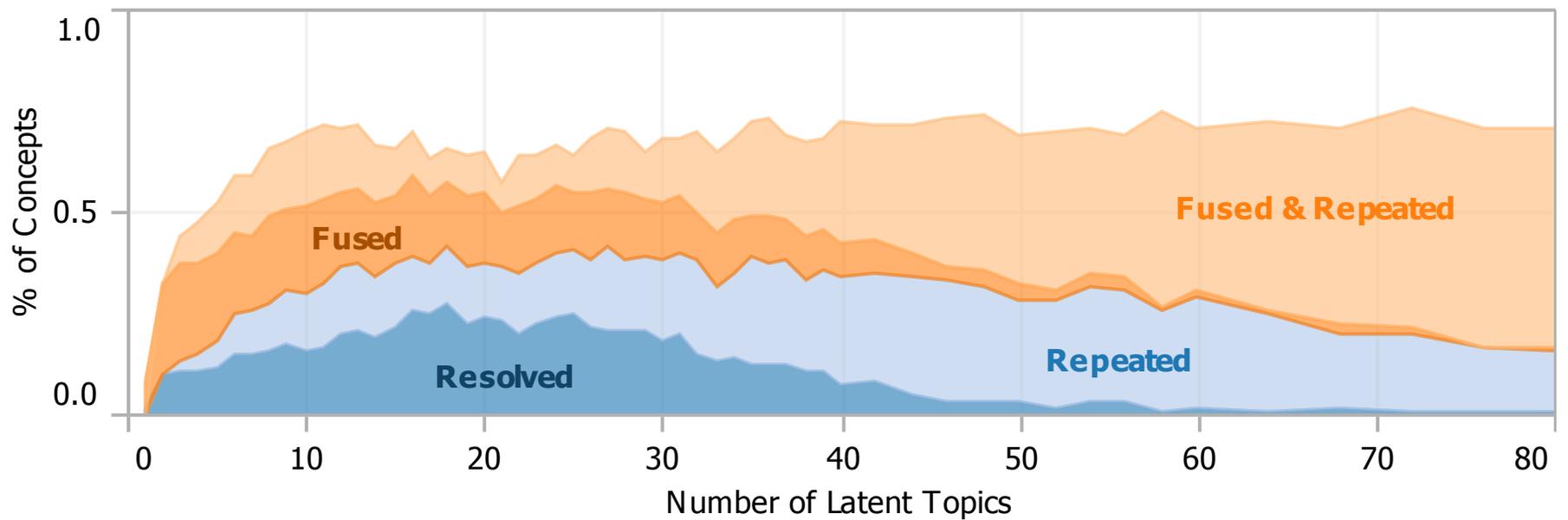
## Resolved/Fused Concepts vs. Number of Latent Topics

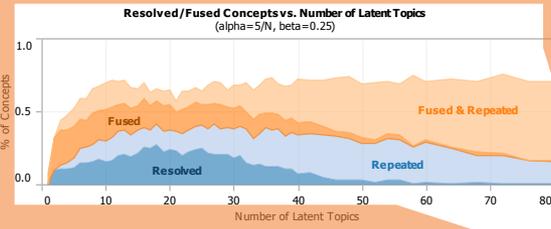
( $\alpha=5/N$ ,  $\beta=0.25$ )



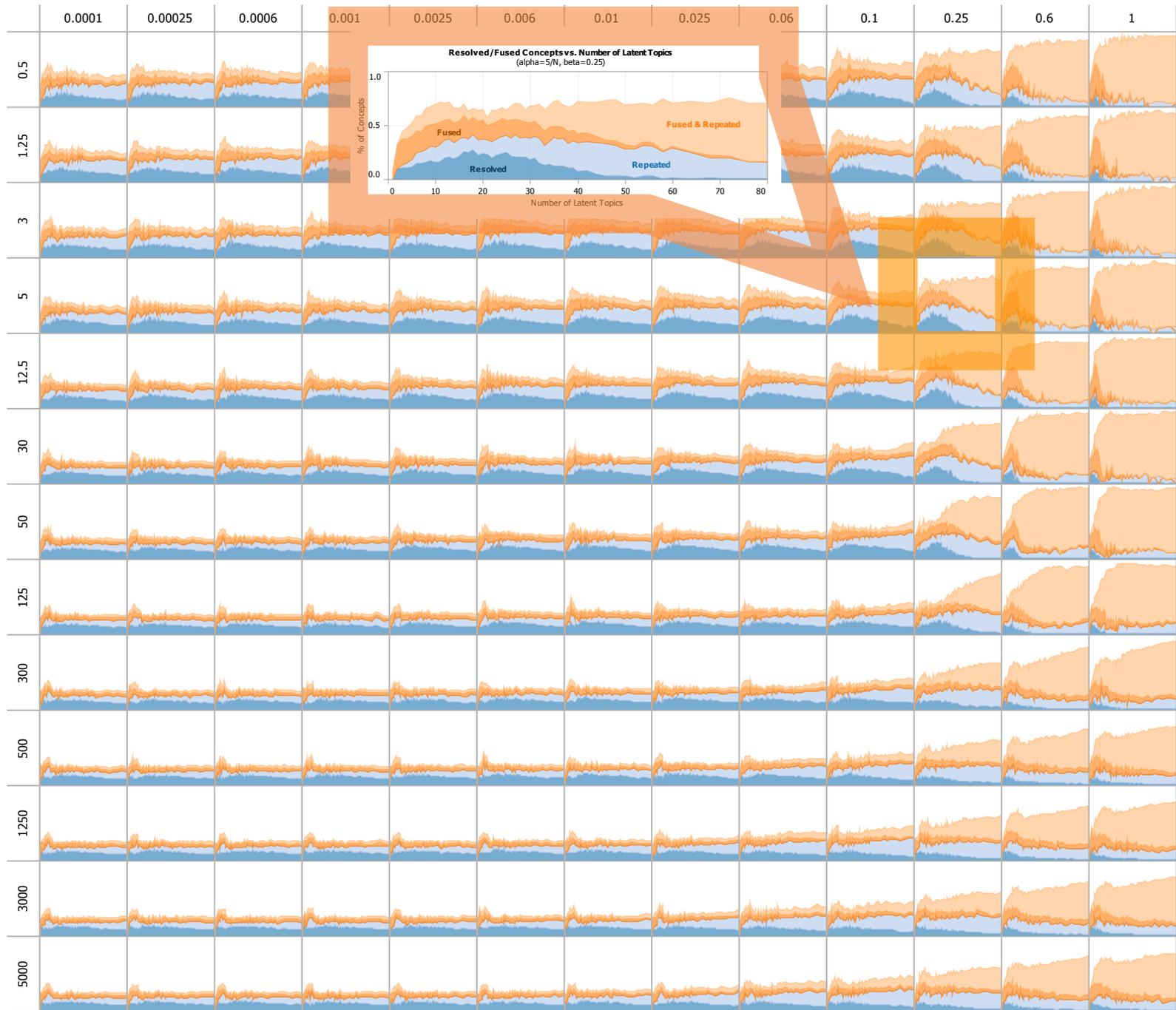
## Resolved/Fused Concepts vs. Number of Latent Topics

( $\alpha=5/N$ ,  $\beta=0.25$ )

