

Text Visualization



Why Visualize Text?

Understanding – get the “gist” of a document

Grouping – cluster for overview or classification

Comparison – compare document collections, or inspect evolution of collection over time

Correlation – compare patterns in text to those in other data, e.g., correlate with social network

What is Text Data?

Documents

Articles, books and novels

E-mails, web pages, blogs

Tags, comments

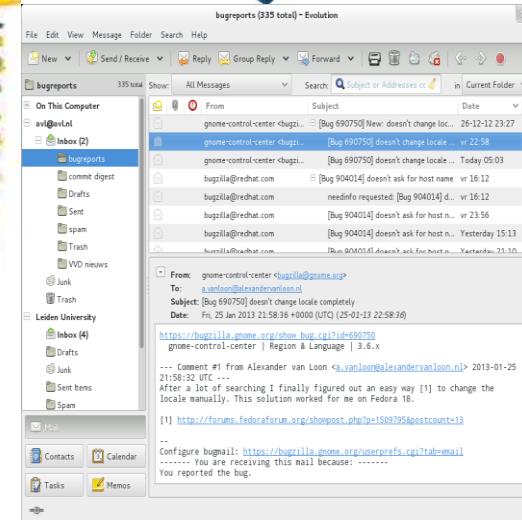
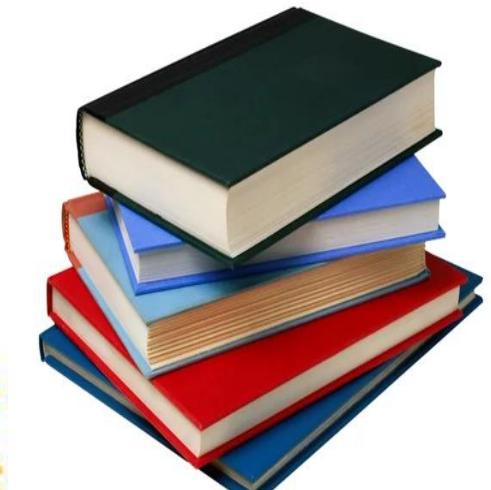
Computer programs, logs

Collections of Documents

Messages (e-mail, blogs, tags, comments)

Social networks (personal profiles)

Academic collaborations (publications)



Example: Text Analysis and Visualization

Scenario of Malaysia's Economy

Text Data

News articles

Market reports

What questions might you want to answer?
What visualizations might help?

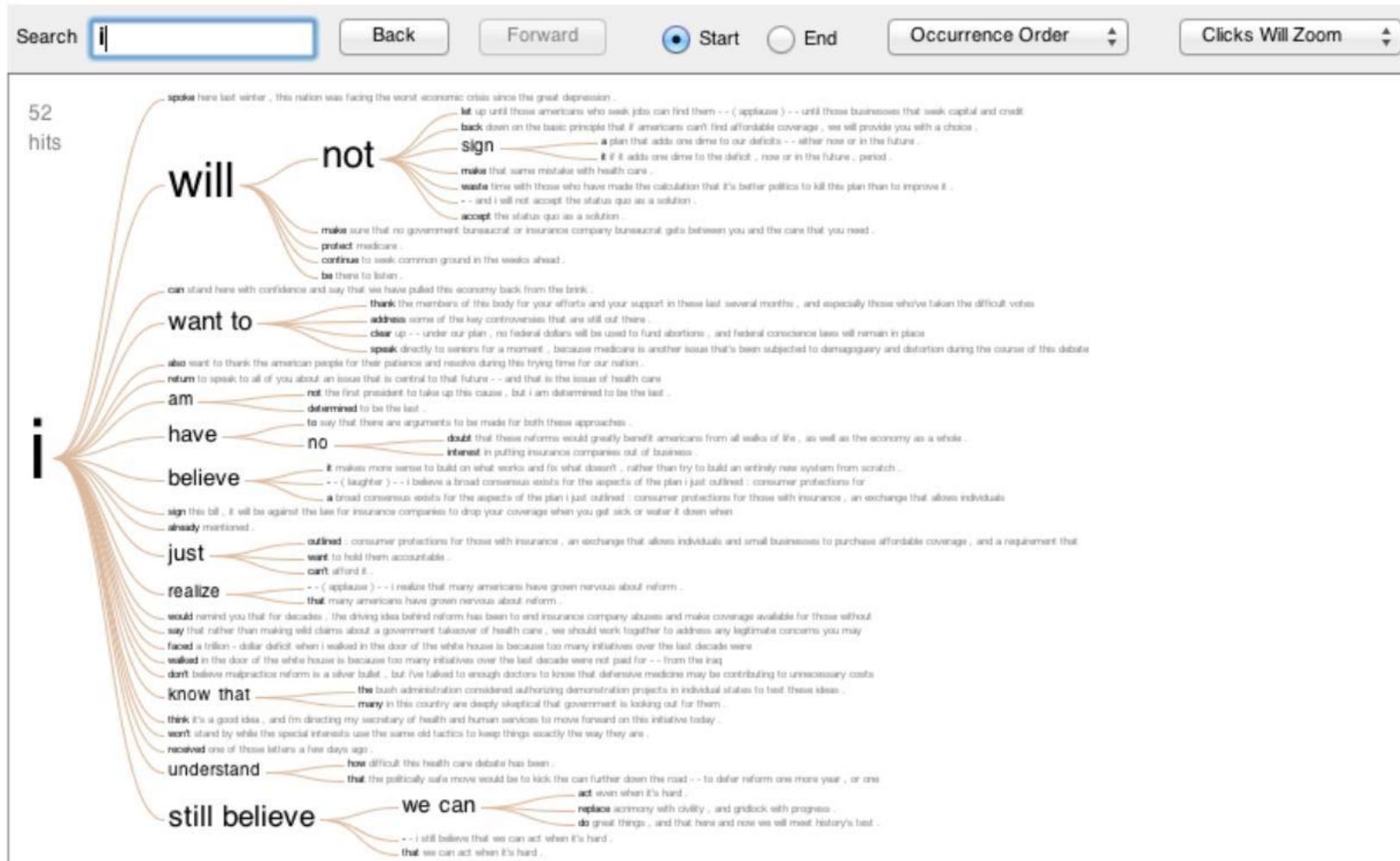
Tag Clouds: Word Count

Malaysia Economy



Word Tree: Word Sequences

Visualizations : Word Tree President Obama's Address to Congress on Health Care



Unstructured Data and Degree of Evaluation

Many text visualizations do not represent the text directly. They represent the output of a **language model** (word counts, word sequences, term links, term trends, etc.).

- Can you interpret the visualization? How well does it convey the properties of the model?
- Do you trust the model? How does the model enable us to reason about the text?

Text Visualization Challenges

High Dimensionality

Where possible use text to represent text...
... which terms are the most descriptive?



Context & Semantics

Provide relevant context to aid understanding.
Show (or provide access to) the source text.

Modeling Abstraction

Determine your analysis task.
Understand abstraction of your language models.
Match analysis task with appropriate tools and models.

Text as Data

Words as nominal data?

High dimensional (10,000+)

More than equality tests

Words have meanings and relations

- Correlations: *Hong Kong, San Francisco, Bay Area*
- Order: *April, February, January, June, March, May*
- Membership: *Tennis, Running, Swimming, Hiking, Piano*
- Hierarchy, antonyms & synonyms, entities, ...

Text Pre-processing

1. Tokenization

Segment text into terms.

Remove stop words? *a, an, the, of, to, be*

Numbers and symbols? *#gocard, @stanfordfbball, Beat Cal!!!!!!!*

Entities? *San Francisco, O'Connor, U.S.A.*

Text Pre-processing

1. Tokenization

Segment text into terms.

Remove stop words? *a, an, the, of, to, be*

Numbers and symbols? *#gocard, @stanfordfbball, Beat Cal!!!!!!!*

Entities? *San Francisco, O'Connor, U.S.A.*

2. Stemming

Group together different forms of a word.

Porter stemmer? *visualization(s), visualize(s), visually -> visual*

Lemmatization? *goes, went, gone -> go*

Text Pre-processing

1. Tokenization

Segment text into terms.

Remove stop words? *a, an, the, of, to, be*

Numbers and symbols? *#gocard, @stanfordball, Beat Cal!!!!!!*

Entities? *San Francisco, O'Connor, U.S.A.*

2. Stemming

Group together different forms of a word.

Porter stemmer? *visualization(s), visualize(s), visually -> visual*

Lemmatization? *goes, went, gone -> go*

3. Ordered list of terms

Bag of Words Model

Ignore ordering relationships within the text

A document \approx vector of term weights

- Each dimension corresponds to a term (10,000+)
- Each value represents the relevance
 - For example, simple term counts

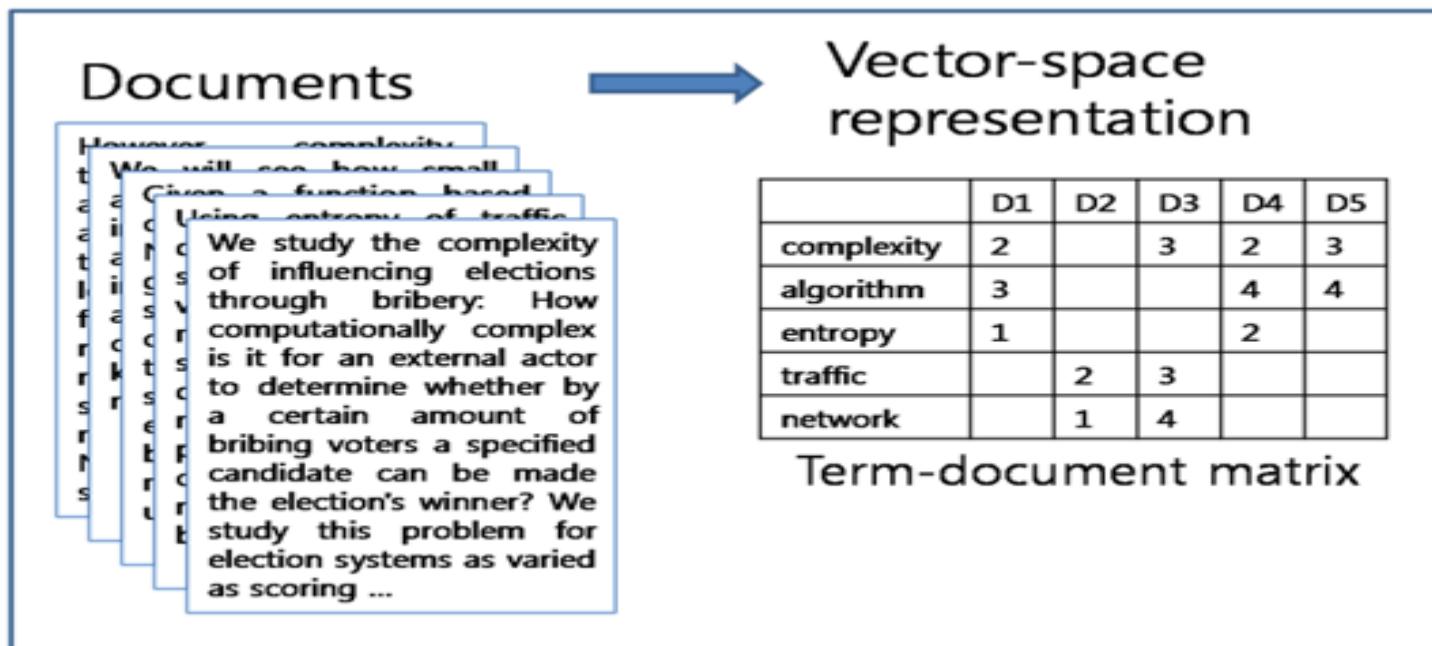
Aggregate into a document-term matrix

- Document vector space model

Document-Term Matrix

Each document is a vector of term weights

Simplest weighting is to just count occurrences



Visualizations : Wordle of Sarah Palin RNC 9/3/2008 Speech

Creator: Anonymous

Tags:

Edit Language Font Layout Color



Data file: Sarah Palin speaks at the Republican National Convention, 9/3/2008

Data source: SFGate / AP



This data set
has not yet been rated

Tag Clouds or Word Clouds

Strengths

Can help with gisting and initial query formation.

Weaknesses

Sub-optimal visual encoding (size vs. position)

Inaccurate size encoding (long words are bigger)

May not facilitate comparison (unstable layout)

Term frequency may not be meaningful

Does not show the structure of the text

Tips: Descriptive Phrases

Understand the limitations of your language model.

