

Text Visualization

Maneesh Agrawala

CS 448B: Visualization
Fall 2020

1

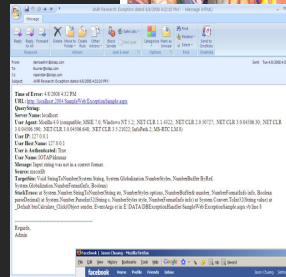
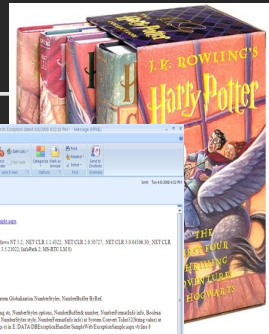
Text as data

Documents

- Articles, books and novels
- Computer programs
- E-mails, web pages, blogs
- Tags, comments

Collection of documents

- Messages (e-mail, blogs, tags, comments)
- Social networks (personal profiles)
- Academic collaborations (publications)



2

Announcements

3

Final project

Data analysis/explainer or conduct research

- **Data analysis:** Analyze dataset in depth & make a visual explainer
- **Research:** Pose problem, Implement creative solution

Deliverables

- **Data analysis/explainer:** Article with multiple interactive visualizations
- **Research:** Implementation of solution and web-based demo if possible
- **Short video (2 min)** demoing and explaining the project

Schedule

- Project proposal: **Thu 10/29**
- Design Review and Feedback: **Tue 11/17 & Thu 11/19**
- Final code and video: **Sat 11/21 11:59pm**

Grading

- Groups of **up to 3 people**, graded individually
- Clearly report responsibilities of each member

4

Class Schedule

Guest Lecture Th Nov 12

Jessica Hullman on
Visualizing Uncertainty



5

Design Feedback (Next Week)

Signup for a ~10 min slot

Will post signups on Piazza later this week

Plan to give a 5 min presentation (mostly demo) of work so far. We will give oral feedback.

6

Text Visualization

8

Why visualize text?

9

Why Visualize Text?

Understanding: get the “gist” of a document

Grouping: cluster for overview or classification

Compare: compare document collections, or inspect evolution of collection over time

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

10

Example: Health Care Reform

Background

Initiatives by President Clinton

Overhaul by President Obama

Text data

News articles

Speech transcriptions

Legal documents

What questions might you want to answer?

What visualizations might help?

11



16

Gulf of Evaluation

Many (most?) text visualizations do not represent text directly. They represent the output of a **language model** (word counts, word sequences, 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?

17

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

18

Topics

Text as Data

Visualizing Document Content

Visualizing Conversation

Document Collections

20

Text as Data

21

Words as nominal data?

High dimensional (10,000+)

More than equality tests

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

Words have meanings and relations

22

Text Processing Pipeline

Tokenization

Segment text into terms.

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

Numbers and symbols? *#cardinal, @Stanford, OMG!!!!!!!*

Entities? *Palo Alto, O'Connor, U.S.A.*

23

Text Processing Pipeline

Tokenization

Segment text into terms.

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

Numbers and symbols? *#cardinal, @Stanford, OMG!!!!!!!*

Entities? *Palo Alto, O'Connor, U.S.A.*

Stemming

Group together different forms of a word.

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

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

24

Text Processing Pipeline

Tokenization

Segment text into terms.

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

Numbers and symbols? *#cardinal, @Stanford, OMG!!!!!!!!*

Entities? *Palo Alto, O'Connor, U.S.A.*

Stemming

Group together different forms of a word.

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

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

Ordered list of terms

25

The Bag of Words Model

Ignore ordering relationships within the text

A document \approx vector of term weights

Each term corresponds to a dimension (10,000+)

Each value represents the relevance

- For example, simple term counts

Aggregate into a document \times term matrix

Document vector space model

26

Document x Term matrix

Each document is a vector of term weights
Simplest weighting is to just count occurrences

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	157	73	0	0	0	0
Brutus	4	157	0	1	0	0
Caesar	232	227	0	2	1	1
Calpurnia	0	10	0	0	0	0
Cleopatra	57	0	0	0	0	0
mercy	2	0	3	5	5	1
worser	2	0	1	1	1	0

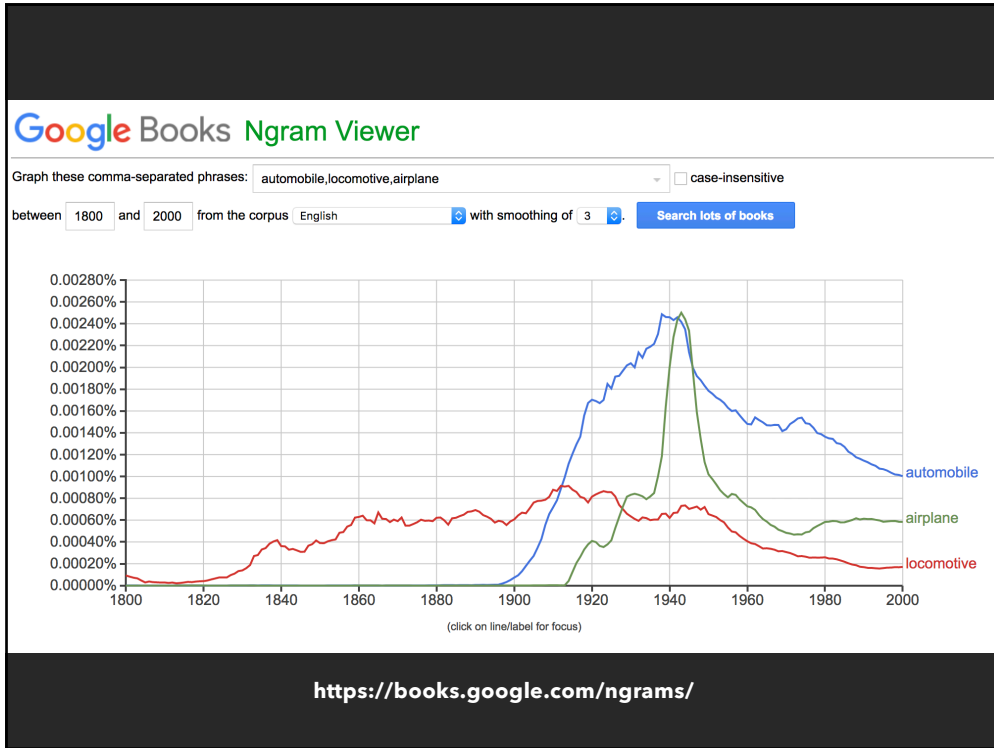
27

WordCount (Harris 2004)