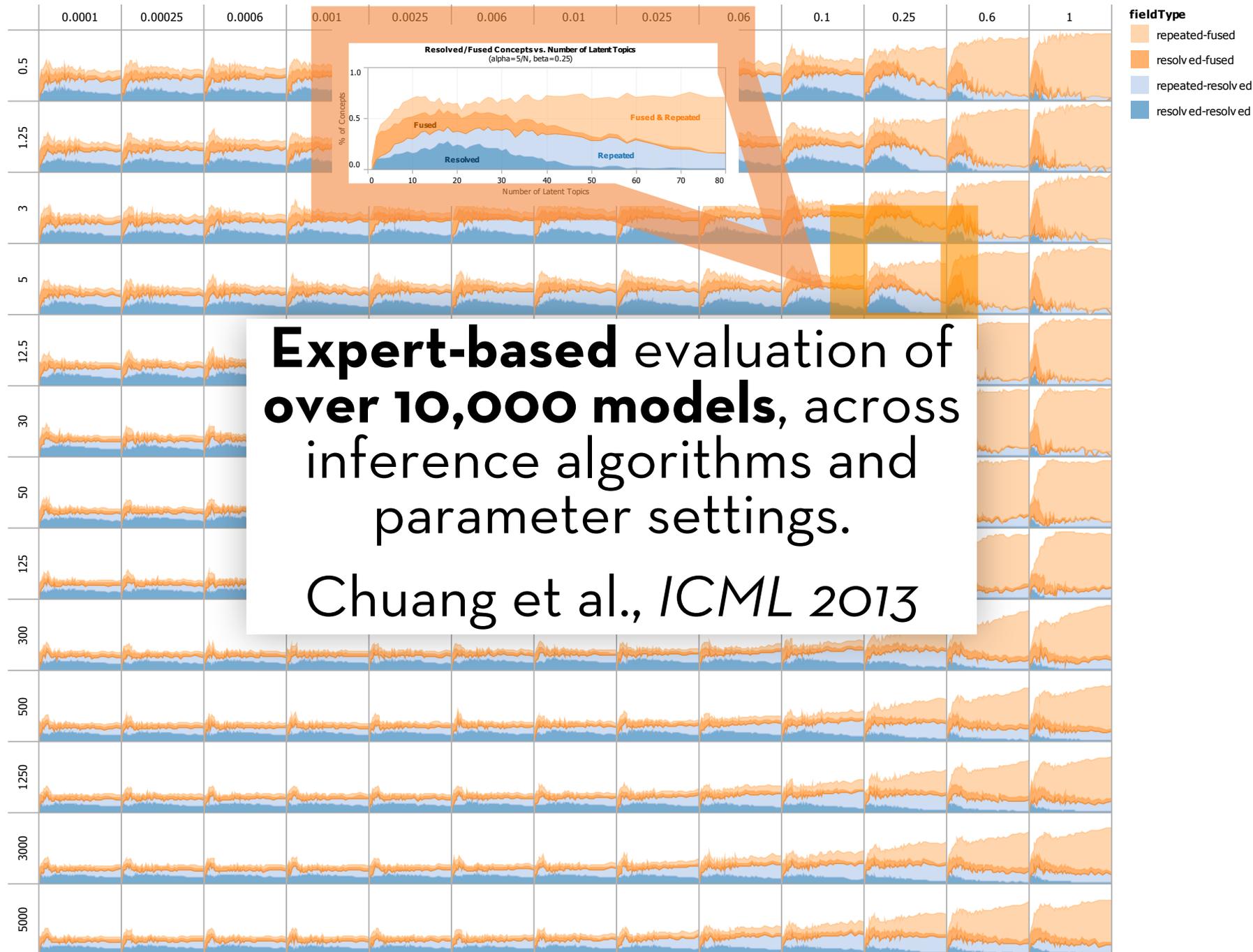




### Alpha vs Beta



Alpha vs Beta



# STEPPING BACK

# STEPPING BACK: A POST-GRADUATE EXERCISE

Your chosen field's biggest blind spot?

# STEPPING BACK: A POST-GRADUATE EXERCISE

Your chosen field's biggest blind spot?



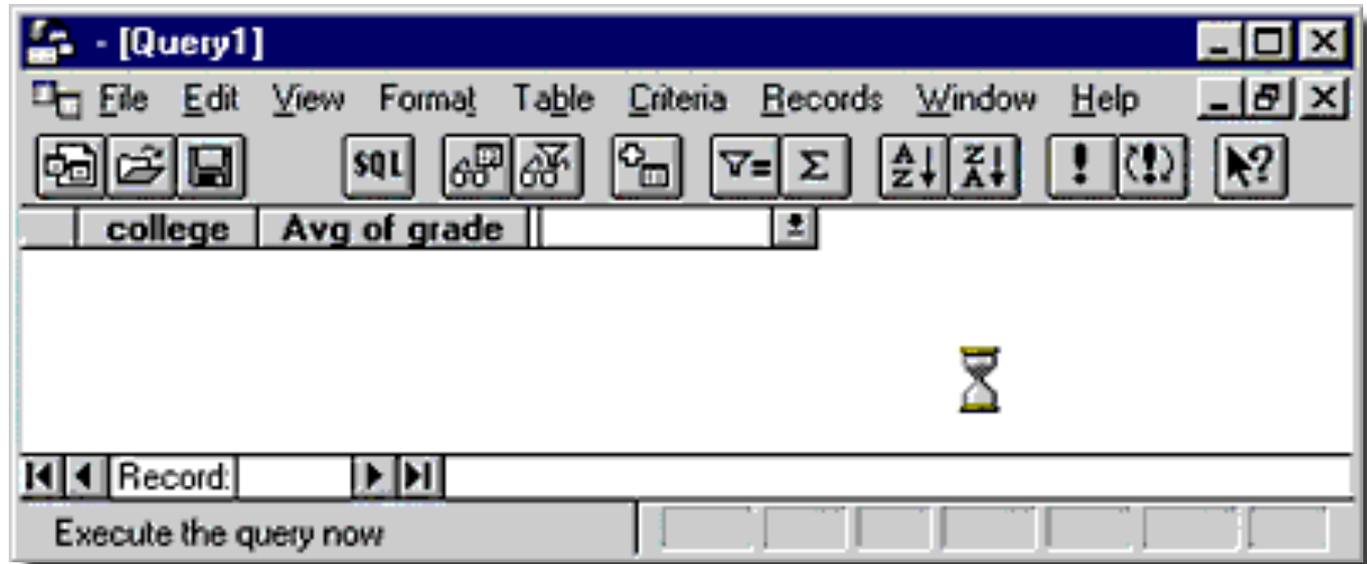
# STEPPING BACK: A POST-GRADUATE EXERCISE

Your chosen field's biggest blind spot?

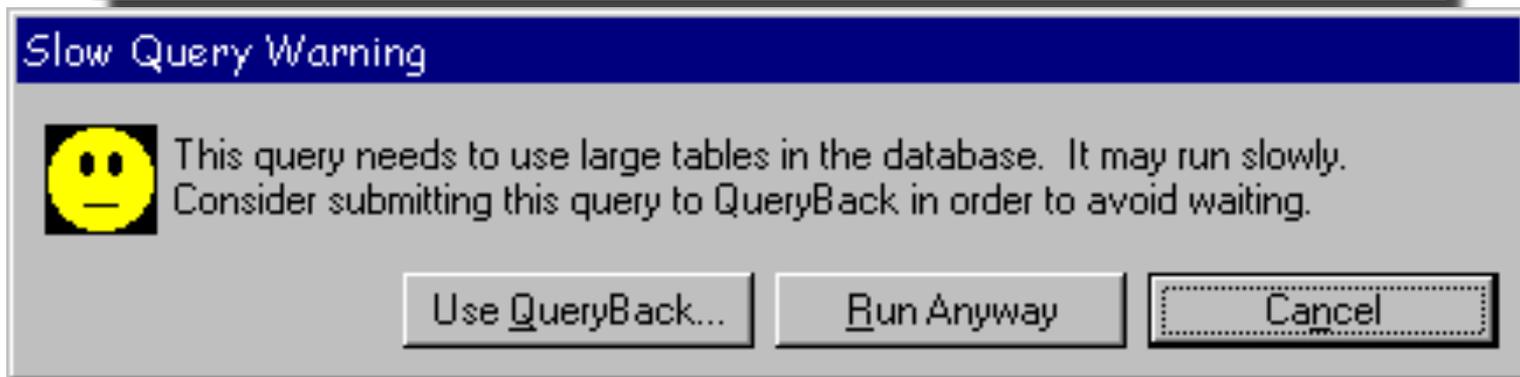
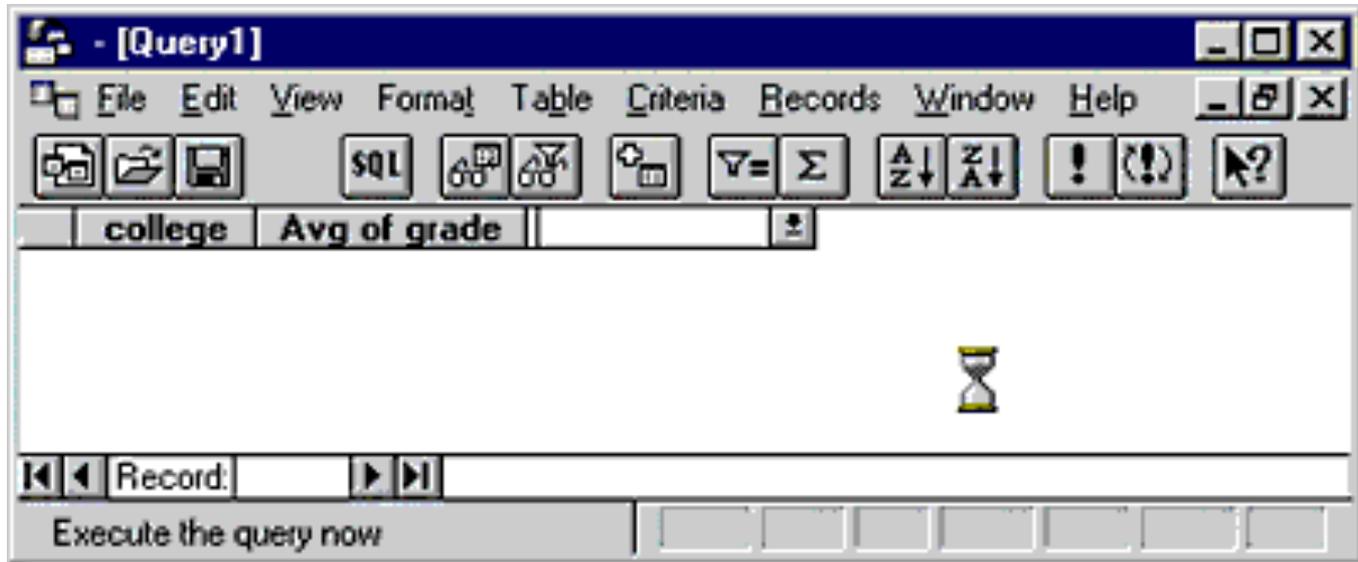