Bag of words:

- Easy to compute

- Single words

- Loss of word ordering

Select appropriate model and visualization

- Generate longer, more meaningful phrases

- Adjective-noun word pairs for reviews

- Show keyphrases within source text

Document Content

Information Retrieval

Search for documents

Match query string with documents

Visualization to **contextualize results**

The screenshot shows a Google Scholar search results page. The search bar at the top contains the query "acronym resolution". Below the search bar are various filters: "Scholar" (selected), "Articles and patents", "anytime", "include citations", and a "Create email alert" button. The search results are listed below:

- A supervised learning approach to acronym identification**
D Nadeau, P Turney - The Eighteenth Canadian ..., 2005 - nparc.cisti-icist.nrc-cnrc.gc.ca
... Recently the fields of Genetics and Medicine have become especially interested in acronym resolution (Pustejovsky et al., 2001, Yu et al. 2002). ... Pustejovsky et al.'s acronym resolution technique searches for definitions of acronyms within noun phrases. ...
[Cited by 48](#) - [Related articles](#) - [All 16 versions](#)
- Biomedical term mapping databases**
JD Wren, JT Chang, J Pustejovsky ... - Nucleic acids ..., 2005 - Oxford Univ Press
... the prevalence of polyonyms, or acronyms with multiple definitions. An important part of any high-throughput effort to tie experimental findings to published knowledge within the scientific literature involves acronym resolution. ...
[Cited by 41](#) - [Related articles](#) - [All 22 versions](#)
- Anthropogenic climate change over the Mediterranean region simulated by a global variable resolution model**
AL Gibelin... - Climate Dynamics, 2003 - Springer
... The long simulations CC and CS are split into two 30-year datasets CC1 and CS1 for the period 1960–1989 and CC2 and CS2 for the period 2070–2099 Full name Acronym Resolution Period Coupled Coupled control CC T63 1950–2099 Yes ...
[Cited by 197](#) - [Related articles](#) - [BL Direct](#) - [All 5 versions](#)
- Metaphrase: an aid to the clinical conceptualization and formalization of patient problems in healthcare enterprises.**
MS Tuttle, NE Olson, KD Keck, WG Cole... - Methods of information ..., 1998 - ukpmc.ac.uk
... Title not supplied (PMID:10566483). Concept definition and manipulation are supported through

The 2007 State of the Union Address

Over the years, President Bush's State of the Union address has averaged almost 5,000 words each, meaning the the President has delivered over 34,000 words. Some words appear frequently while others appear only sporadically. Use the tools below to analyze what Mr. Bush has said.

Search

or choose a word here.

Use of the phrase "Tax" in past State of the Union Addresses

2001*	2002	2003	2004	2005	2006	2007
29	7	13	-	11	10	10



Next Instance of 'Tax'

The word in context

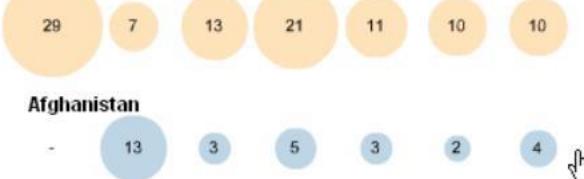
I believe in local control of schools. We should not, and we will not, run public schools from Washington, D.C. Yet when the federal government spends **TAX** dollars, we must insist on results. Children should be tested on basic reading and math skills every year between grades three and eight. Measuring is the only way to know whether all our children are learning. And I want to know, because I refuse to leave any child behind in America.

-- 2001 (Paragraph 14 of 73)

Compared with other words

2001*	2002	2003	2004	2005	2006	2007	
Tax	29	7	13	21	11	10	10

Afghanistan



Economy(ic)



Insurance



Iraq/Iraqi(s)



Iran



Oil



Social Security



* As a newly elected president, Mr. Bush did not deliver a formal State of the Union address in 2001. His Feb. 27 speech to a joint session of Congress was analogous to the State of the Union, but without the title.

Concordance

What is the common local context of a term?

The screenshot shows the Larkin.Concordance application window. The menu bar includes File, Text, Search, Edit, Headwords, Contexts, View, Tools, and Help. The toolbar contains icons for file operations like Open, Save, Print, and zoom. The main area has two panes: a left pane listing headwords with their counts, and a right pane displaying contextual snippets for the selected word 'HEART'.

Headword	No.
HEAR	15
HEARD	9
HEARING	7
HEARS	3
HEARSE	1
HEART	25
HEART'S	2
HEART-SHAPED	1
HEARTH	1
HEARTS	7
HEARTY	1
HEAT	6
HEAT-HAZE	1
HEATH	1
HEATS	1
HEAVE	1
HEAVEN	4
HEAVEN-HOLDING	1
HEAVIER-THAN...	1
HEAVIEST	2
HEAVILY	2

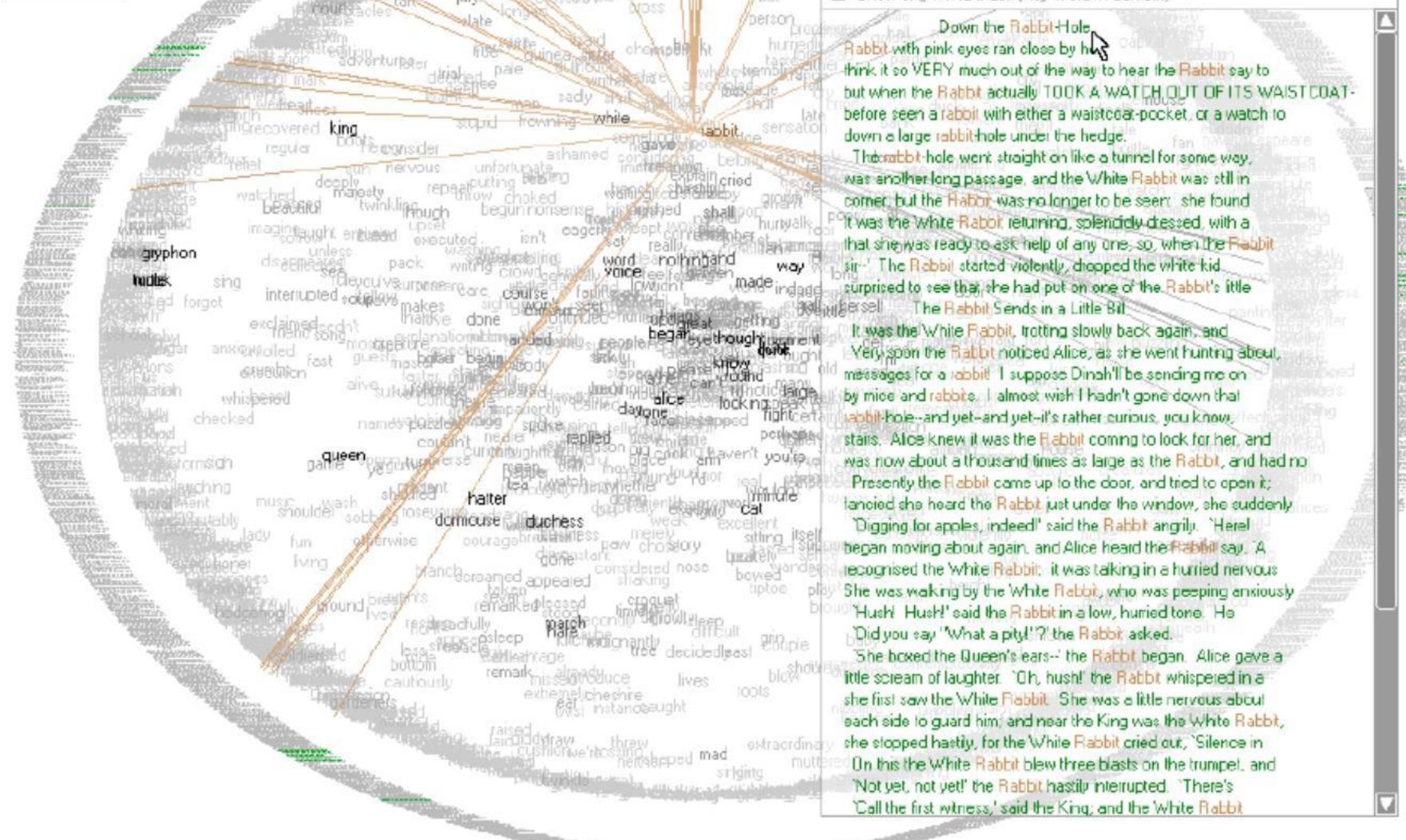
The right pane displays the following context snippets for 'HEART':

Context...	Word	...Context	Reference
That my own	heart	drifts and cries, having no...	Deep Analysis
By the shout of the	heart	continually at work	And the wave
Nothing to adapt the skill of the	heart	to, skill	And the wave
The tread, the beat of it, it is my own	heart	,	Träumerei
Because I follow it to my own	heart		Many famous
	My	is ticking like the sun:	I am washed i
	The vague	heart	sharpened to a candid co...
	Contract my	heart	by looking out of date.
	Having no	heart	to put aside the theft
	And the boy puking his	heart	out in the Gents
	A harbour for the	heart	against distress.
	These I would choose my	heart	to lead
	Time in his little cinema of the	heart	
	This petrified	heart	
How should they sweep the girl clean...	heart	has taken,	A Stone Churc
	Hands that the	heart	I see a girl dra
	For the	heart	Heaviest of flo
	With the unguessed-at	heart	Dawn
	If hands could free you,	heart	One man walk
	That overflows the	heart	If hands could
		,	Pour away tha

At the bottom, there are buttons for Words (7318), Tokens (37070), At word (2990), Deleted lines (1 [24]), Word sort (Asc alpha [string]), Context sort (Asc occurrence order), and a status bar.

Down the Rabbit-Hole

Down the Rabbit-Hole

[Hide text!](#)[Show concordance](#)[Show thesaurus lookup](#)[Show story line](#)

if love be rough with you , be rough with love .

if love be blind , love cannot hit the mark .

if love be blind , it best agrees with night .

if love be

blind ,

rough with you , be rough with love .