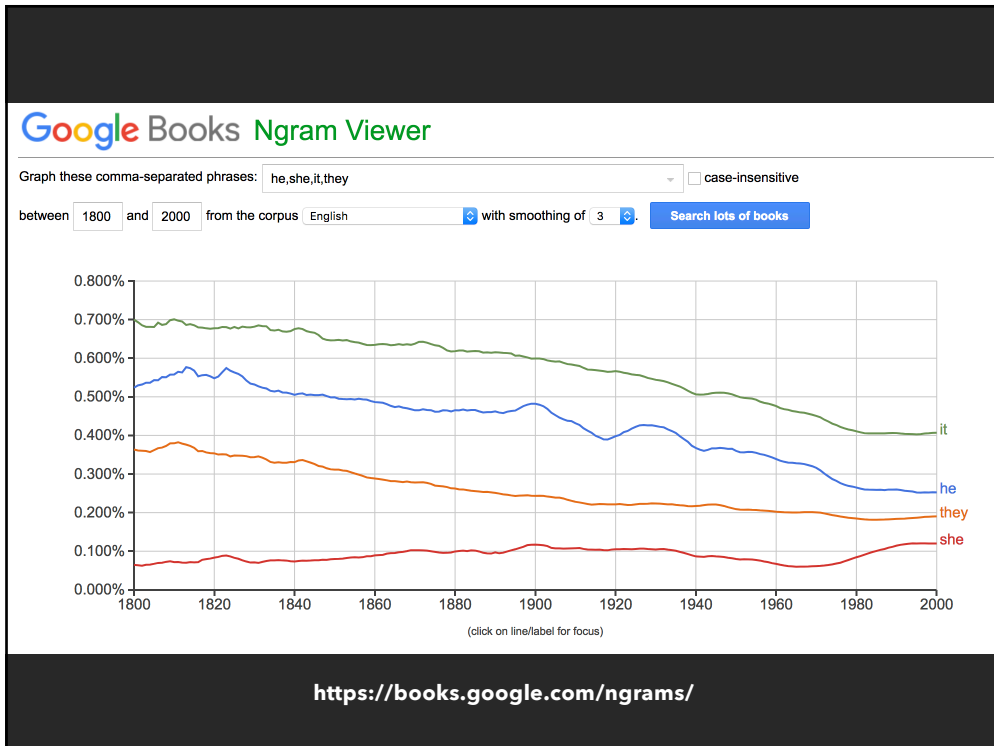
The screenshot shows the WordCount website interface. At the top right, there is a button labeled "WORDCOUNT". Below it, there are navigation links for "PREVIOUS WORD" and "NEXT WORD". The main display area shows a horizontal bar chart representing the frequency distribution of words. The word "the" is highlighted in large black text, with a red "1" below it indicating its rank. Other words are shown in smaller, lighter text: "of" (rank 2), "and" (rank 3), "to" (rank 4), "ain" (rank 5), "that" (rank 6), "is" (rank 7), "was" (rank 8), "for" (rank 9), "on" (rank 10), "you" (rank 11), "help" (rank 12), "with" (rank 13), "by" (rank 14), "the" (rank 15), "at" (rank 16), "the" (rank 17), "in" (rank 18), "the" (rank 19), "the" (rank 20), "the" (rank 21), "the" (rank 22), "the" (rank 23), "the" (rank 24), "the" (rank 25). Below the chart, there is a section for "CURRENT WORD" and a search bar. The search bar contains the text "FIND WORD:" followed by a dropdown menu, "BY RANK:" followed by a dropdown menu, and "REQUESTED WORD: THE". Below the search bar, it says "RANK: 1". At the bottom right, it says "86800 WORDS IN ARCHIVE" and "ABOUT WORDCOUNT".

<http://wordcount.org>

28



29



30

Given a text, what are the best descriptive words?