“A little disdain is not  
amiss;  
a little scorn is alluring.”  
— William Congreve, 1670–1729

# BATCH PROCESSING AND UI



# BATCH PROCESSING AND UI



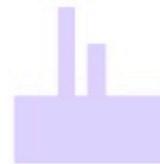
# THE CONTROL PROJECT: 1996-2001

<http://control.cs.berkeley.edu>



## CONTROL

- **On-Line processing of large datasets**
  - constant, useful feedback for long-running (data-intensive) operations
  - progressive refinement of answers
  - online user control
- **A blend of**
  - data processing, statistics, UI
  -



# THE CONTROL PROJECT: 1996-2001

<http://control.cs.berkeley.edu>



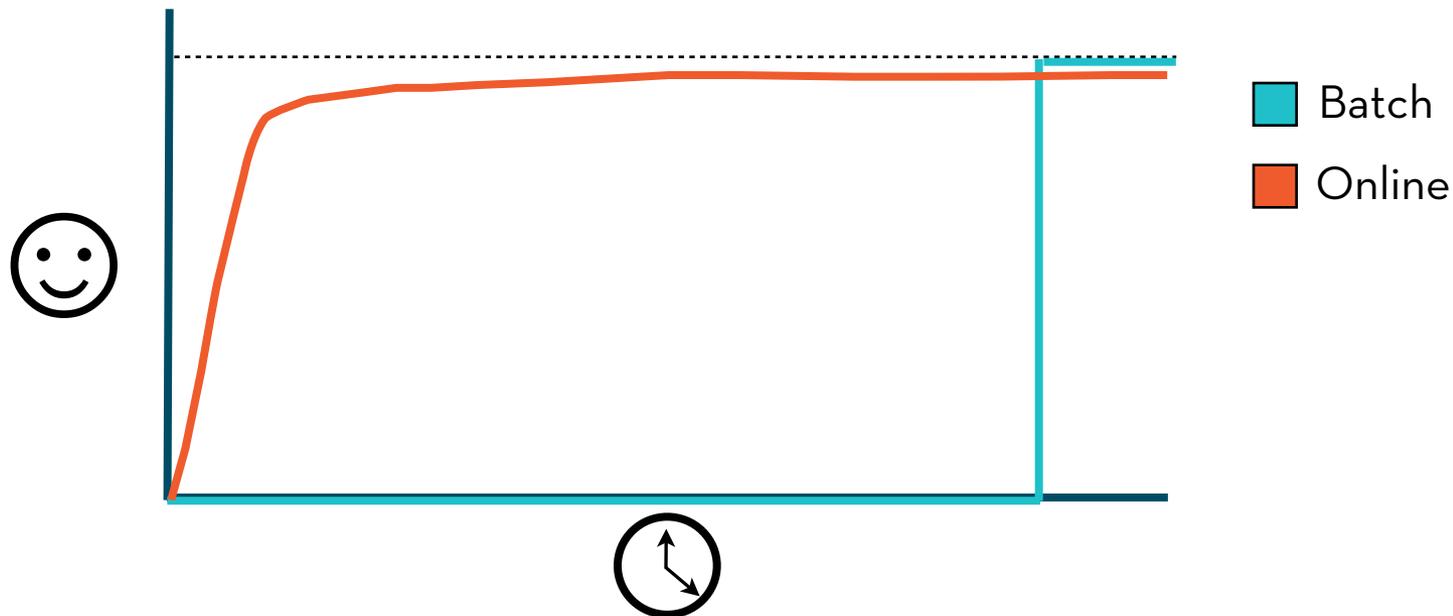
## CONTROL

- **On-Line processing of large datasets**
  - constant, useful feedback for long-running (data-intensive) operations
  - progressive refinement of answers
  - online user control
- **A blend of**
  - data processing, statistics, UI

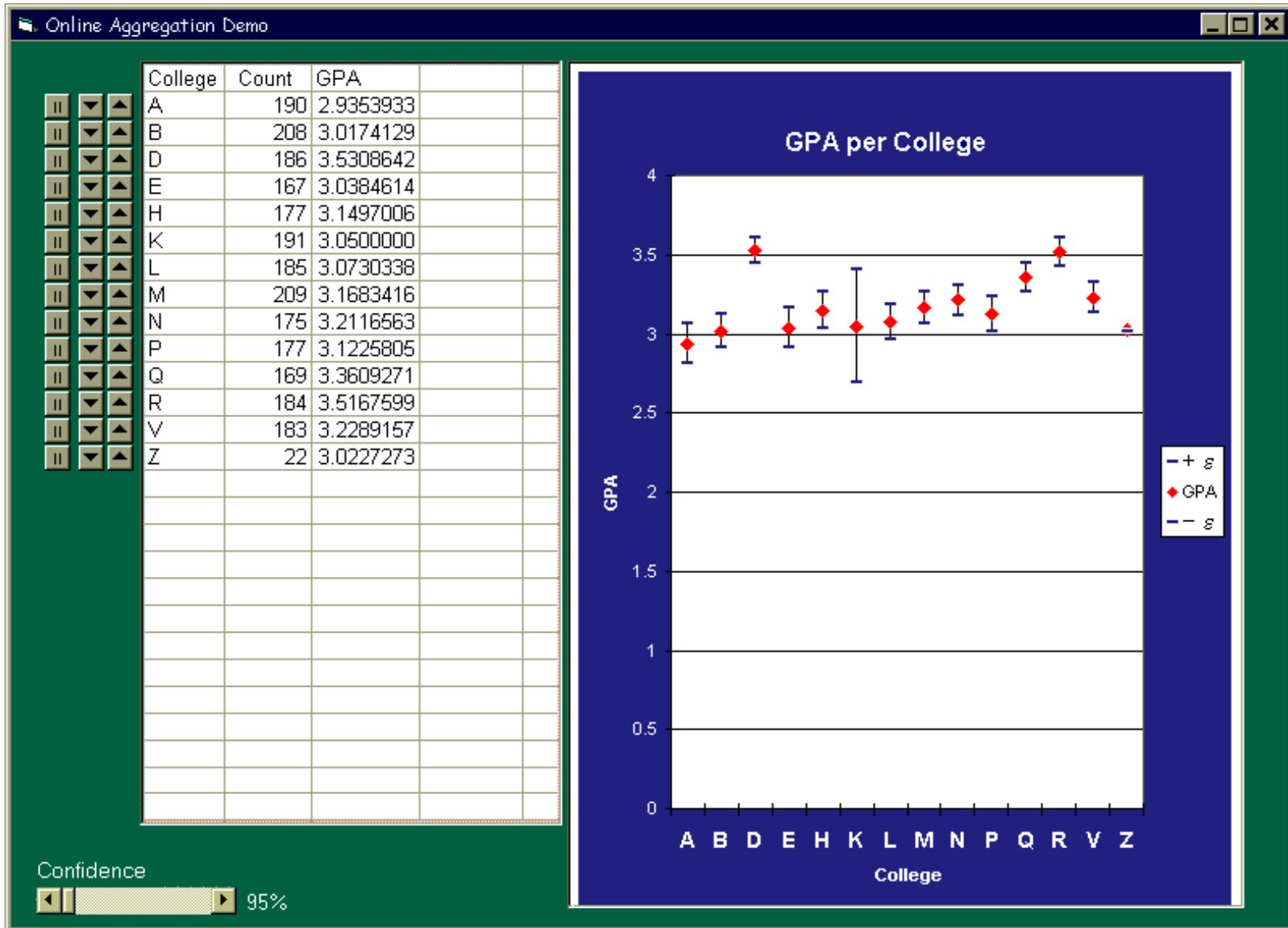


# GOALS FOR ONLINE PROCESSING

- Maximize 1st derivative of the “mirth index”
- Mirth subject to dynamic redefinition
- Need FEEDBACK and CONTROL