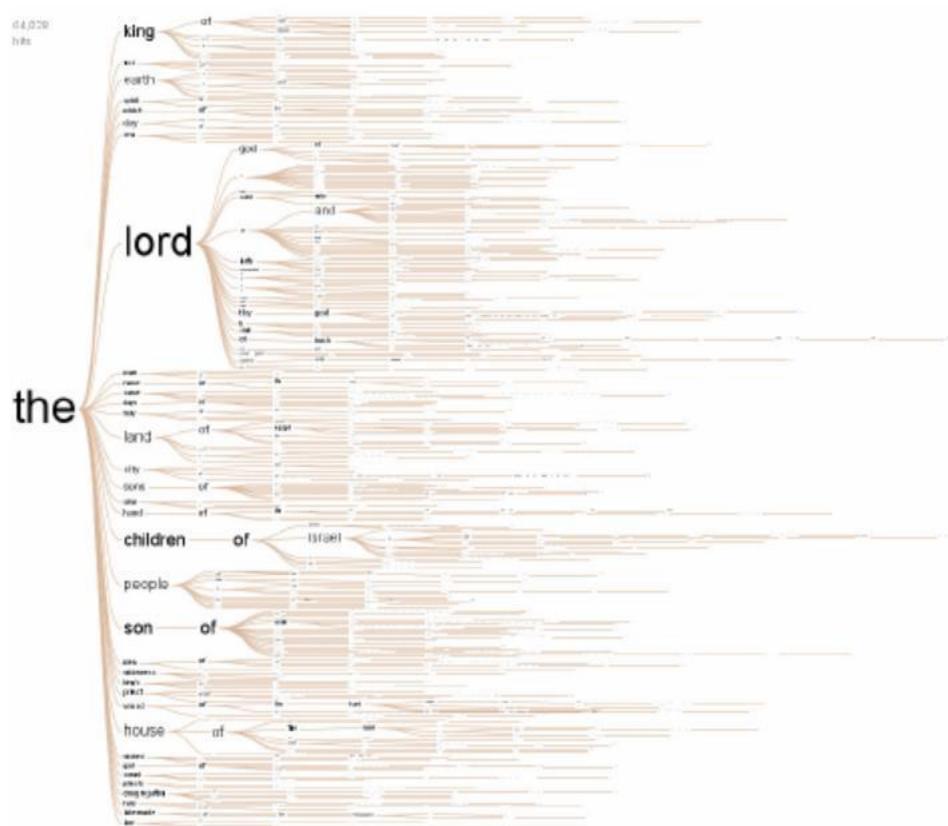
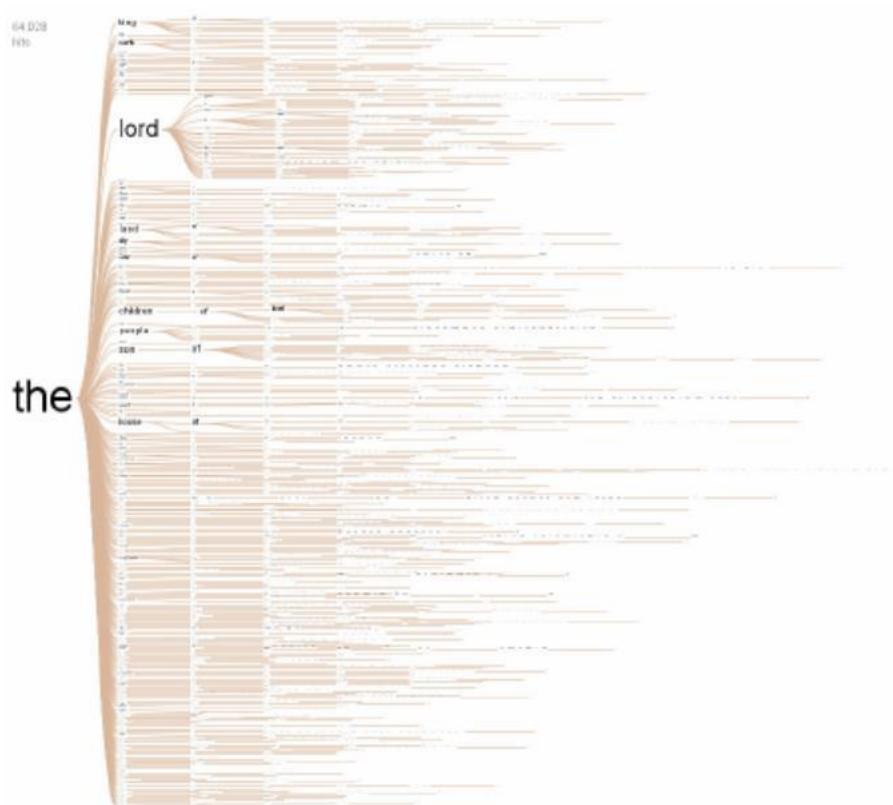
love cannot hit the mark .

it best agrees with night .

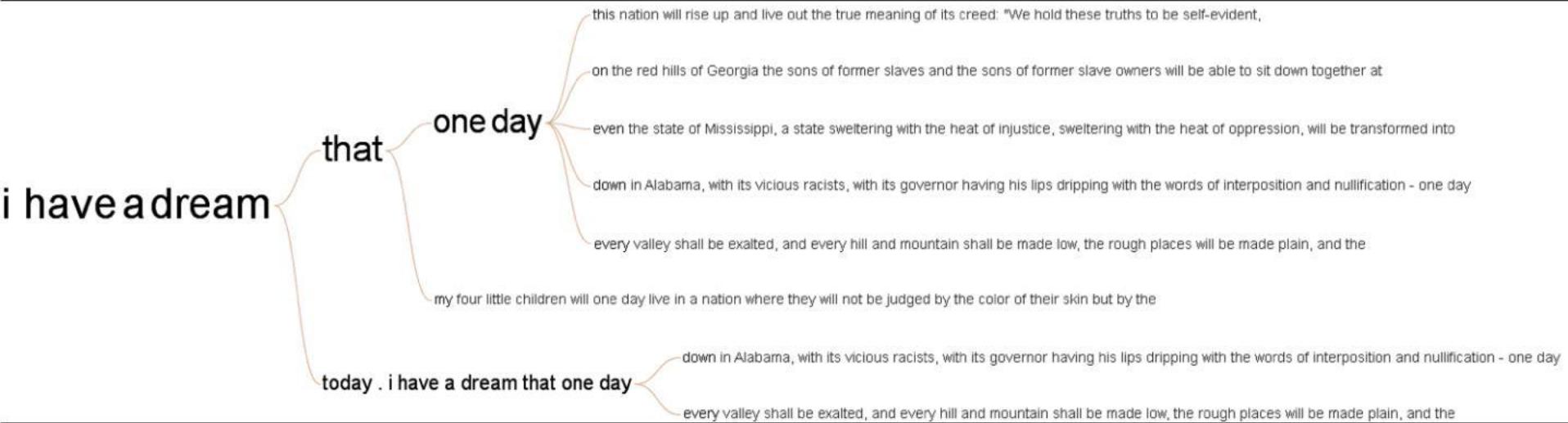
Word Tree [Wattenberg et al.]

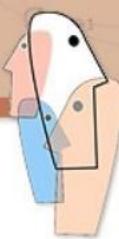


Filter Infrequent Runs



Recurrent Themes in Speeches





Visualizations : Word tree / Alberto Gonzales

Creator: Martin Wattenberg

Tags:

Search

[Back](#)

[Forward](#)

Start

End

Occurrence Order

Clicks Will Zoom

118

hits

i don't

want

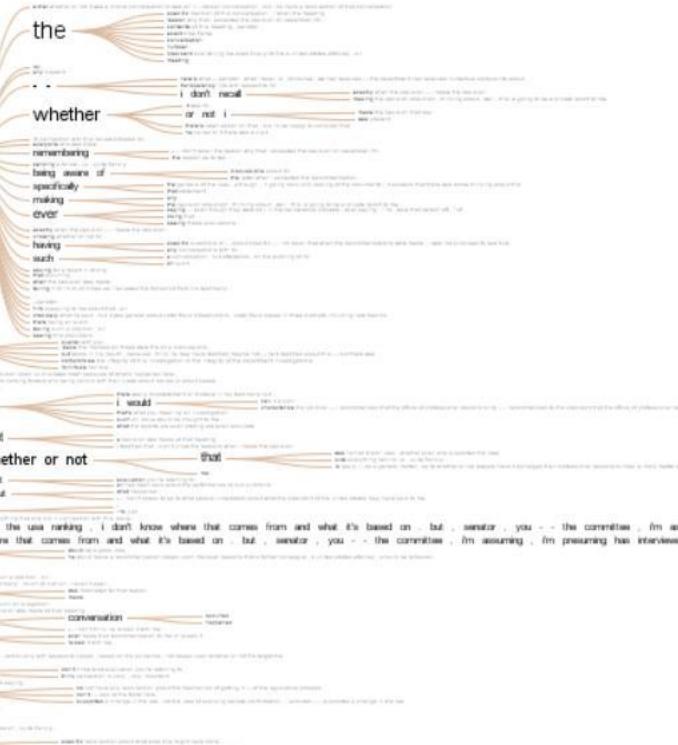
to

know

believe

think

have



Data file: Word in testimony from Gonzales, 4/19/2007

Data source: CQ Transcript Wire via the Washington Post

This data set

has not yet been rated



Comments (4)

currently showing



This visualization has 4 positive and 0 negative

Glimpses of Structure...

Concordances show local, repeated structure
But what about other types of patterns?

Lexical: <A> at

Syntactic: <Noun> <Verb> <Object>

Phrase Nets [van Ham et al.]

Look for specific **linking patterns** in the text:

'A and B', 'A at B', 'A of B', etc

Could be output of regexp or parser.

Visualize patterns in a node-link view

Occurrences -> Node size

Pattern position -> Edge direction

Select a phrase

Showing 73 of 1719 terms

- word1 and word2
- word1 's word2
- word1 of the word2
- word1 the word2
- word1 a word2
- word1 at word2
- word1 is word2
- word1 [space] word2