33

Keyword Weighting

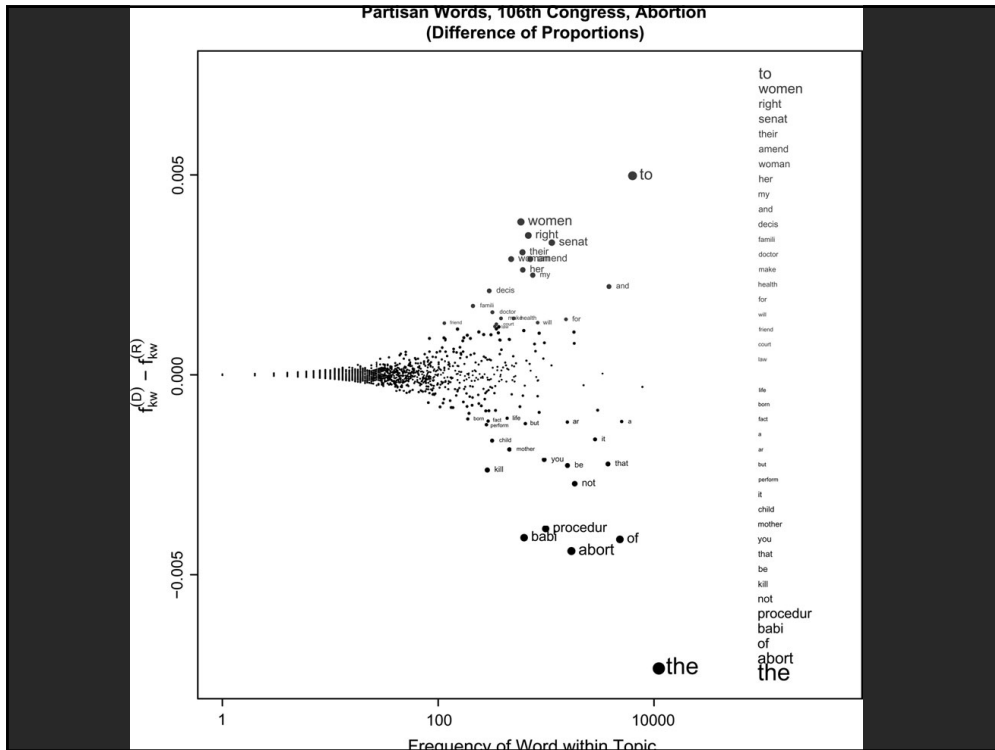
Term Frequency

$tf_{td} = \text{count}(t) \text{ in } d$

Can take log frequency: $\log(1 + tf_{td})$

Can normalize to show proportion: $tf_{td} / \sum_t tf_{td}$

34



35

Keyword Weighting

Term Frequency

$$tf_{td} = \text{count}(t) \text{ in } d$$

TF.IDF: Term Freq by Inverse Document Freq

$$tf.idf_{td} = \log(1 + tf_{td}) \times \log(N/df_i)$$

$$df_i = \# \text{ docs containing } t; \quad N = \# \text{ of docs}$$

36

Limitations of Frequency Statistics

Typically focus on unigrams (single terms)

Often favors frequent (TF) or rare (IDF) terms

Not clear that these provide best description

“Bag of words” ignores additional info.

Grammar / part-of-speech

Position within document

Recognizable entities

42

How do people describe text?

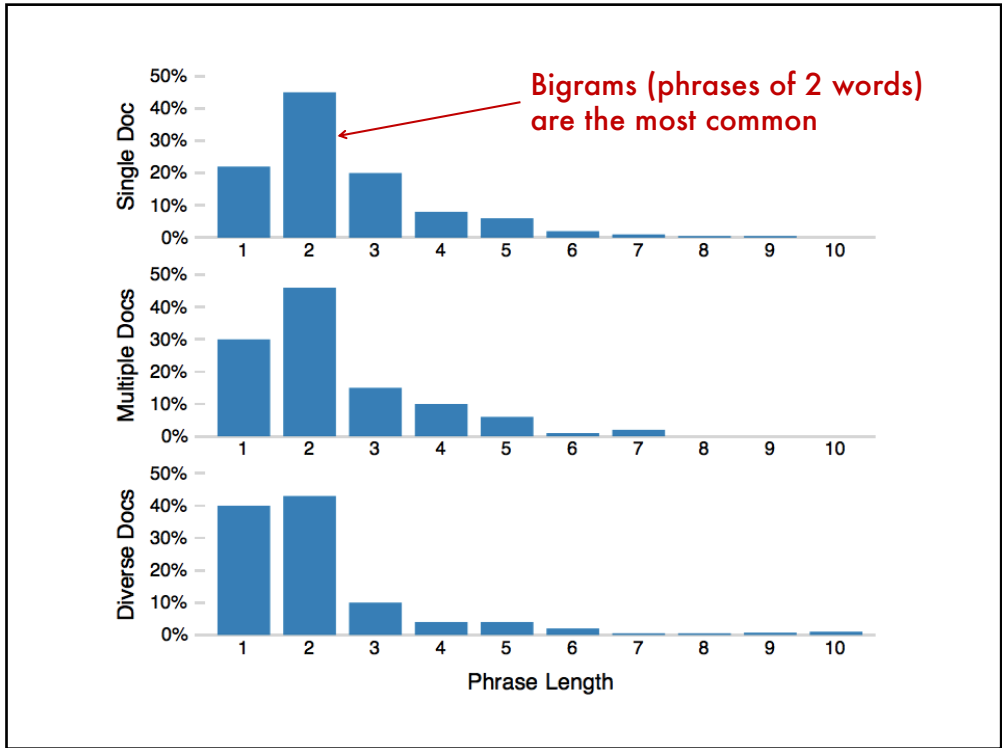
Asked 69 graduate students to read and describe dissertation abstracts

Each given 3 documents in sequence; summarized each using keyphrases, then summarized the 3 together as a whole using keyphrases