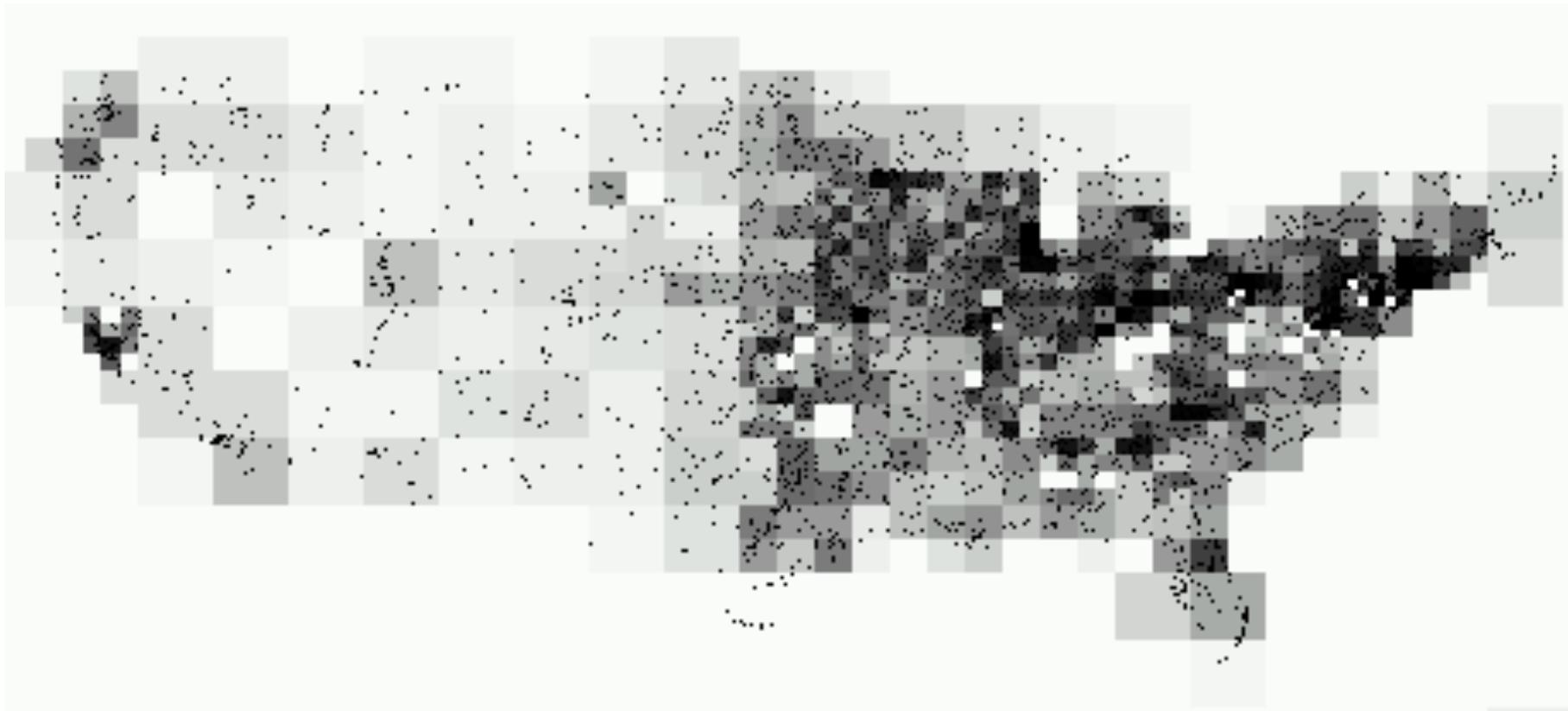
# ONLINE AGGREGATION [SIGMOD '97, '99]



# CLOUDS [IEEE COMPUTER 8/99]



# CLOUDS [IEEE COMPUTER 8/99]



# POTTER'S WHEEL [VLDB 01]

The screenshot shows a database application window with a menu bar (File, Classify, Cluster, Transform, Discrepancies, Sort, Show Buffer, 10000, 100%) and a data table. A context menu is open over the 'Split Column' option, which has further sub-options: 'Split Values' (with sub-options 'Split into 2 columns', 'Split by Example', and 'Split into many columns') and 'Divide Values'.

Delay	Carrier	Number	Source	Dest	Sch	Dept_Act	Arr_Sch	Arr_Act	Status	Random	
-4	AMERIC...	0039	SFO						NORMAL	101684	
42	AMERIC...	0849	ORD to						NORMAL	101898	
4	AMERIC...	0624	ORD	MCO	1997/04/...				NORMAL	101928	
0	AMERIC...	1291	ORD	MIA	1997/12/...	W	06:00		CANCEL...	102111	
8	AMERIC...	0407	ORD	SJC	1997/01/...	W	09:00	09:30 11:47	11:55	NORMAL	102172
-19	AMERIC...	2205	ORD to ...		1998/04/...	W	07:00	06:57 09:20	09:01	NORMAL	102203
-7	AMERIC...	1041	ORD to ...		1997/07/...	W	11:50	11:52 13:16	13:09	NORMAL	102416
13	AMERIC...	1112	ORD	BOS	1998/10/...	W	17:00	16:58 20:28	20:41	NORMAL	10284
-23	AMERIC...	2209	ORD to ...		1998/03/...	W	08:00	07:57 10:31	10:08	NORMAL	102844
-19	AMERIC...	1765	ORD to L...		1997/09/...	W	08:40	08:36 11:16	10:57	NORMAL	103027
-6	AMERIC...	2366	ORD	ROC	1998/05/...	W	20:40	20:41 23:22	23:16	NORMAL	103118
30	AMERIC...	0265	JFK	SEA	1997/10/...	W	08:20	08:17 11:07	11:37	NORMAL	103515
15	AMERIC...	2267	ORD	DFW	1998/02/...	W	18:15	18:41 20:50	21:05	NORMAL	10375
0	AMERIC...	1891	ORD to L...		1998/03/...	W	13:10	13:10 15:42	15:42	CANCEL...	103790
-3	AMERIC...	1754	ORD to ...		1997/01/...	W	16:10	16:11 19:07	19:04	NORMAL	100280
-7	AMERIC...	0218	ORD to ...		1997/07/...	W	06:50	06:45 09:45	09:38	NORMAL	104583
33	AMERIC...	1565	ORD	SNA	1998/11/...	W	19:00	19:00 21:20	21:53	NORMAL	104675
58	AMERIC...	1984	ORD	ROC	1997/10/...	W	10:08	11:25 12:47	13:45	NORMAL	100280
-12	AMERIC...	1609	ORD to ...		1997/11/...	W	18:35	18:34 20:04	19:52	NORMAL	104919
57	AMERIC...	0552	ORD to ...		1998/03/...	W	13:25	14:21 17:25	18:22	NORMAL	105041
-7	AMERIC...	1536	ORD to L...		1998/08/...	W	10:25	10:23 13:30	13:23	NORMAL	105102
9	AMERIC...	1754	ORD to ...		1997/02/...	W	16:10	16:09 18:56	19:05	NORMAL	105194
-20	AMERIC...	1856	ORD to ...		1997/02/...	W	06:20	06:17 09:12	08:52	NORMAL	105499
7	AMERIC...	0655	JFK	STT	1998/01/...	W	08:00	08:01 12:44	12:51	NORMAL	105651
0	AMERIC...	1285	ORD	SFO	1997/09/...	W	15:00	15:04 17:29	17:29	NORMAL	106109
-48	AMERIC...	0117	JFK to LAX		1998/09/...	W	15:00	14:58 17:49	17:01	NORMAL	106170
-4	AMERIC...	0267	ORD to T...		1997/08/...	W	19:25	19:22 21:12	21:08	NORMAL	106231

# HADOOP WINTER

# HADOOP WINTER AND A NEW SPRING

# HADOOP WINTER AND A NEW SPRING

## Scalable Approximate Query Processing With The DBO Engine

Christopher Jermaine, Subramanian Arumugam, Abhijit Pol, Alin Dobra  
CISE Department, University of Florida  
Gainesville, FL, USA  
{cjermain, sa2, apol, adobra}@cise.ufl.edu

## MapReduce Online

Tyson Condie, Neil Conway, Peter Alvaro, Joseph M. Hellerstein  
*UC Berkeley*

Khaled Elmeleegy, Russell Sears  
*Yahoo! Research*

## Online Aggregation for Large MapReduce Jobs

Niketan Pansare<sup>1</sup>, Vinayak Borkar<sup>2</sup>, Chris Jermaine<sup>1</sup>, Tyson Condie<sup>3</sup>  
<sup>1</sup>Rice University, <sup>2</sup>UC Irvine, <sup>3</sup>Yahoo! Research  
np6@rice.edu, vborkar@ics.uci.edu, cmj4@rice.edu, tcondie@yahoo-inc.com