or enter your own

* and *

Submit

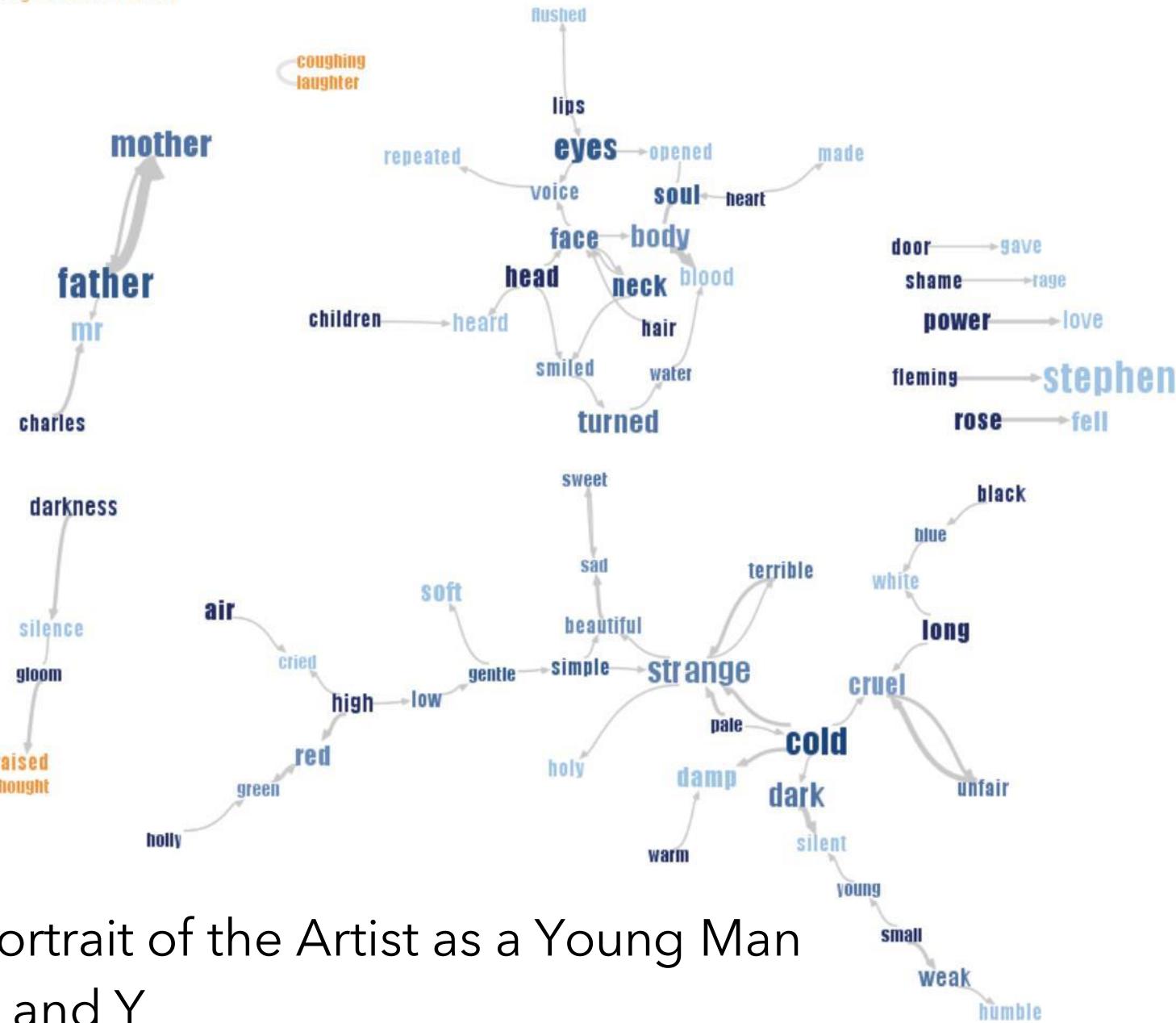
Filters

Show top: 100

Hide common words

Zoom

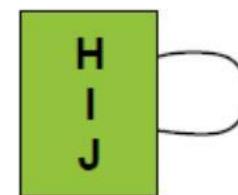
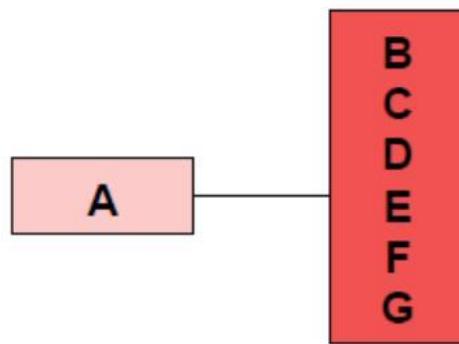
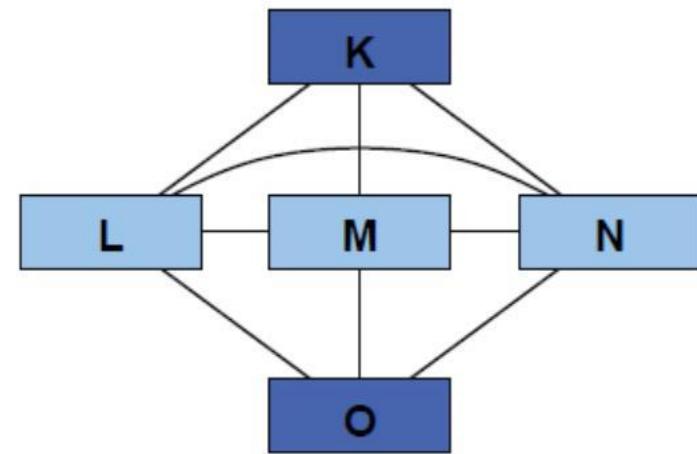
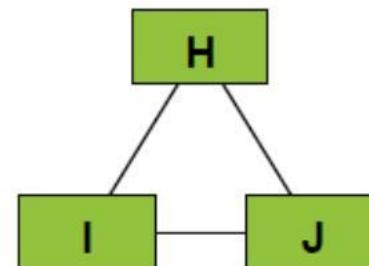
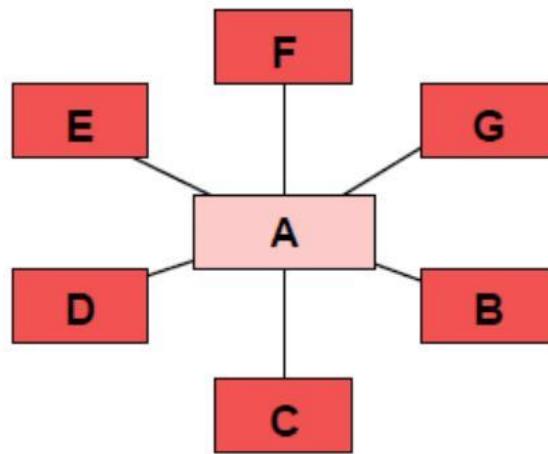
In Out Reset



Portrait of the Artist as a Young Man

X and Y

Node Grouping



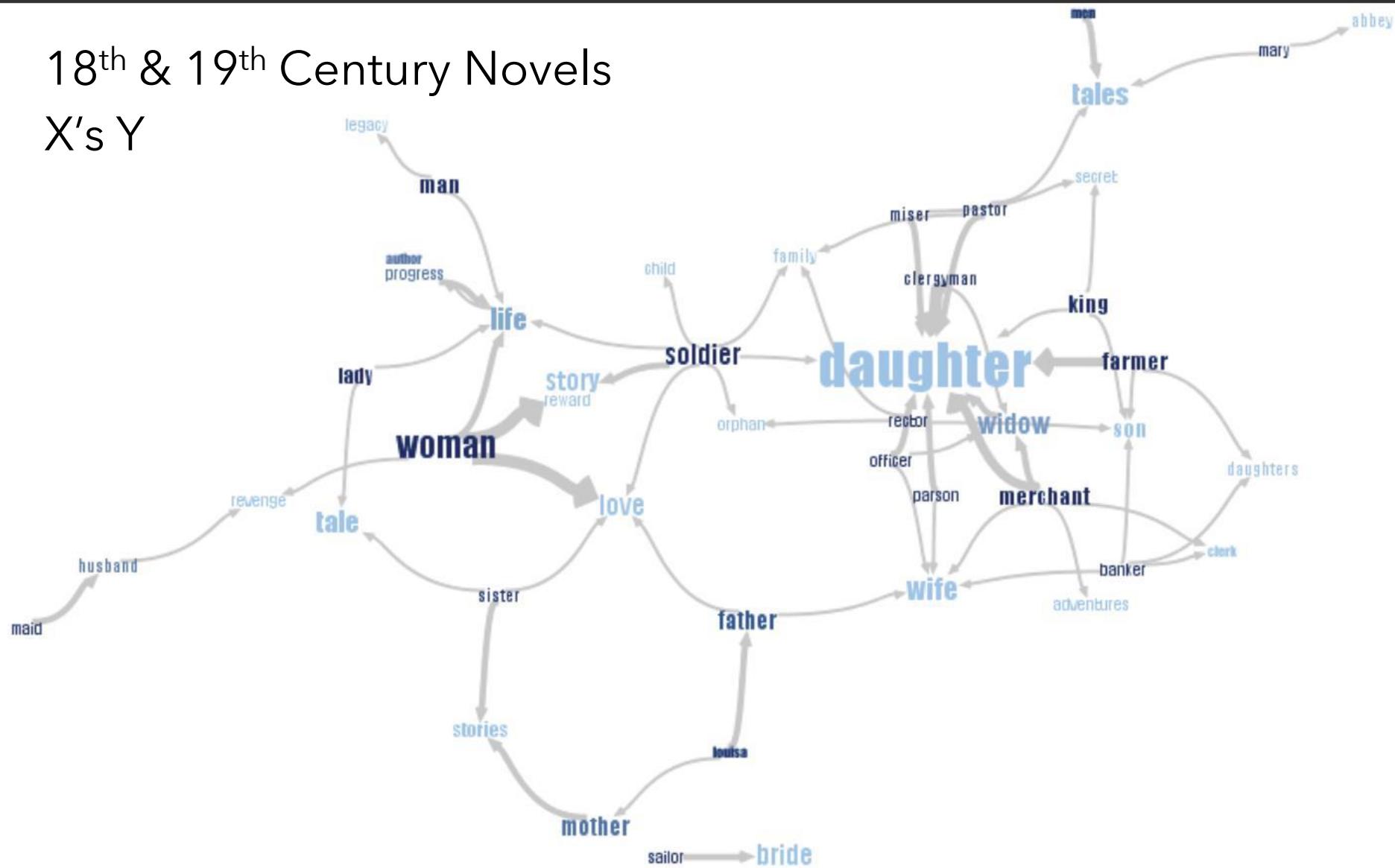
(a)

(b)

(c)