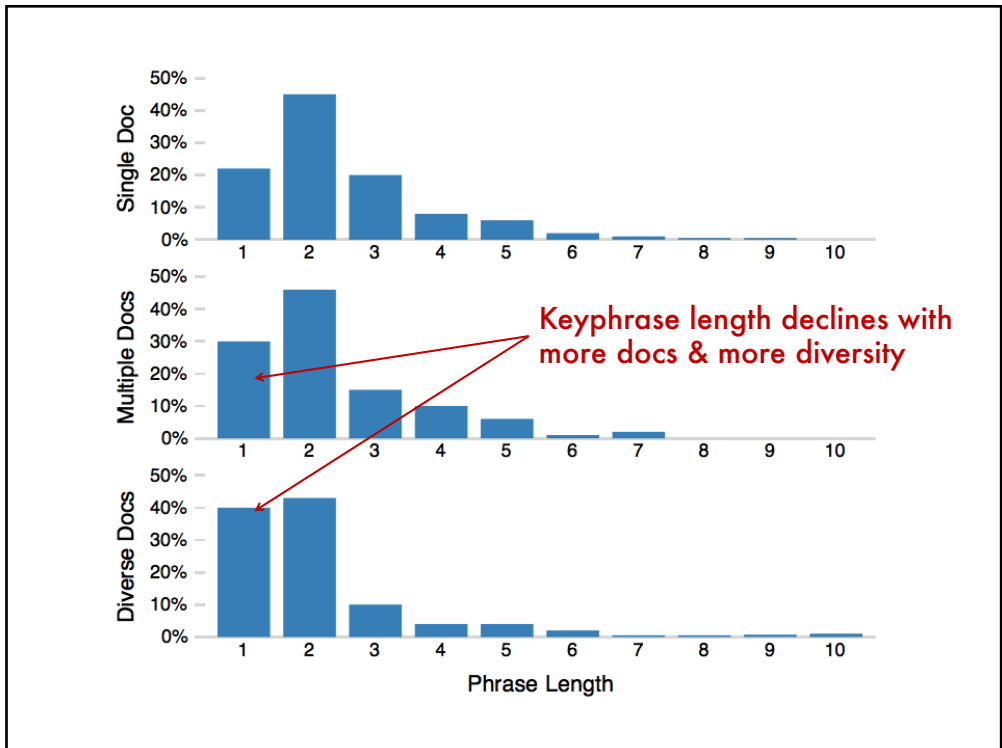
Were matched to both *familiar* and *unfamiliar* topics; *topical diversity* within a collection was varied systematically

[Chuang 2012]

43



44



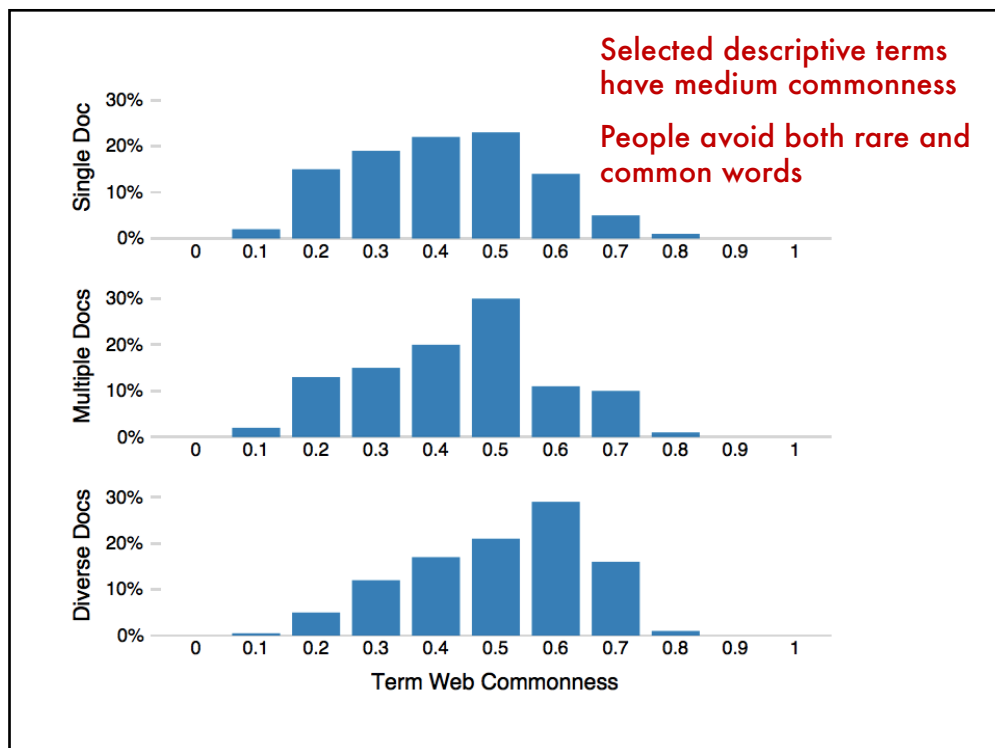
45

Term Commonness

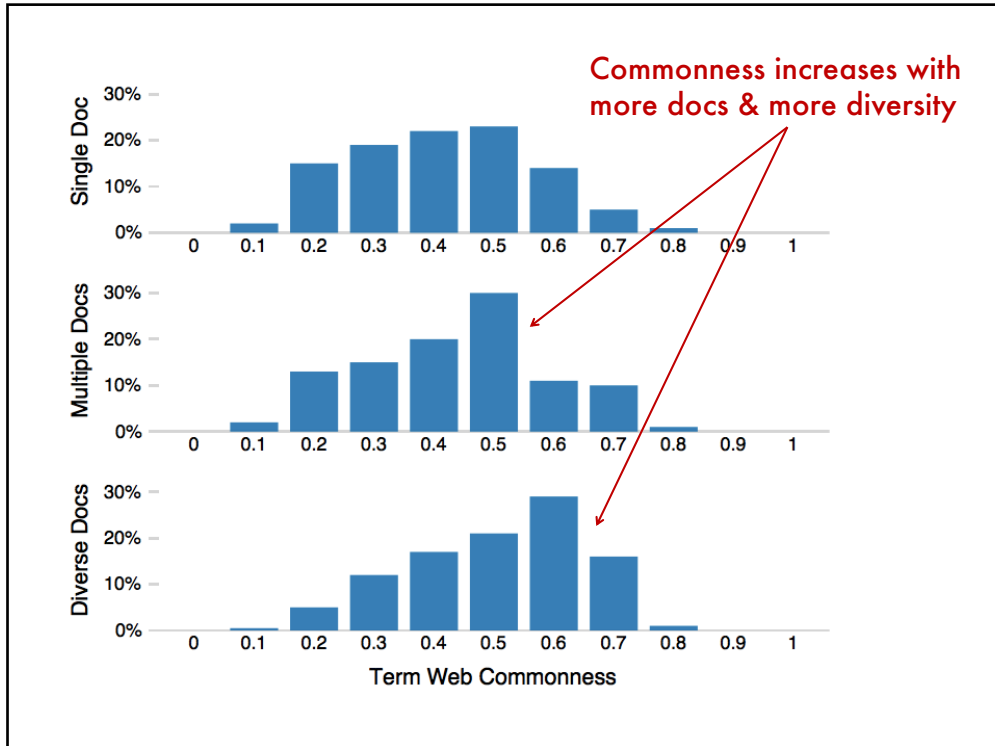
$$\log(\text{tf}_w) / \log(\text{tf}_{\text{the}})$$

The normalized term frequency relative to the most frequent n-gram, e.g., the word "the".