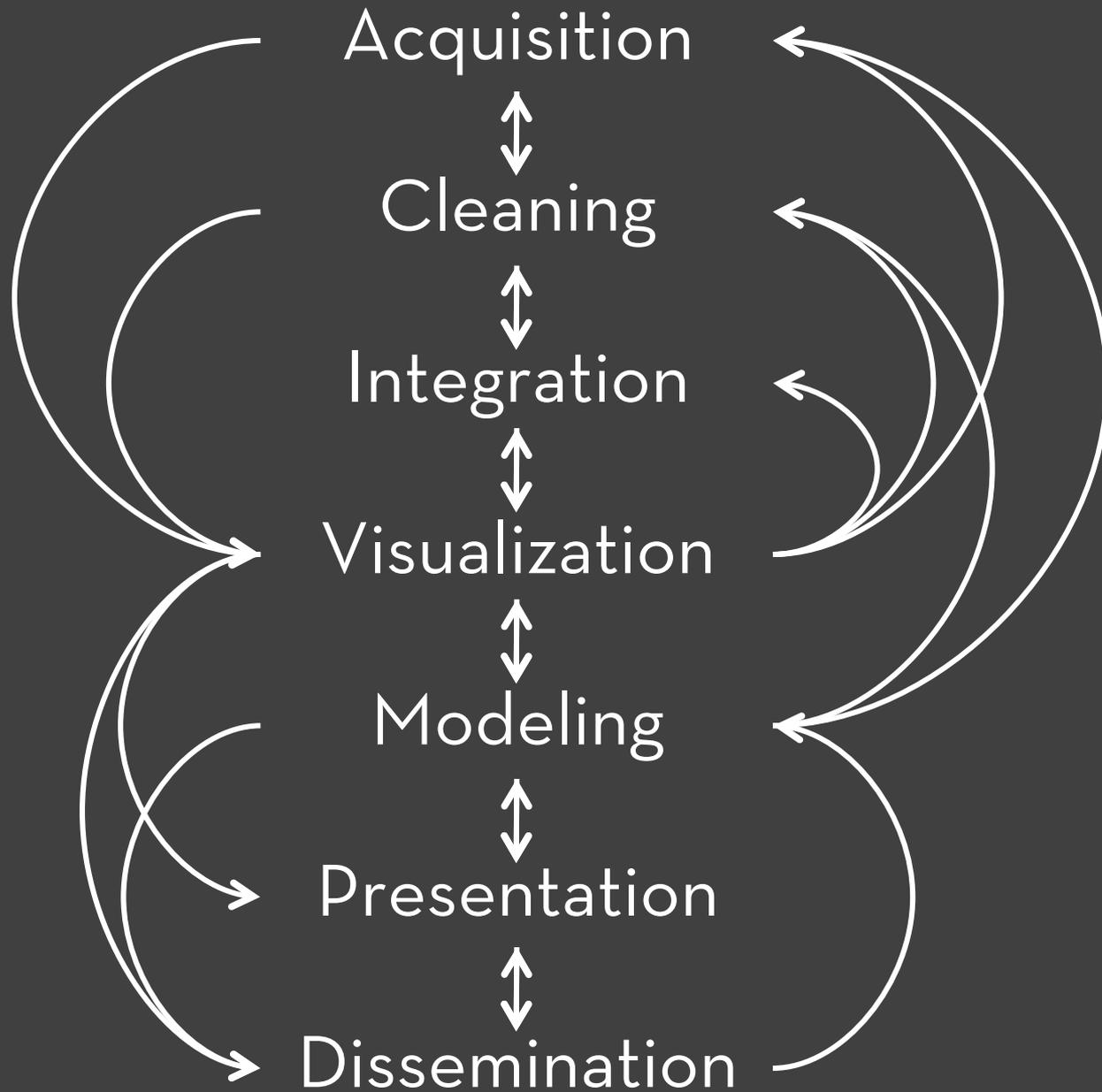
## Incremental, Approximate Database Queries and Uncertainty for Exploratory Visualization

Danyel Fisher  
Microsoft Research



## Stat! – An Interactive Analytics Environment for Big Data

Mike Barnett<sup>1</sup>, Badrish Chandramouli<sup>1</sup>, Robert DeLine<sup>1</sup>, Steven Drucker<sup>1</sup>, Danyel Fisher<sup>1</sup>,  
Jonathan Goldstein<sup>1</sup>, Patrick Morrison<sup>2</sup>, John Platt<sup>1</sup>  
<sup>1</sup>Microsoft Research  
Redmond, Washington, USA  
{mbarnett, badrishc, rdeline, sdrucker,  
danyelf, jongold, jplatt}@microsoft.com  
<sup>2</sup>North Carolina State University  
Raleigh, North Carolina, USA  
pjmorris@ncsu.edu



# ENGAGING PRACTITIONERS

To gain insight, we interviewed **35 analysts**:

# ENGAGING PRACTITIONERS

To gain insight, we interviewed **35 analysts**:

*25 Companies*

Healthcare

Retail, Marketing

Social networking

Media

Finance, Insurance

*Various titles*

Data analyst

Data scientist

Software engineer

Consultant

Chief technical officer

# ENGAGING PRACTITIONERS

To gain insight, we interviewed **35 analysts**:

*25 Companies*

Healthcare

Retail, Marketing

Social networking

Media

Finance, Insurance

*Various titles*

Data analyst

Data scientist

Software engineer

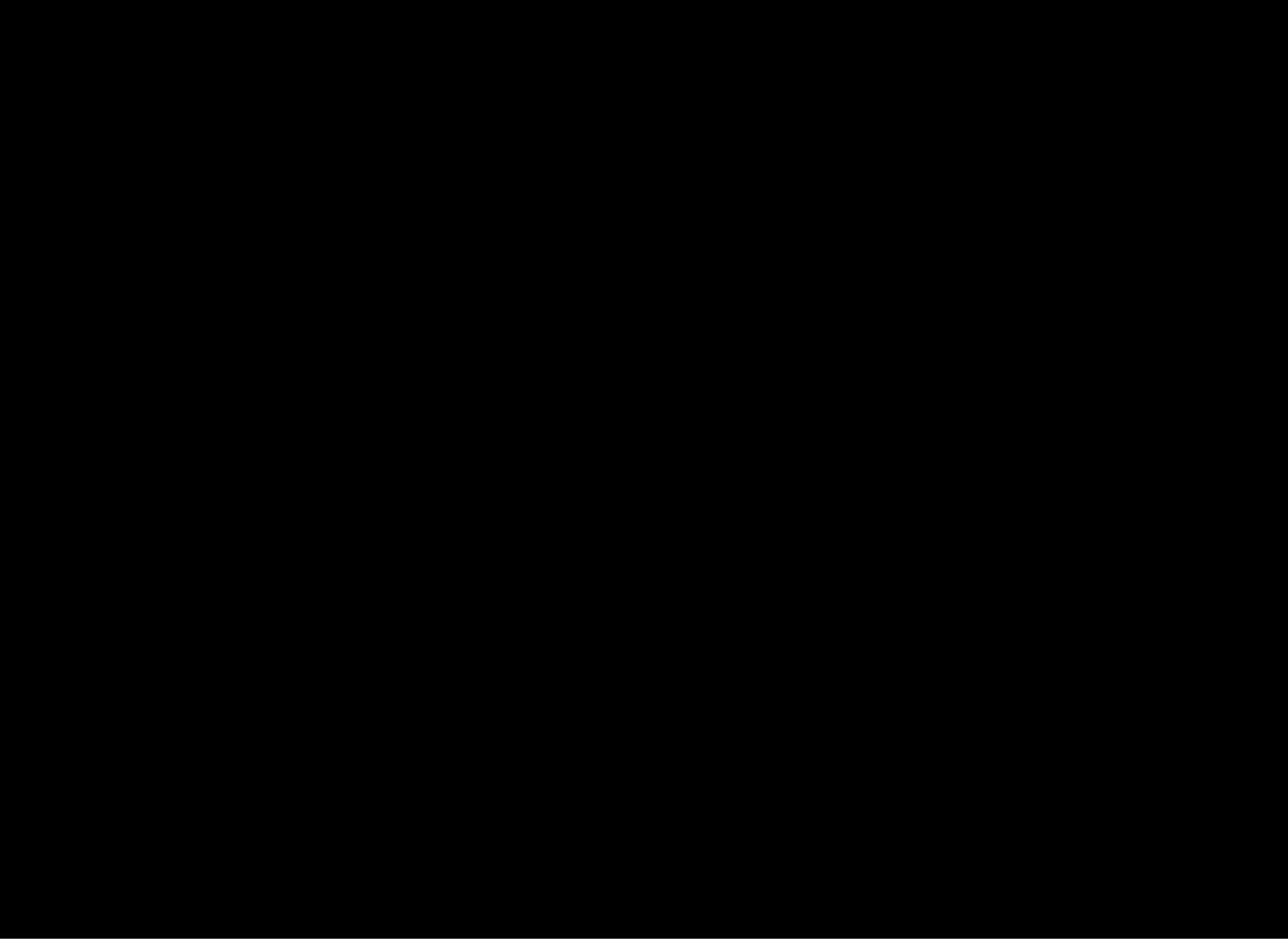
Consultant

Chief technical officer

**Enterprise Data Analysis and Visualization: An Interview Study**

Sean Kandel, Andreas Paepcke, Joseph Hellerstein, Jeffrey Heer  
*IEEE Visual Analytics Science & Technology (VAST), 2012*





“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.”

“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.”

**Lost productivity**

“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.”

**Lost productivity**

“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.”

## **Lost productivity**

“Most of the time once you transform the data ... the insights can be scarily obvious.”

“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.”

**Lost productivity**

“Most of the time once you transform the data ... the insights can be scarily obvious.”

**Lost accessibility**

“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.”

**Lost productivity**

“Most of the time once you transform the data ... the insights can be scarily obvious.”

**Lost accessibility**

“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.”

**Lost productivity**

“Most of the time once you transform the data ... the insights can be scarily obvious.”