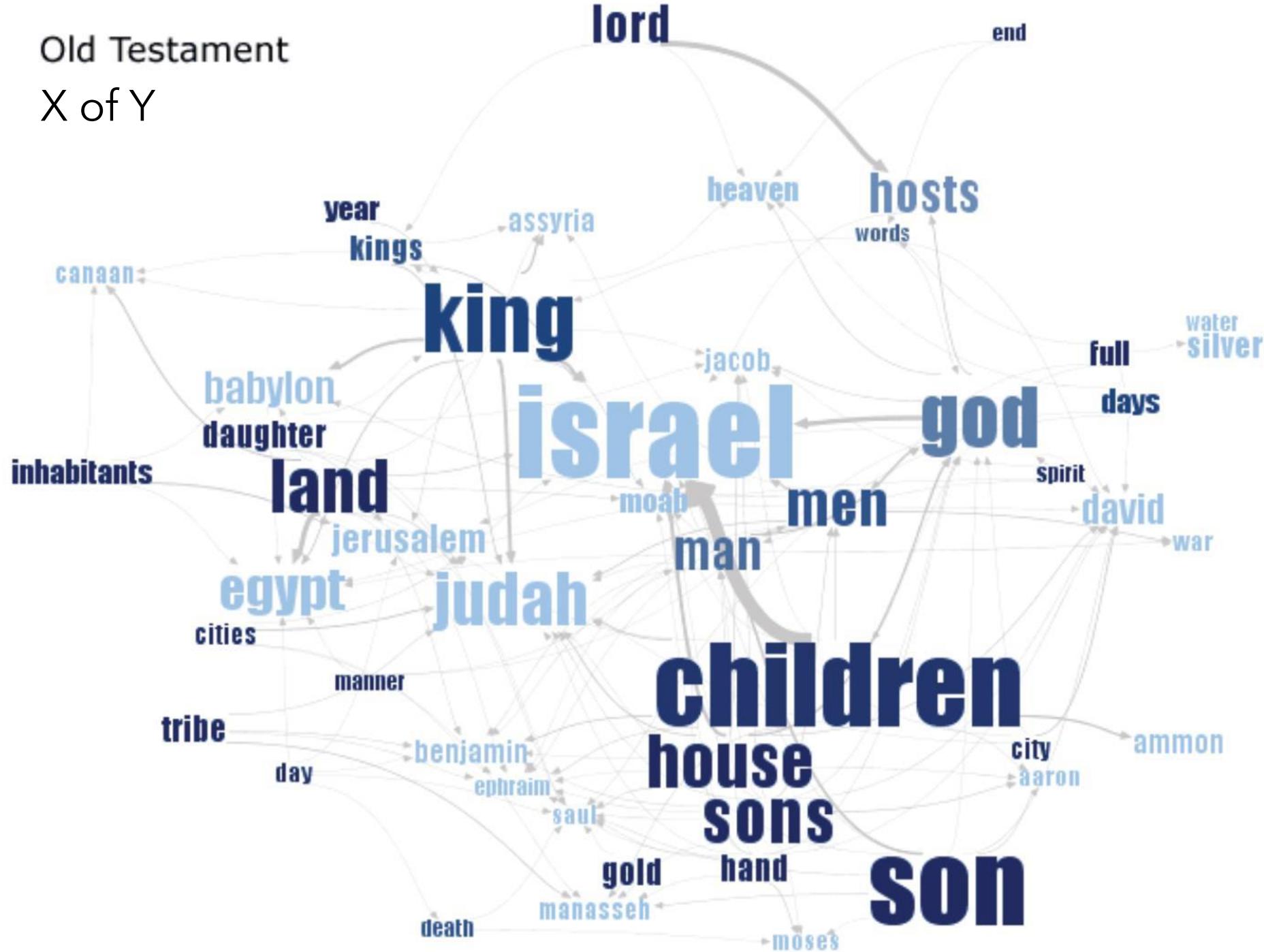
18th & 19th Century Novels

X's Y



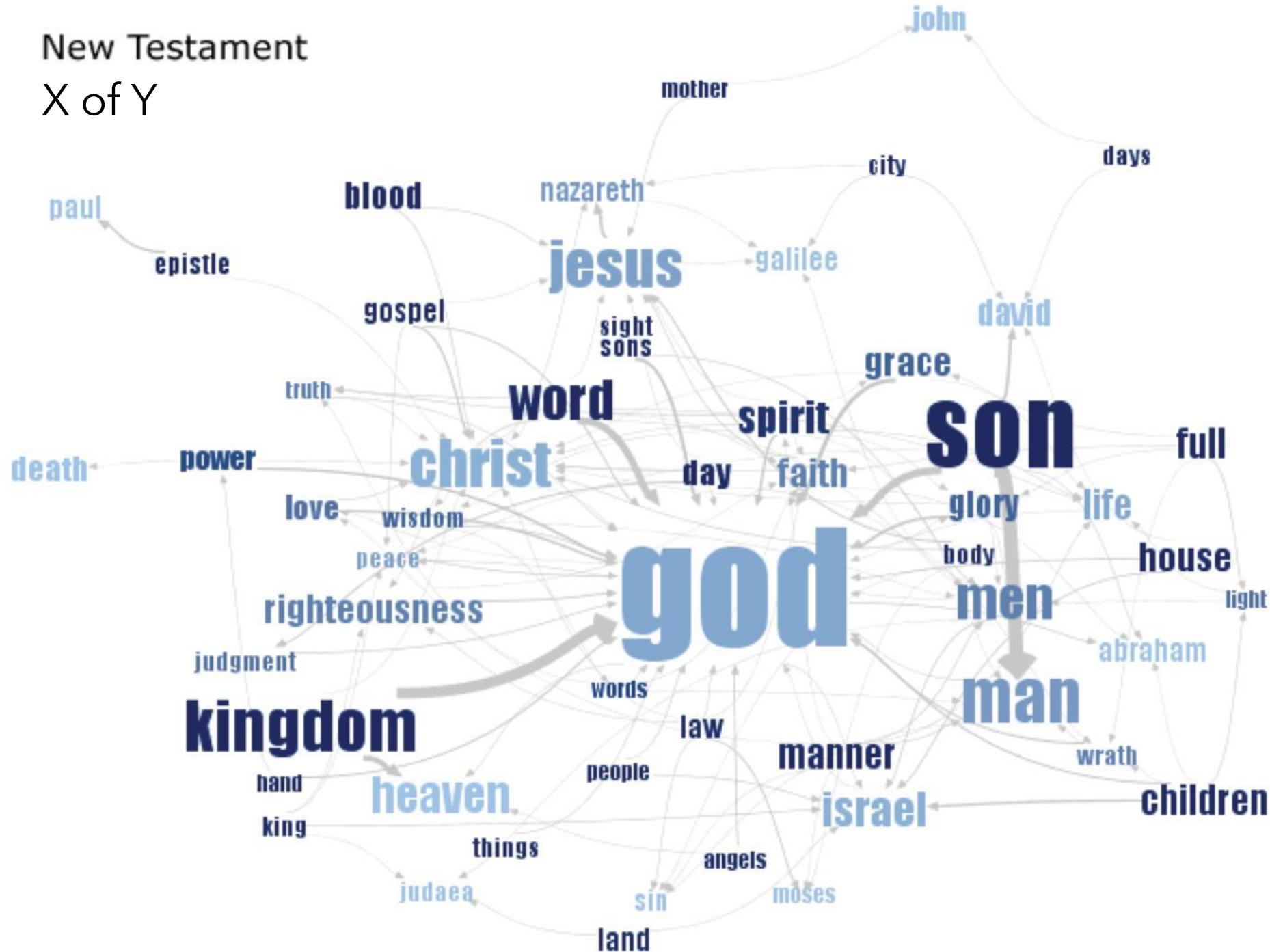
Old Testament

X of Y



New Testament

X of Y



Document Content

Understand Your Analysis Task

Visually: Word position, browsing, brush & link

Semantically: Word sequence, hierarchy, clustering

Both: Spatial layout reflects semantic relationships

The Role of Interaction

Language model supports visual analysis cycles

Allow modifications to the model: custom patterns
for expressing contextual or domain knowledge

Conversations

Visualizing Conversation

Many dimensions to consider:

Who (senders, receivers)

What (the content of communication)

When (temporal patterns)

Interesting cross-products:

What x When -> Topic “Zeitgeist”

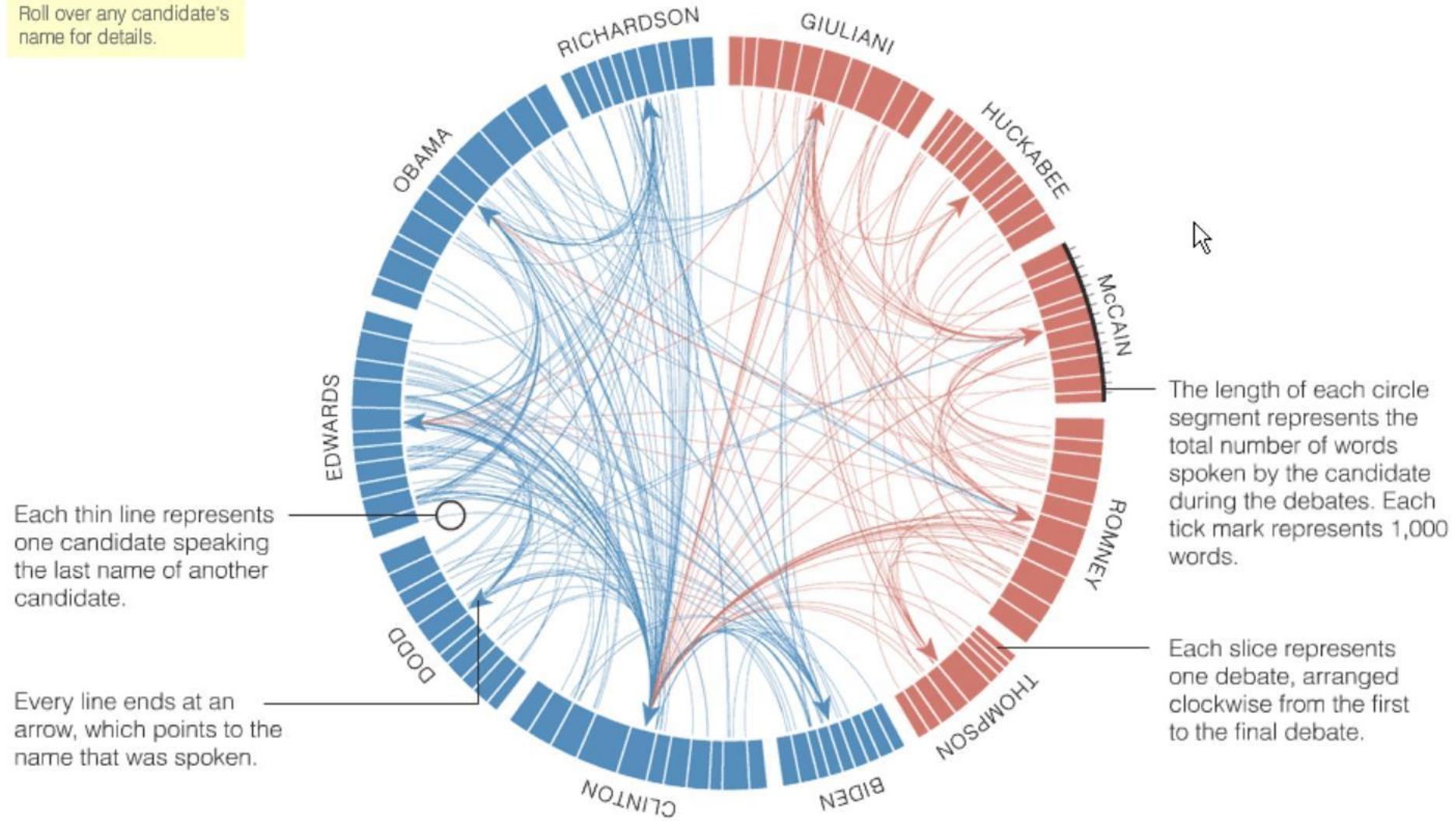
Who x Who -> Social network

Who x Who x What x When -> Information flow

Naming Names

Names used by major presidential candidates in the series of Democratic and Republican debates leading up to the Iowa caucuses.

Roll over any candidate's name for details.



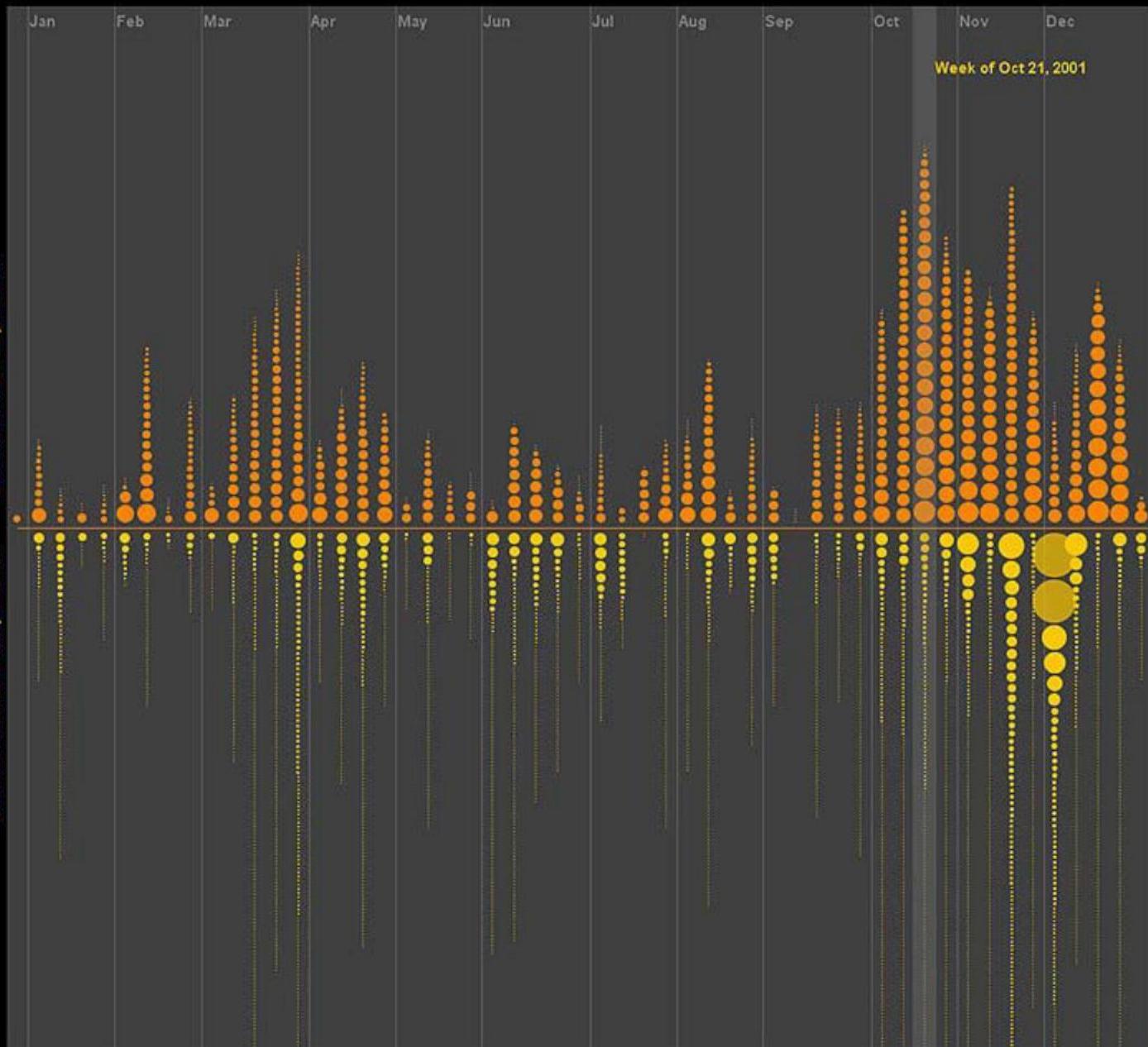
Usenet Visualization [Viegas & Smith]

Show correspondence patterns in text forums

Initiate vs. reply; size and duration of discussion



author: jillyb@mail.com

[back to newsgroups](#)

Week of Oct 21, 2001

subject	# of posts
Wednesday Spooker ASF	21
WET #3 Anyone for breakfast	20
Sunny Side Up ASF)	18
Saturday Ensemble and WET	18
Oh no! Watch out! ASF	18
Thursday Combo-Post WET #	16
The Yellow Rose Inn... A gift to	16
WET #1 JBP The First Time	15
We Love the Earth ASF	15
Monday Spooker "The Sight"	15
C'mon!!!	14
Theberge "Le Vent Se Lève"	14
Holiday Tog #3)	13
Spooker du Jour)	13
Beginning ASF Short and	13
Second Try A Katie for Suzy	12
Come On a Safari With Me	11
Tuesday Spooker ASF	11
Curses, Foiled Again....ASF	10
Halloween Togs Take Two)	9
Beauty of the Fury Jim Warren	9
I thought I saw _____? ASF	7
Wednesday Evening at the Con	4
Second Try A Katie for Suzy	2
Frank Was A Monster ASF	1

subject	# of posts
Sunday Twofer ASF)	9
Chopsticks/A Jilly fake	8
Oh no! Trouble in Discworld!	7
WET - your thirst! ASF	6
A pretty for you...Reposted fro	5
Saturday Spooker ASF	5
Sample Previous install Upgr	4
Tennessee weather tonite	4
WET - Well I am not smiling!	4
Somethin' mushy <asf>	3
Getting seasonal with workin...	3
A Haunted House)	3
do you wonder what debt's be...	3
Question: Ethics of posters in	3
For Jerry	3
Olu's Tribe - slightly rated	3
WET - Glass Bottles	3
Peace Train<ASF>	2
Arrival at Stewart Island II	2
WET195 Wrap-up	2
Cat O'Lantern	2
I Put a Spell on You (Happy H...	2
Goodbye to Summer - A Timel...	2
Two Pumpkins In A Strange B...	2
Still Heading South !	2
WET- Frank Sinatra - The Man...	2
WET Autumn	2
Purple Martin ASF	2
Opposites Attract...	2
Time	2