46

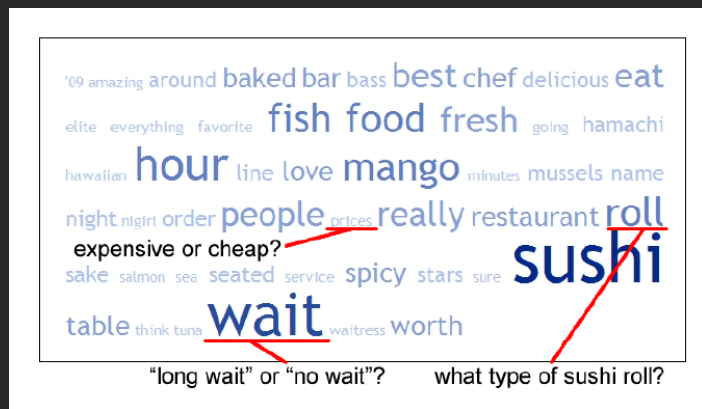


47



48

Yelp: Review Spotlight [Yatani 2011]



52

Yelp: Review Spotlight [Yatani 2011]



53

Tips: Descriptive Keyphrases

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

54

Visualizing Document Content

55

Information Retrieval

Search for documents
Match query string with documents
Visualization to **contextualize results**

The screenshot shows a Google Scholar search for 'acronym resolution'. The search bar at the top contains the query and a magnifying glass icon. Below the search bar, it indicates 'About 154,000 results (0.04 sec)'. The results are listed in a table-like format with columns for article details, abstracts, and PDF links. The first result is 'A supervised learning approach to acronym identification' by D. Nadeau and P. D. Turney, published in the 'Conference of the Canadian Society for ...' in 2005. The second result is 'Leveraging PubMed to Create a Specialty-Based Sense Inventory for Spanish Acronym Resolution' by A. Pomares-Qumbaya and P. López-Ubeda, published in 'Studies in health ...' in 2020. The third result is 'Using word embeddings for unsupervised acronym disambiguation' by J. Charbonnier and C. Wartena, published in 'servis.bib.hs-hannover.de' in 2018. The fourth result is 'SLD: a folk acronym?' by G. A. Ringwood, published in 'ACM Sigplan Notices' in 1989. The left sidebar contains filters for 'Articles', 'Any time', 'Sort by relevance', and 'Create alert'.

Articles	About 154,000 results (0.04 sec)
<p>Any time Since 2020 Since 2019 Since 2016 Custom range...</p> <p>Sort by relevance Sort by date</p> <p><input checked="" type="checkbox"/> include patents <input checked="" type="checkbox"/> include citations</p> <p><input type="checkbox"/> Create alert</p>	<p>A supervised learning approach to acronym identification D. Nadeau, P. D. Turney - Conference of the Canadian Society for ... 2005 - Springer ... 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 110 Related articles All 20 versions</p> <p>[PDF] Leveraging PubMed to Create a Specialty-Based Sense Inventory for Spanish Acronym Resolution A. Pomares-Qumbaya, P. López-Ubeda ... - Studies in health ... 2020 - researchgate.net Acronyms frequently occur in clinical text, which makes their identification, disambiguation and resolution an important task in clinical natural language processing. This paper contributes to acronym resolution in Spanish through the creation of a set of sense ... ☆ All 4 versions</p> <p>[PDF] Using word embeddings for unsupervised acronym disambiguation J. Charbonnier, C. Wartena - 2018 - servis.bib.hs-hannover.de ... Thus, although the goal of our work is acronym expansion, the work is more related to word sense disambiguation (WSD) than to typical work on acronym resolution. The main difference with WSD is that we do not have dictionaries with description of possible senses ... ☆ Cited by 11 Related articles All 6 versions</p> <p>SLD: a folk acronym? G. A. Ringwood - ACM Sigplan Notices, 1989 - dl.acm.org ... 1974]. In view of its relation to Linear Input resolution and confusion over its heritage</p>

56

User Query
(Enter words for different topics on different lines.)

osteoporosis

prevention

research

Run Search
New Query
Quit

Search Limit: 50 100 **250** 500 1000

Number of Clusters: 3 4 **5** 8 10

Mode: TileBars

Cluster
Titles
Backup

FR88513-0157

AP: Groups Seek \$1 Billion a Year for Aging Research

SJMN: WOMEN'S HEALTH LEGISLATION PROPOSED C...

AP: Older Athletes Run For Science

FR: Committee Meetings

FR: October Advisory Committees; Meetings

FR88120-0046

FR: Chronic Disease Burden and Prevention Models; Program...