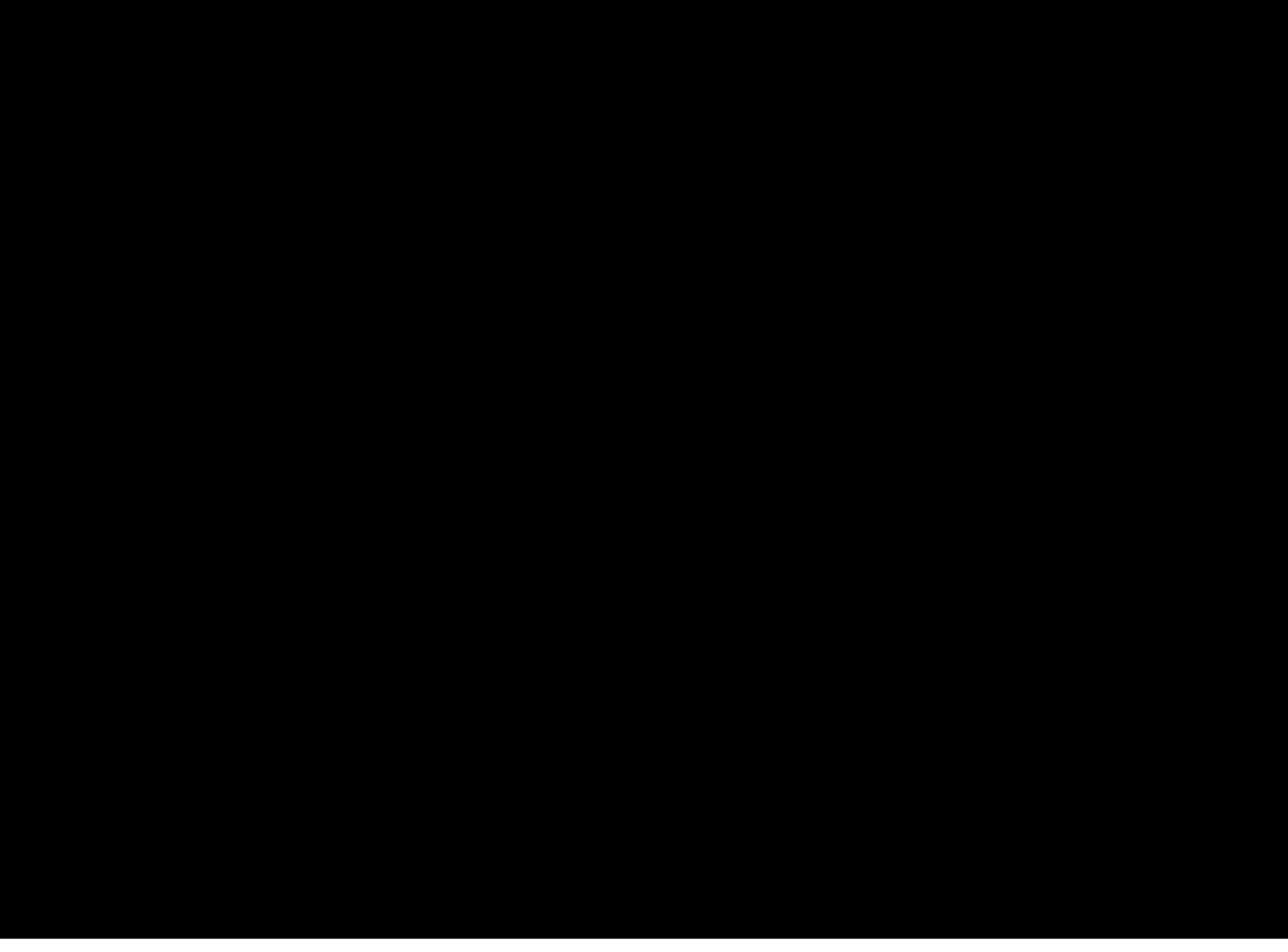
**Lost accessibility**

“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.”

**Lost productivity**

“Most of the time once you transform the data ... the insights can be scarily obvious.”

**Lost accessibility**



“It’s easy to just think you know what you are doing and not look at data at every intermediary step.

An analysis has 30 different steps. It’s tempting to just do this then that and then this. You have no idea in which ways you are wrong and what data is wrong.”

## **Interactivity and Visualization**



TRIFACTA  
PEOPLE + DATA + COMPUTATION



TRIFACTA  
PEOPLE + DATA + COMPUTATION

Analytic Productivity

Remove drudgery, restore time



TRIFACTA  
PEOPLE + DATA + COMPUTATION

## Analytic Productivity

Remove drudgery, restore time

## Data Accessibility

Enable self-service data manipulation



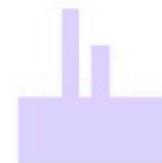
# TRIFACTA

PEOPLE + DATA + COMPUTATION

⋮

## CONTROL

- **On-Line processing of large datasets**
  - constant, useful feedback for long-running (data-intensive) operations
  - progressive refinement of answers
  - online user control
- **A blend of**
  - data processing, statistics, UI