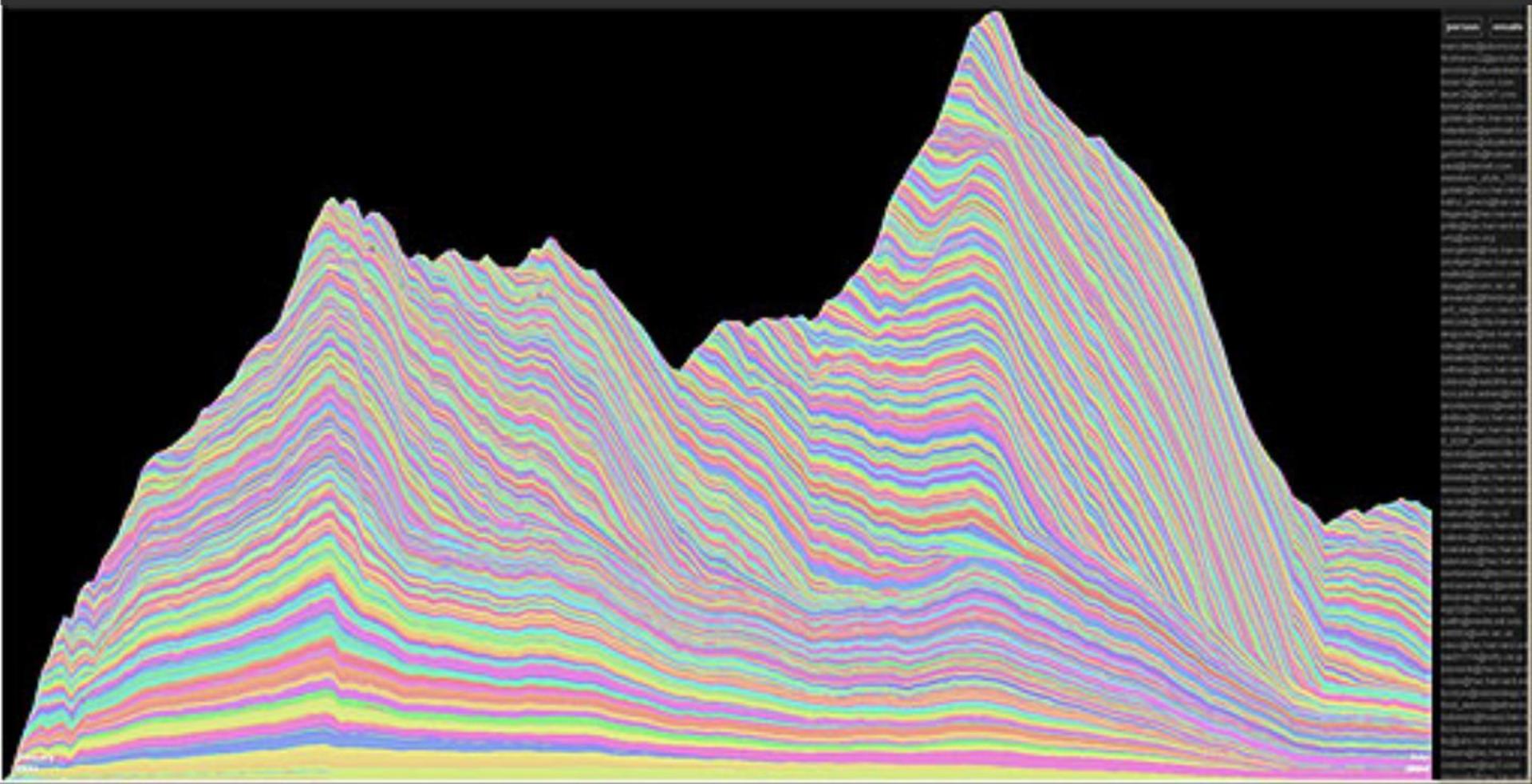
Threads Initiated by Author

Threads not Initiated by Author

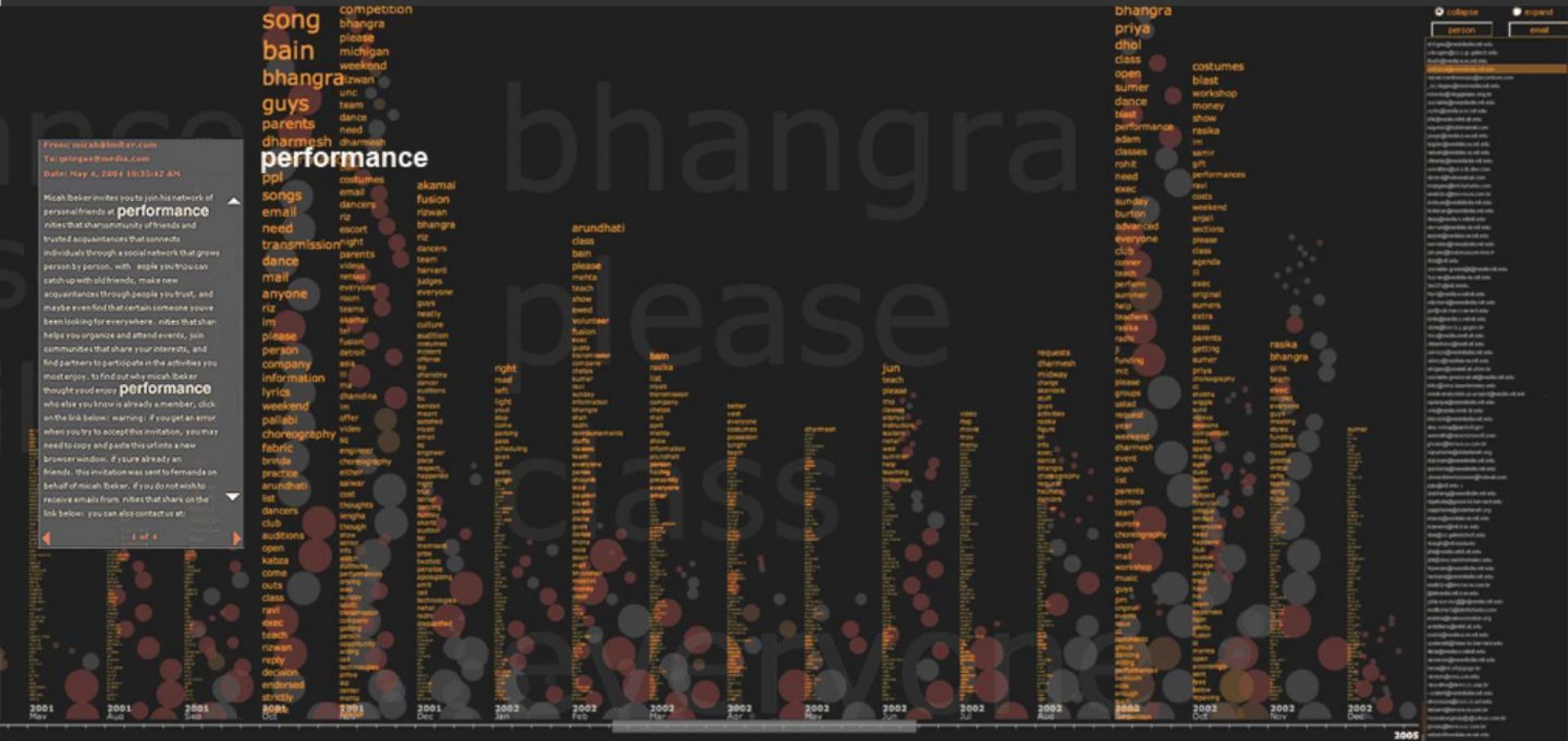
Subject	# of posts
Antihomophobia	245
Evolution Science	86
OLD TESTAMENT	37
Plato's atheist god	30
Scientific Name	24
Why evolution won't go away	23
Creationism	15
Creationists	14
TOMMY COOK	90
Darwin and God	7
Darwin's theory	7
The Intelligent Design movement	7
President Bush vs. God	6
God vs evolution	6
A Darwinized God	6
Answers in Genesis	5
Original Sin/Bad Fallacy	5
Science vs. Intelligent Design	5
SCIENTIFIC PROOF	5
CHRISTIAN SCIENCE	4
Antihomophobia	4
All About Evolution	4
CAR TRUTH	4
EVOLUTION VS. FAITH	4
Vaccines vs. Disease	4

Email Mountain [Viegas]

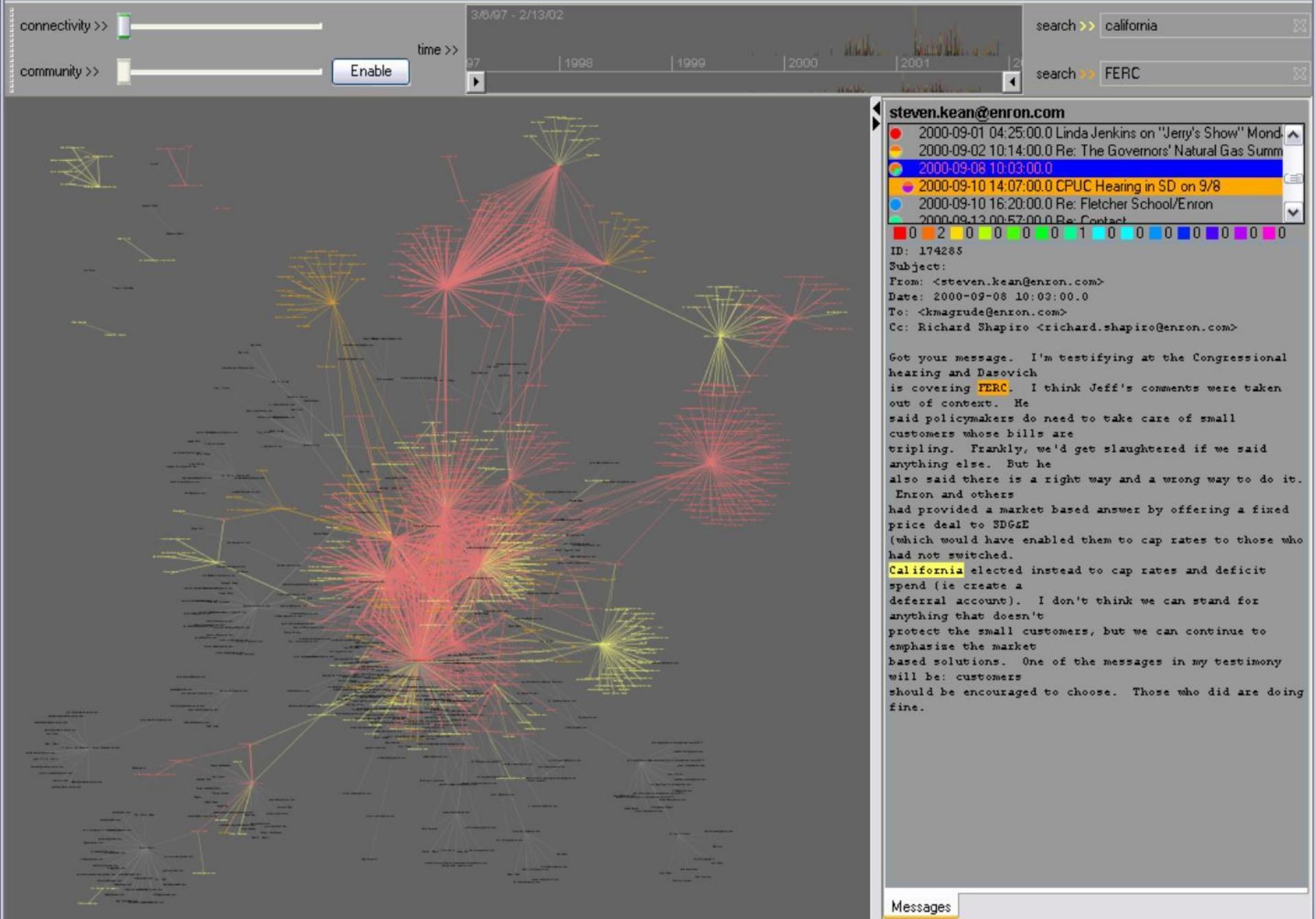
Conversation by person over time (who x when).