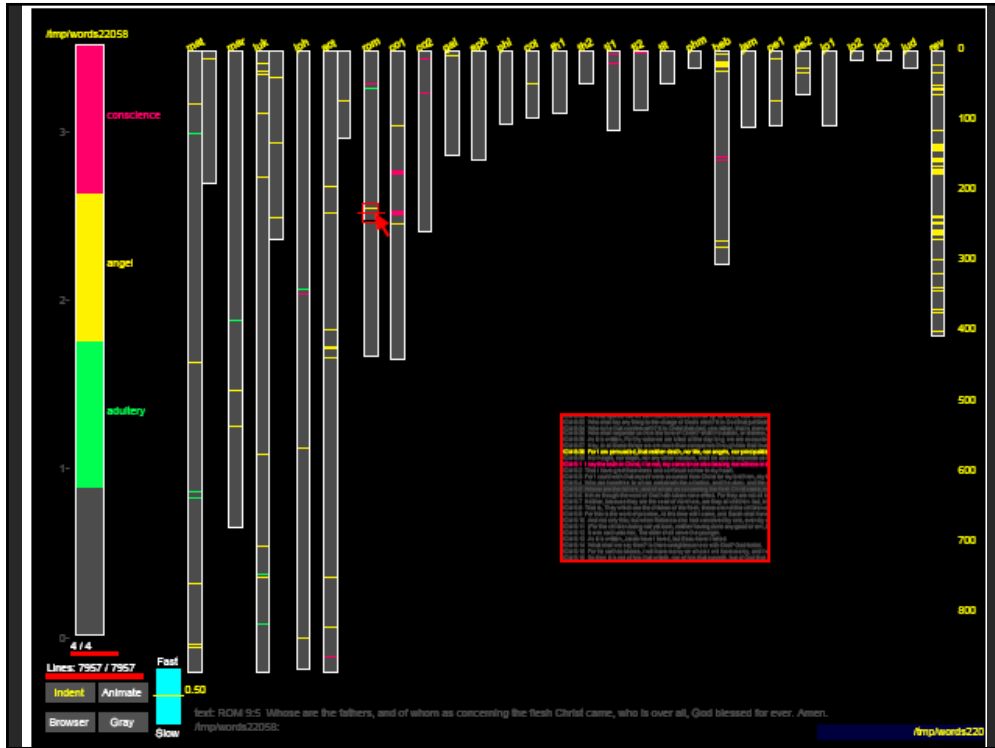
AP: Survey Says Experts Split on Diversion of Funds for AIDS...

FR: Consolidated Delegations of Authority for Policy Developm...

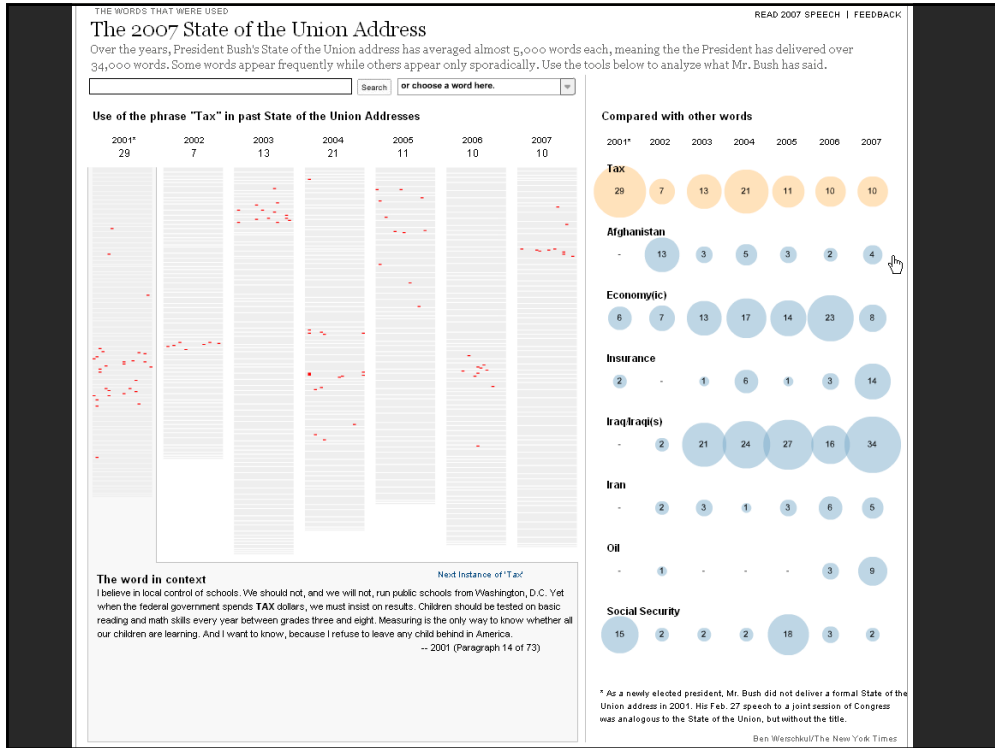
SJMN: RESEARCH FOR BREAST CANCER IS STUCK IN P...

TileBars [Hearst]

57



58



59

Concordance

What is the common local context of a term?

Concordance - Larkin Concordance

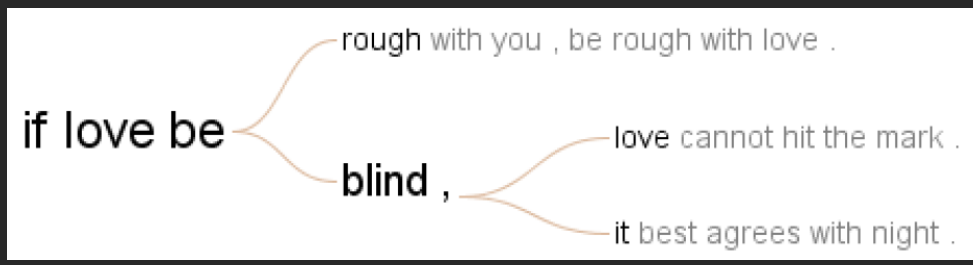
File Text Search Edit Headwords Contexts View Tools Help