# TRIFACTA

PEOPLE + DATA + COMPUTATION

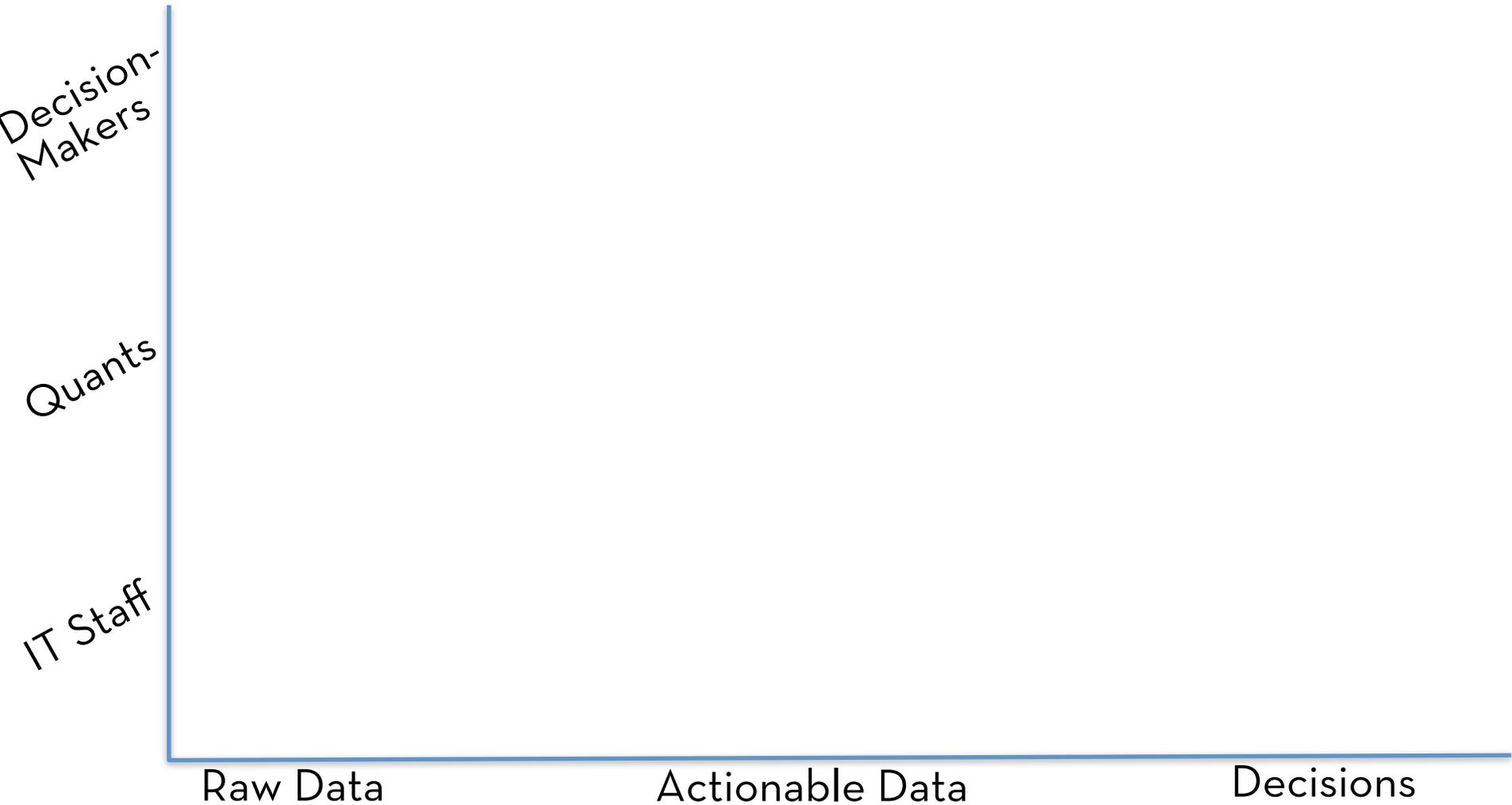
⋮

**CONTROL**

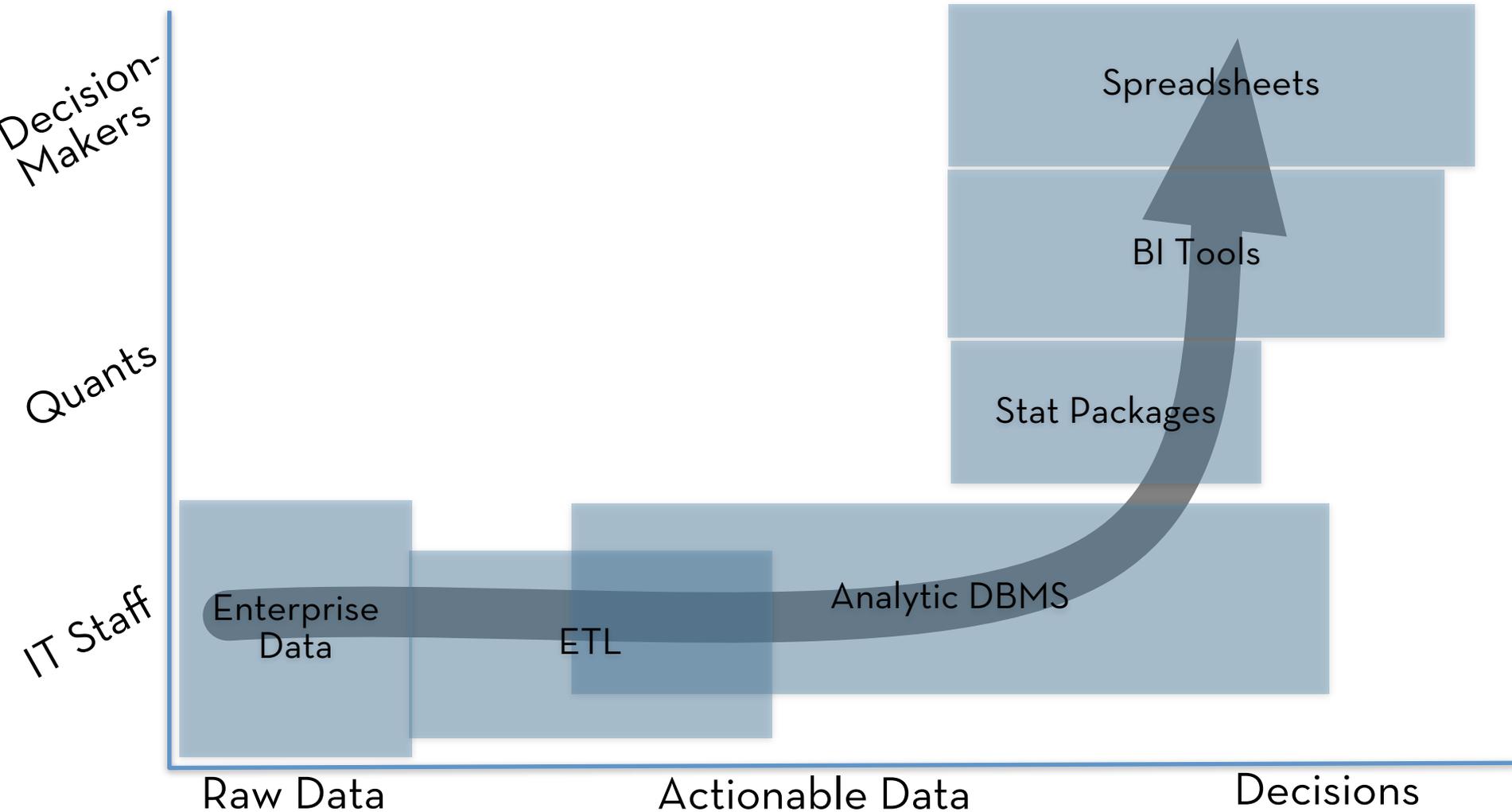
- **On-Line processing of large datasets**
  - constant, useful feedback for long-running (data-intensive) operations
  - progressive refinement of answers
  - online user control
- **A blend of**
  - data processing, statistics, UI