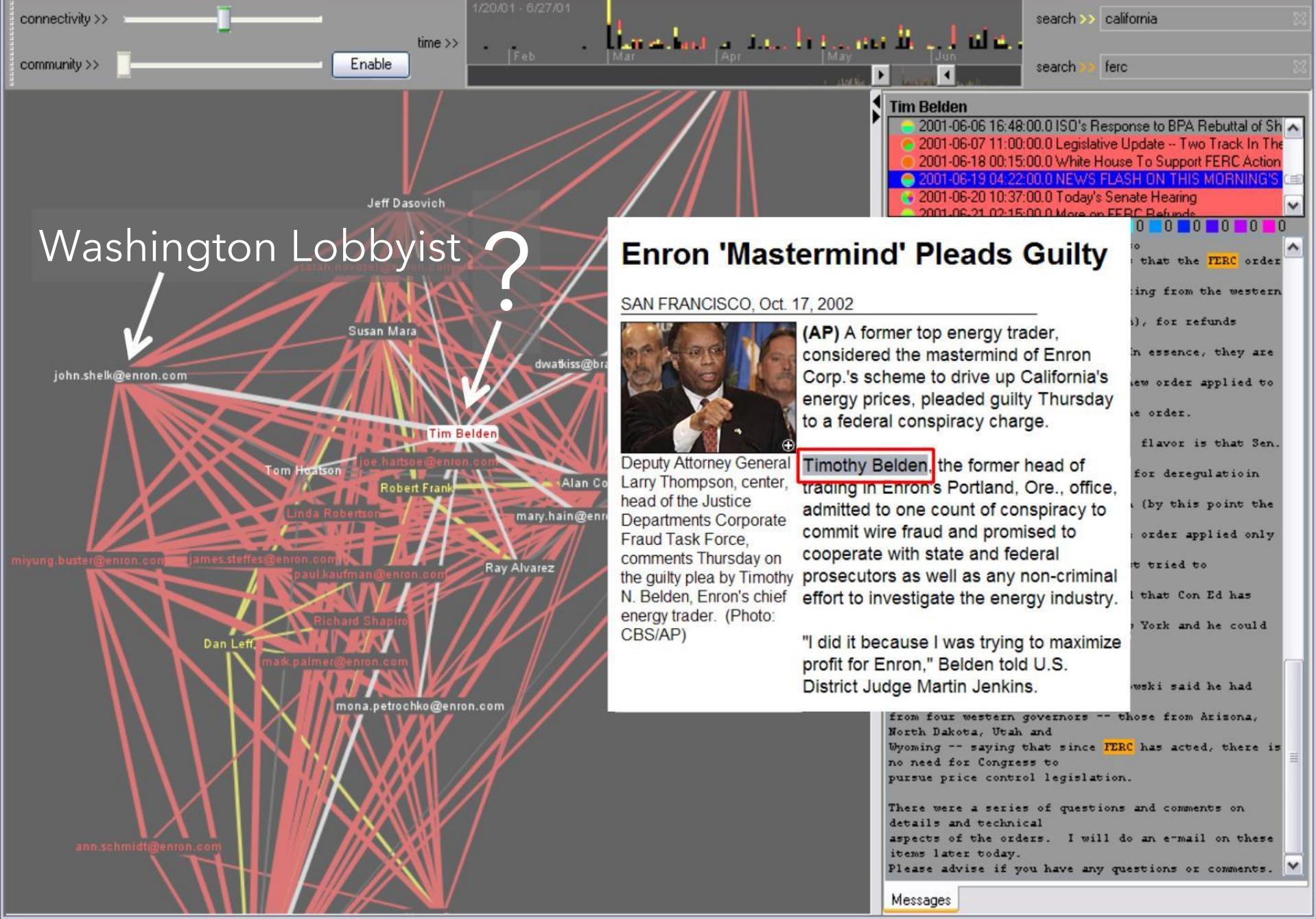


The mail [Viegas]



One person over time, TF.IDF weighted terms





Document Collections

Named Entity Recognition

Label named entities in text:

John Smith -> PERSON

Soviet Union -> COUNTRY

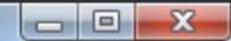
353 Serra St -> ADDRESS

(555) 721-4312 -> PHONE NUMBER

Entity relations: how do the entities relate?

Simple approach: do they co-occur in a small window of text?

List View



Edit View Bookmarks Lists Options

person

Add all

Clear

ABC



Show all connections

place

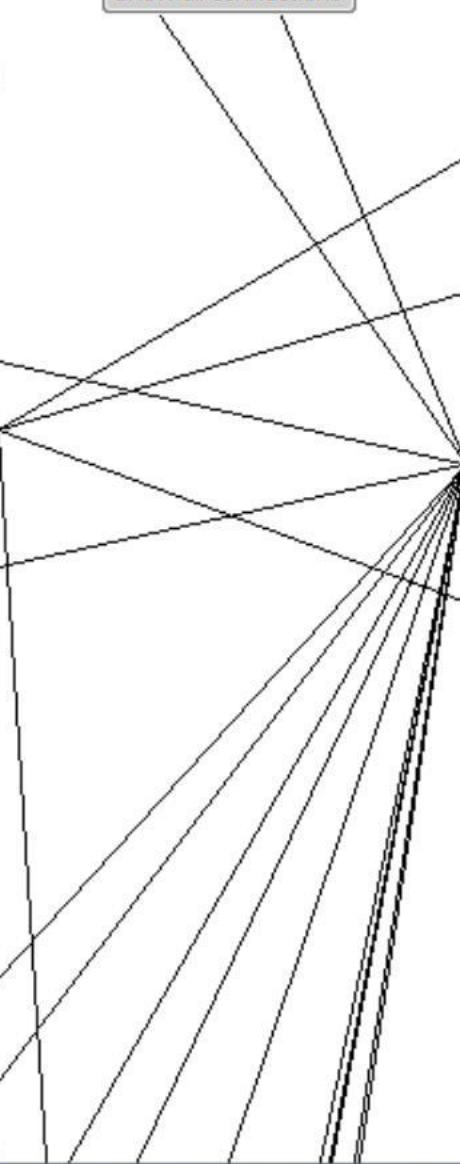
Add all

Clear

ABC

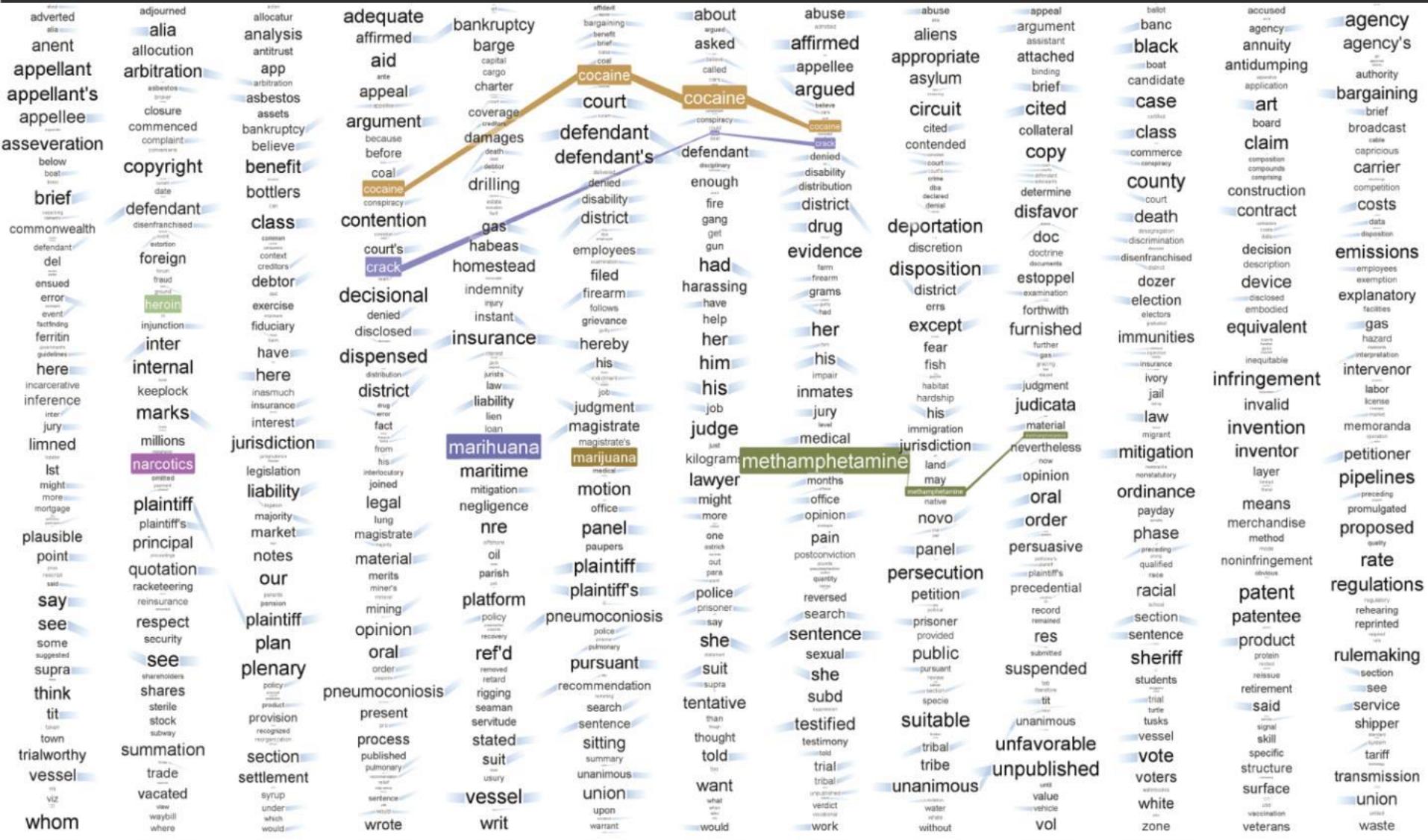


Bugarov
Carlos
Carlos Araneda
Carlos Morales
Castro
Cesar Arze
Charles Wilson
Dan West
Daniel Harris
David Loiseau
Dean Simpson
Dr. Baker
Dustin Marshall
Edgar Spencer
Edward Thompson
Escalante
F. Baker
Felix Baker
Ford
Forrest Wells
Fr. Augustin Dominique
Fred Fisher
George Garcia
Grigory Sizov
Hamid Qatada
Hector Lopez
Herman Fox
Howard Clark
Igor Kolokov
Imad Dahdah
J. T.
Tamat Sveed