Headword	No.	Context...	Word	...Context	Reference
HEAR	15	That my own	heart	drifts and cries, having no...	Deep Analysis
HEARD	9	By the shout of the	heart	continually at work	And the wave
HEARING	7	Nothing to adapt the skill of the	heart	to, skill	And the wave
HEARS	3	The tread, the beat of it, it is my own	heart	,	Träumerei
HEARSE	1	Because I follow it to my own	heart	,	Many famous
HEART	25	My	heart	is ticking like the sun;	I am washed i
HEART'S	2	The vague	heart	sharpened to a candid co...	The March Pa
HEART-SHAPED	1	Contract my	heart	by looking out of date.	Lines on a Yo
HEARTH	1	Having no	heart	to put aside the theft	Home is so Se
HEARTS	7	And the boy puking his	heart	out in the Gents	Essential Bea
HEARTY	1	A harbour for the	heart	against distress.	Bridge for the
HEAT	6	These I would choose my	heart	to lead	After-Dinner F
HEAT-HAZE	1	Time in his little cinema of the	heart	,	Time and Spa
HEATH	1	This petrified	heart	has taken,	A Stone Churc
HEATS	1	How should they sweep the girl clean...	heart	,	I see a girl dra
HEAVE	1	Hands that the	heart	can govern	Heaviest of fk
HEAVEN	4	For the	heart	to be loveless, and as col...	Dawn
HEAVEN-HOLDING	1	With the unguessed-at	heart	riding	One man walk
HEAVIER-THAN...	1	If hands could free you,	heart	,	If hands could
HEAVIEST	2	That overflows the	heart	,	Pour away thi

Words: 7318 Tokens: 37070 At word: 2990 Deleted lines: 1 [24] Word sort: Asc alpha (string) Context sort: Asc occurrence order

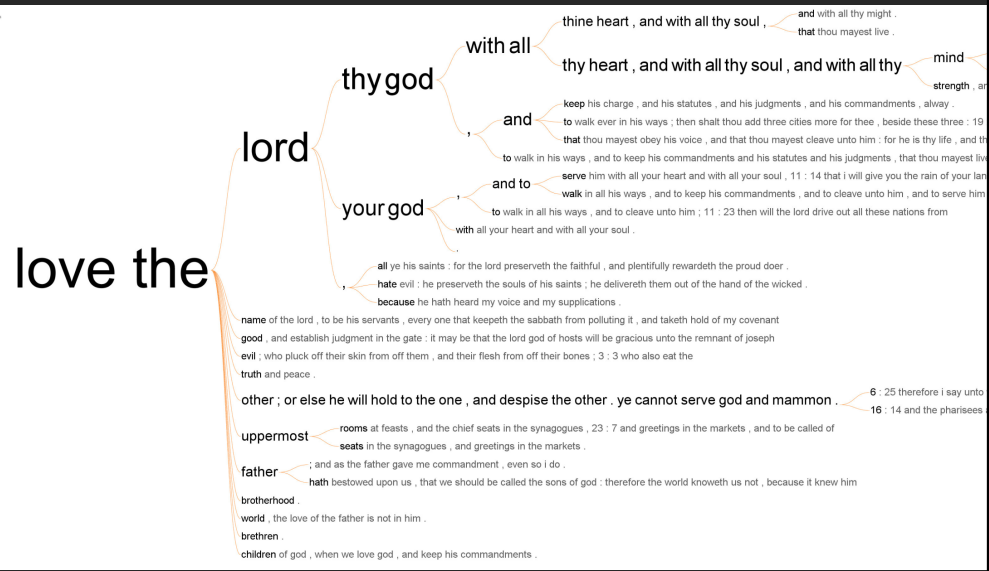
62

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 .



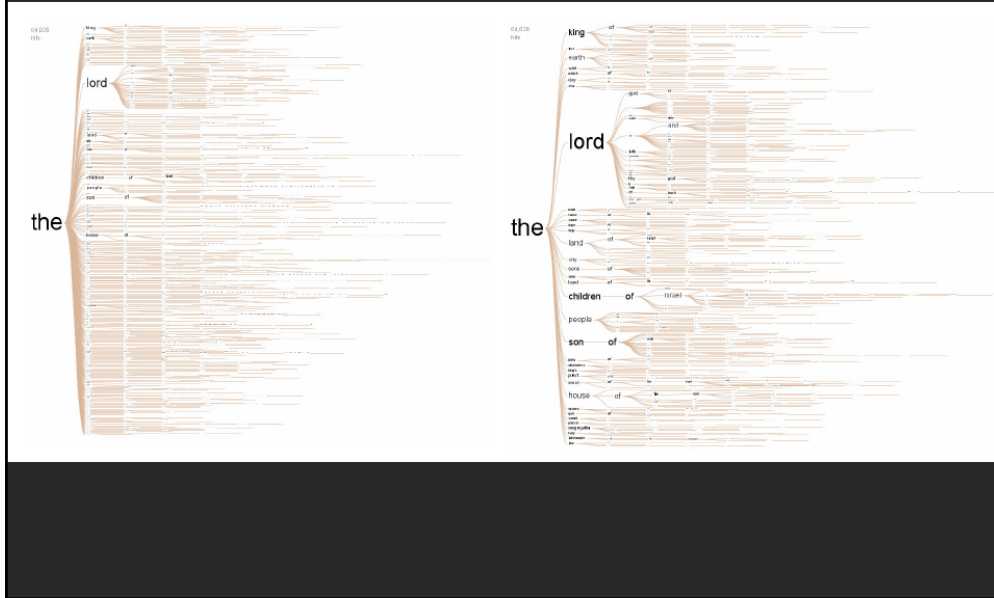
64

WordTree



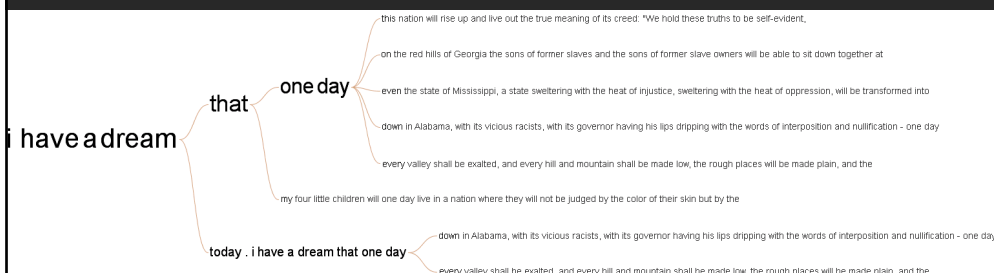
65

Filter infrequent runs



66

Recurrent themes in speech



67

Glimpses of structure

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

For example

Lexical: <A> at

Syntactic: <Noun> <Verb> <Object>

70

Phrase Nets [van Ham 2009]