⋮

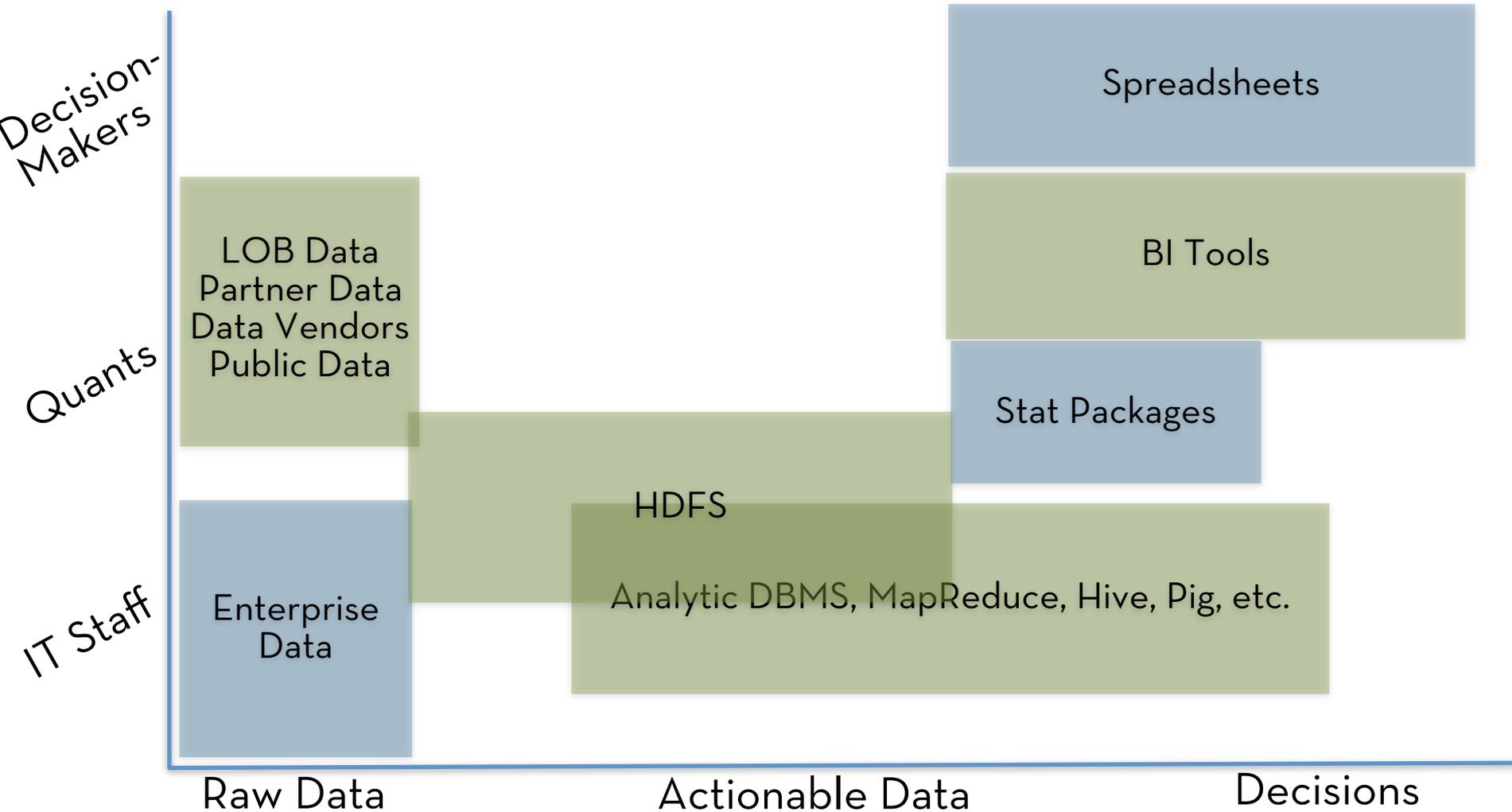
# TRADITIONAL LANDSCAPE



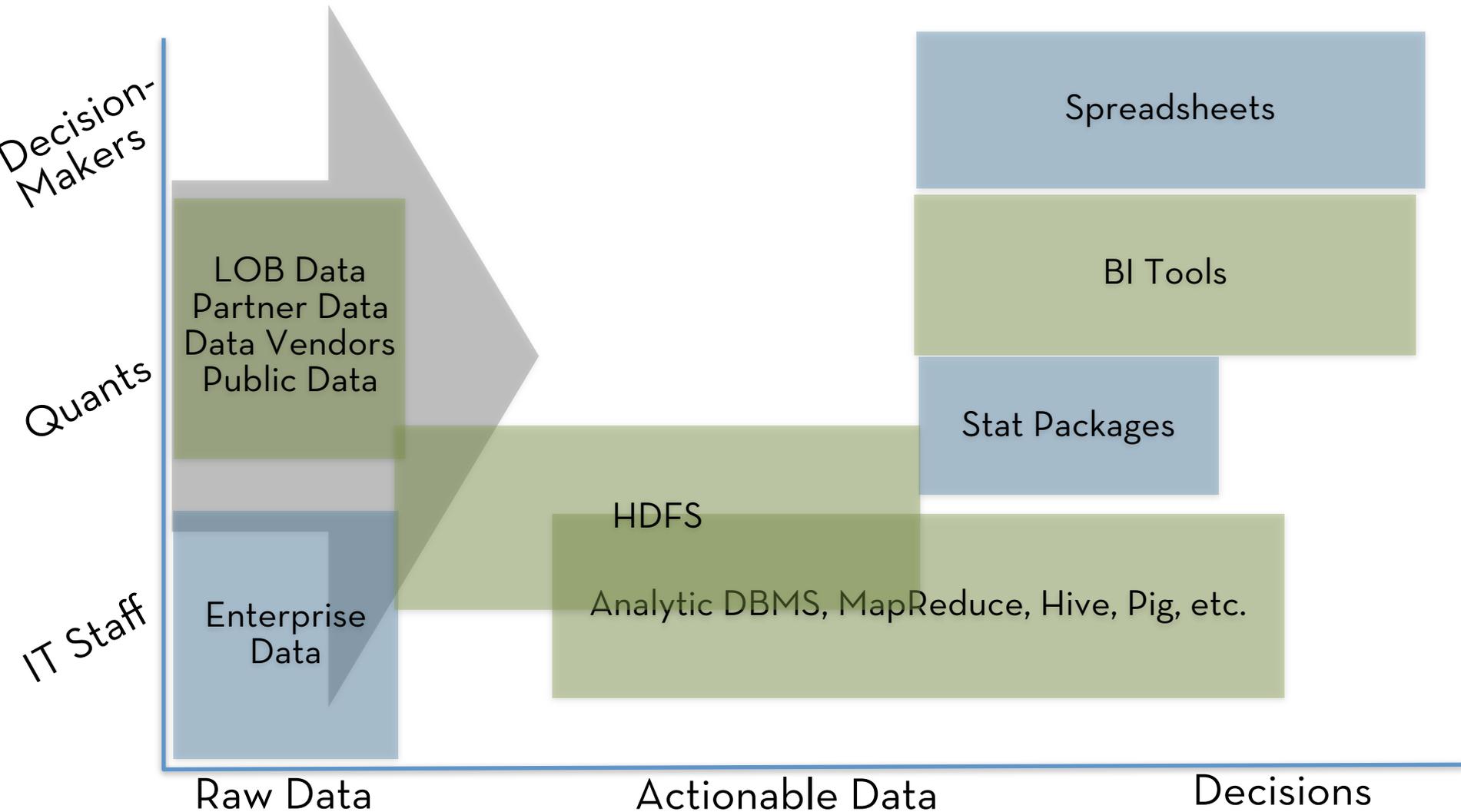
# TRADITIONAL LANDSCAPE



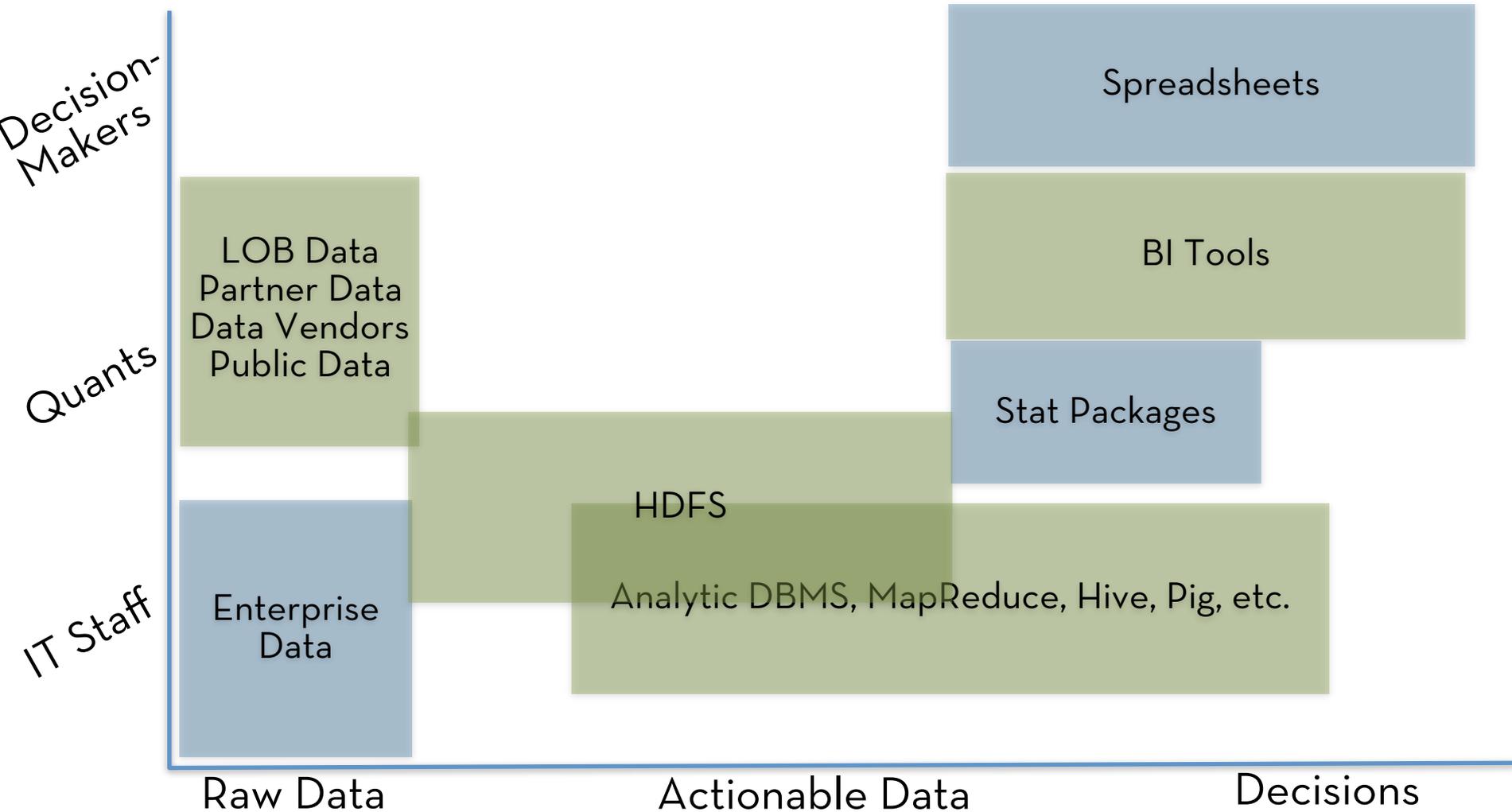
# NEW LANDSCAPE



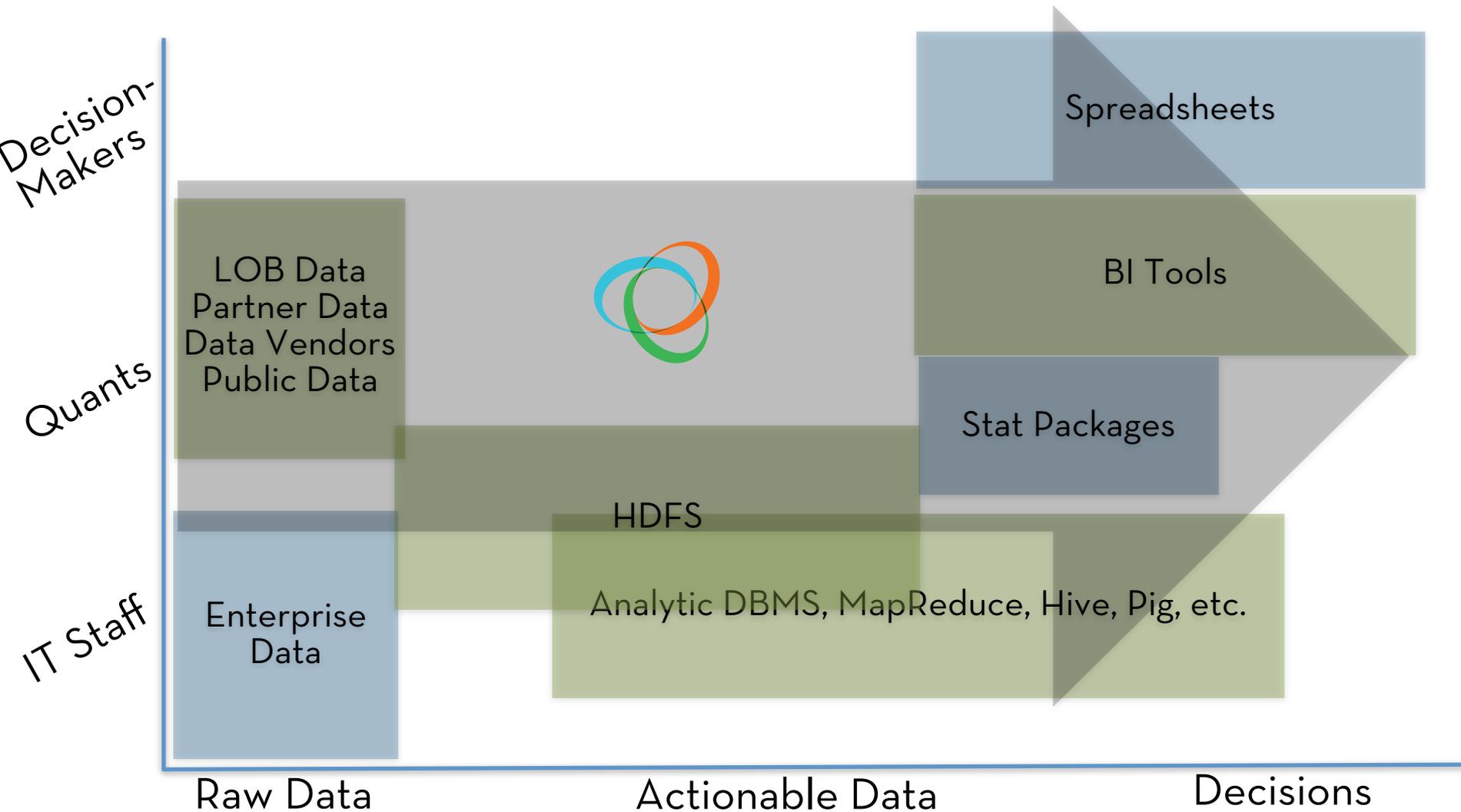
# NEW LANDSCAPE



# NEW LANDSCAPE



# NEW LANDSCAPE



PEOPLE + DATA + COMPUTATION