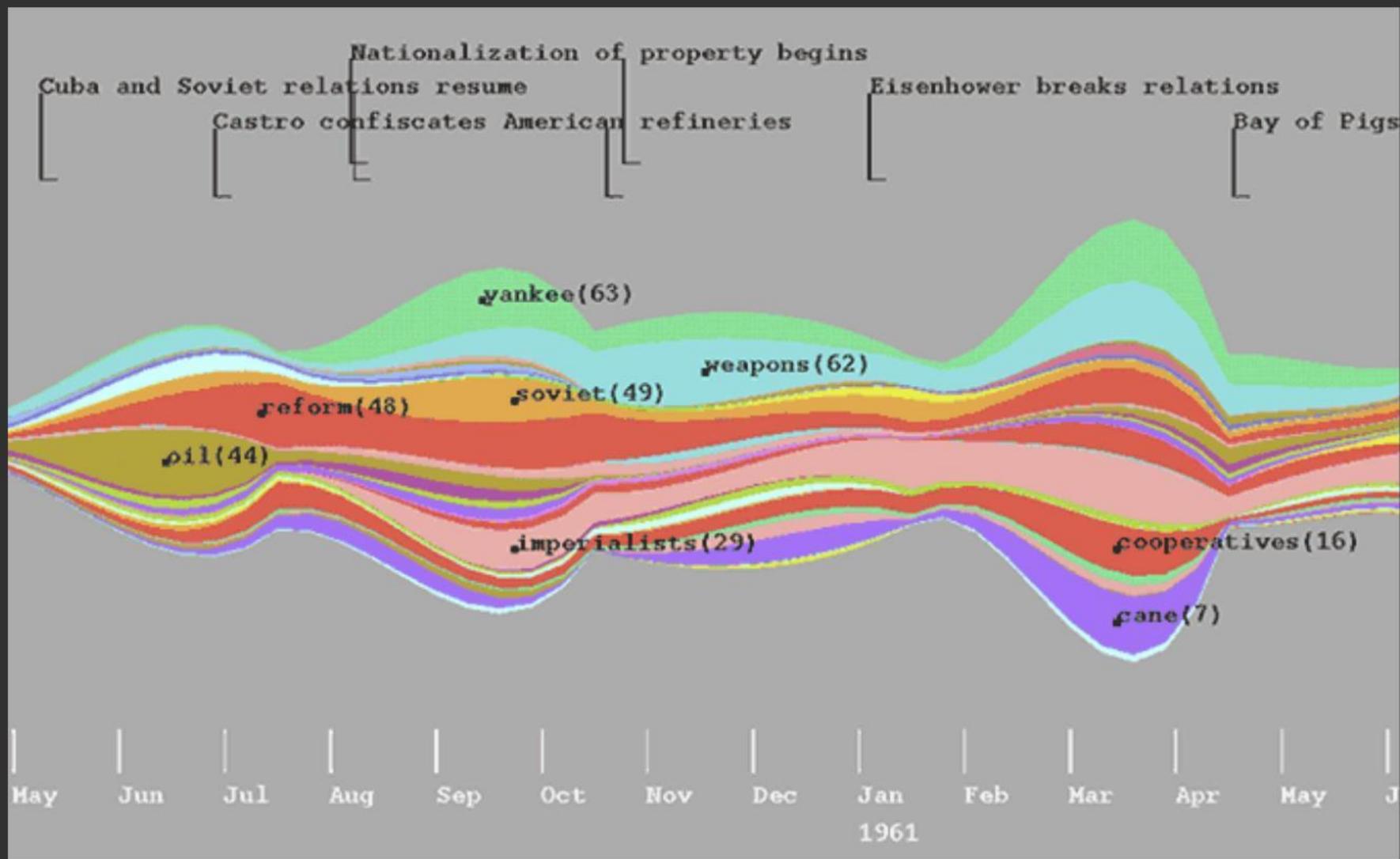


USA
Cuba
Pakistan
Canada
Columbia
Jamaica
Afghanistan
Havana
Detroit
Mexico
Michigan
Montego Bay
Texas
Chitral
Morocco
Peshawar
Russia
Casablanca
Chicago
Illinois
New Jersey
UK
Dominican Republic
Florida
France
London
Moscow
Ontario
Paris
Windsor
Santo Domingo
Virginia

Parallel Tag Clouds [Collins et al.]



Theme River [Havre et al.]



Similarity & Clustering

Compute vector distance among docs

For TF.IDF, typically cosine distance

Similarity measure can be used to cluster

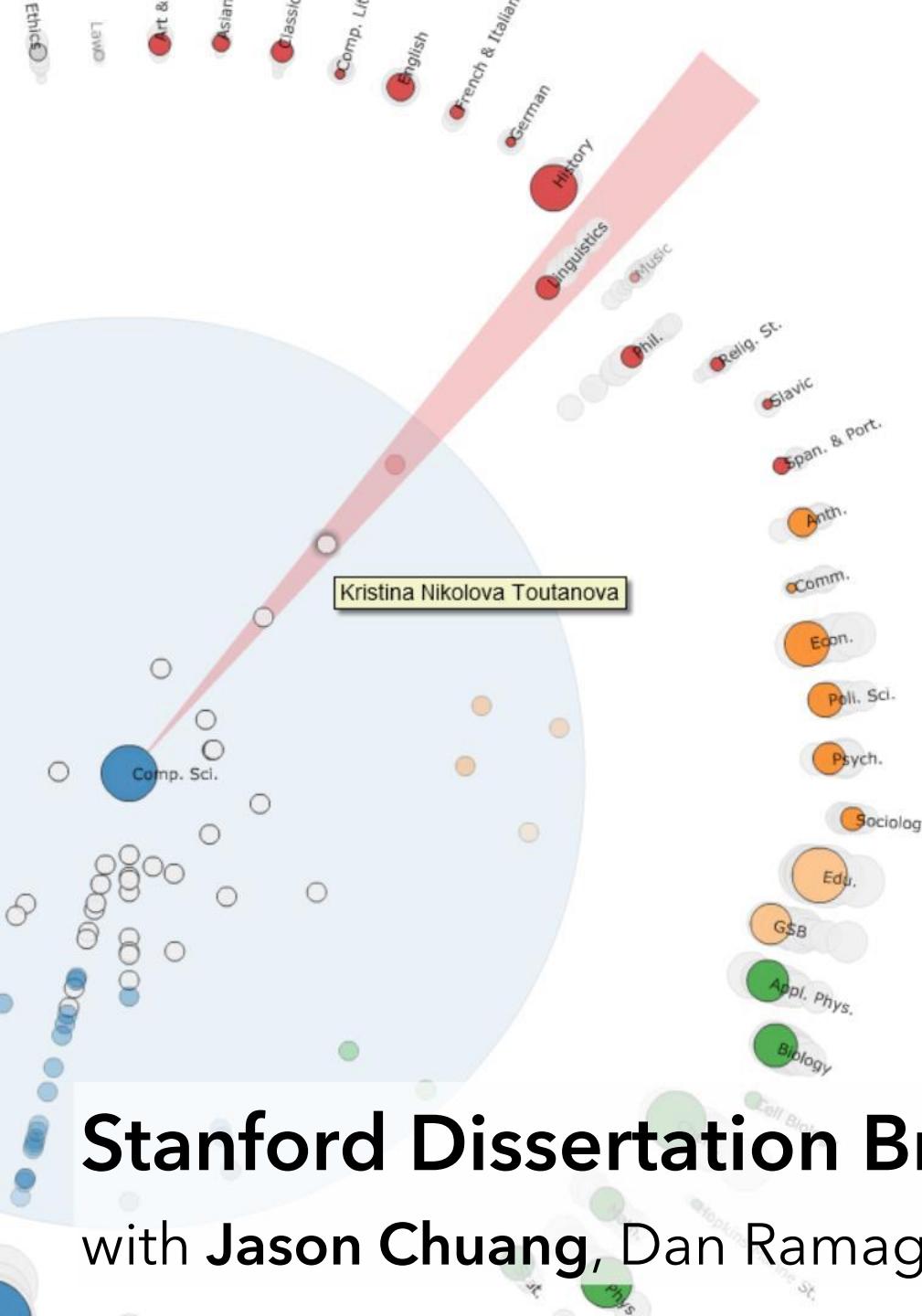
Topic modeling

Assume documents are a mixture of topics

Topics are (roughly) a set of co-occurring terms

Latent Semantic Analysis (LSA): reduce term matrix

Latent Dirichlet Allocation (LDA): statistical model



Stanford Dissertation Browser

with Jason Chuang, Dan Ramage & Christopher Manning

Effective statistical models for syntactic and semantic disambiguation

Student: Kristina Nikolova Toutanova

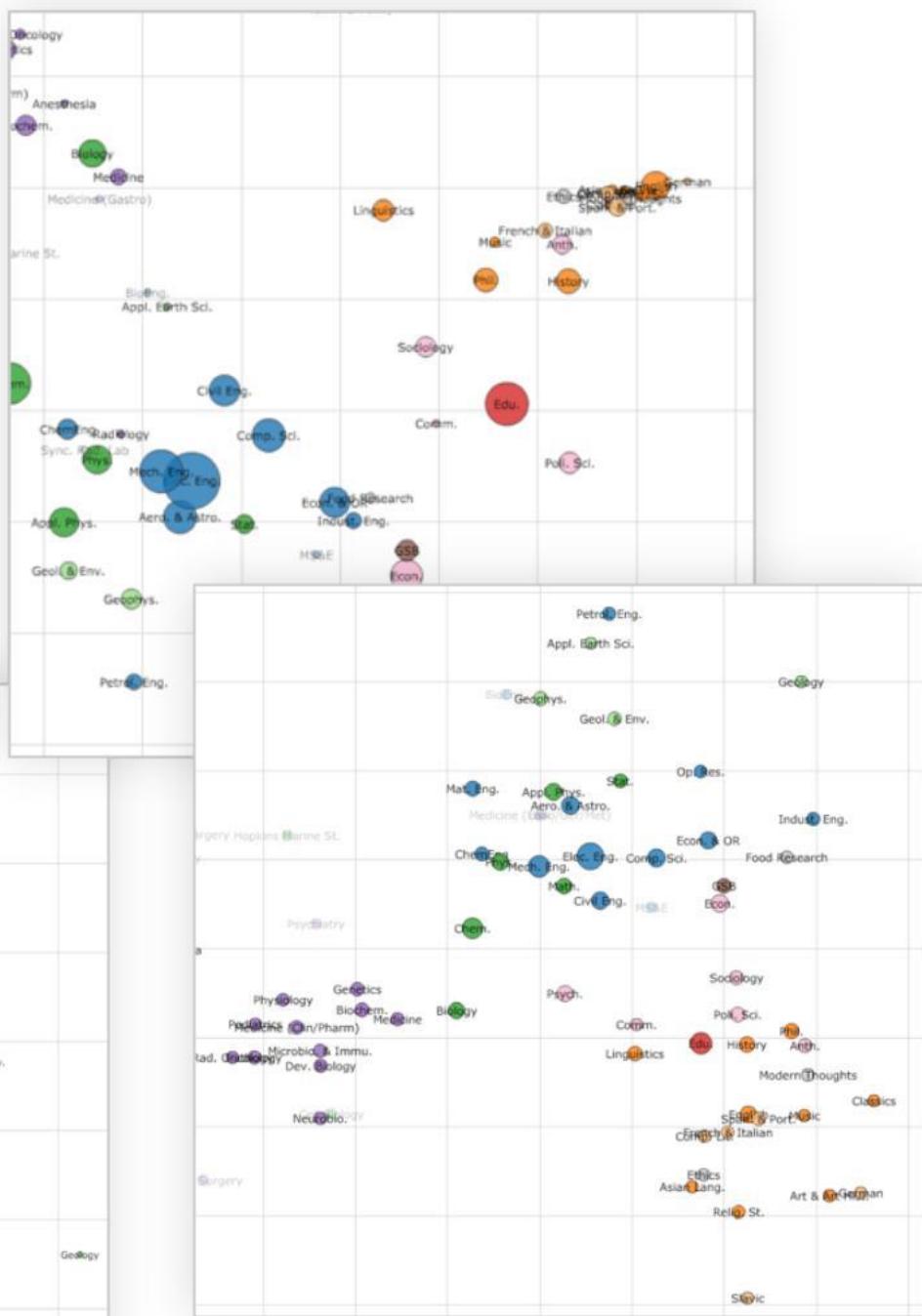
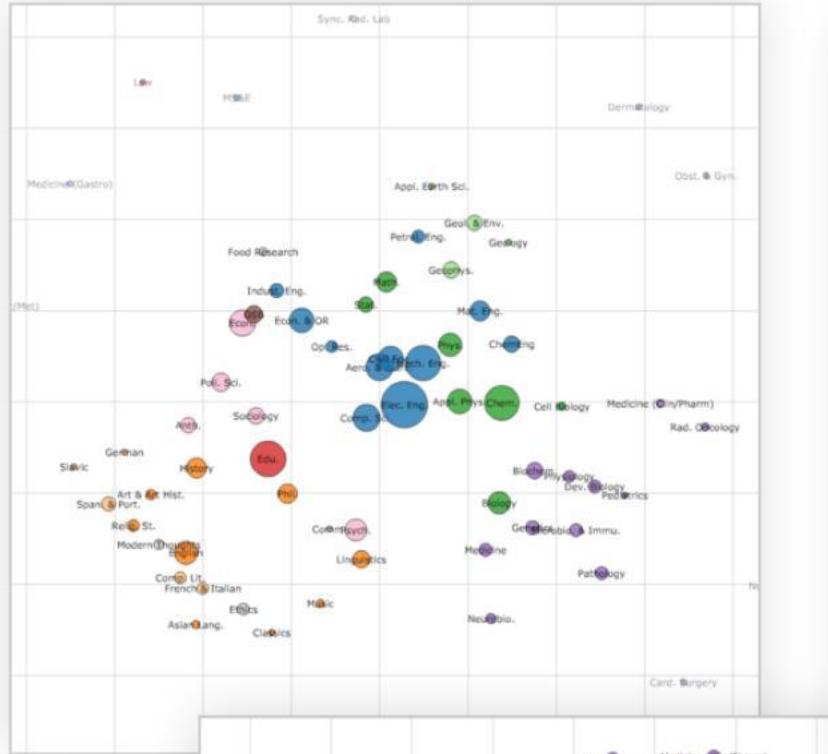
Advisor: Christopher D. Manning

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.



Thank you