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 extracted patterns in a node-link view

Occurrences → Node size

Pattern position → Edge direction

Darker color → higher ratio of out-edges to in-edges

71

Visualizing Conversation

91

Visualizing Conversation

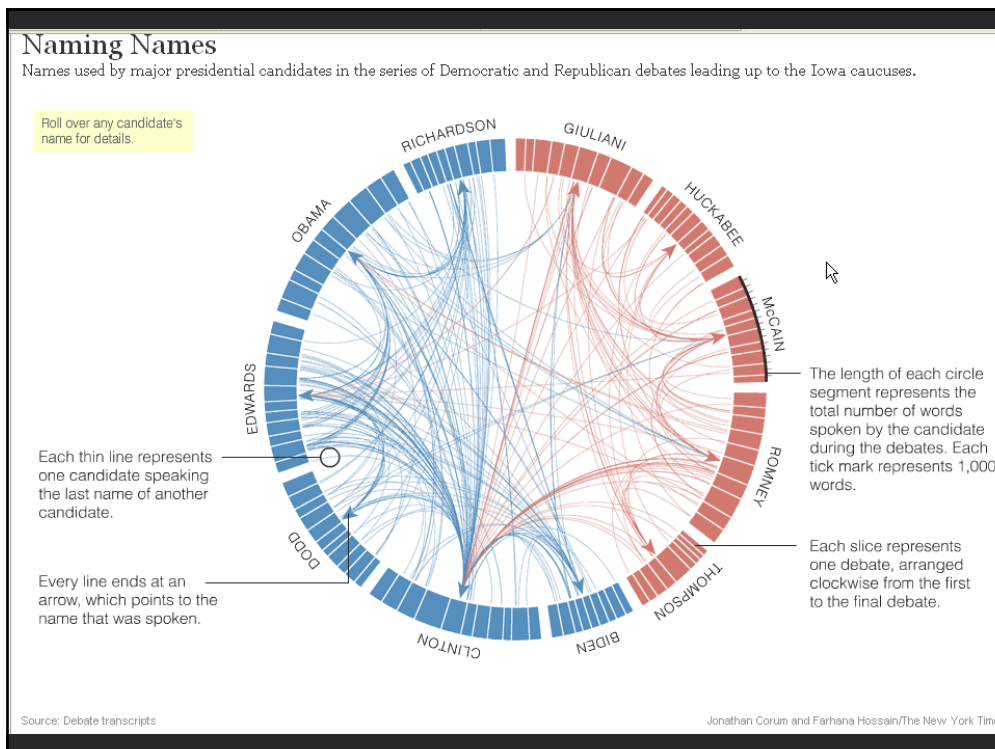
Many dimensions to consider:

- Who (senders, receivers)
- What (the content of communication)
- When (temporal patterns)

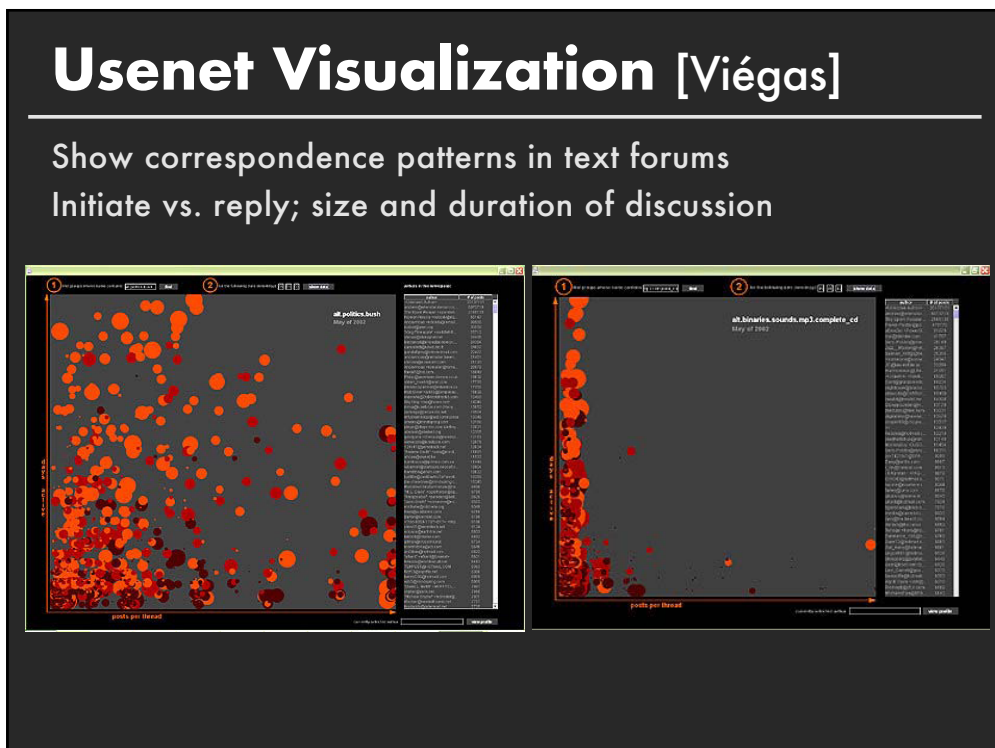
Interesting cross-products:

- What x When → Topic “Zeitgeist”
- Who x Who → Social network
- Who x Who x What x When → Information flow

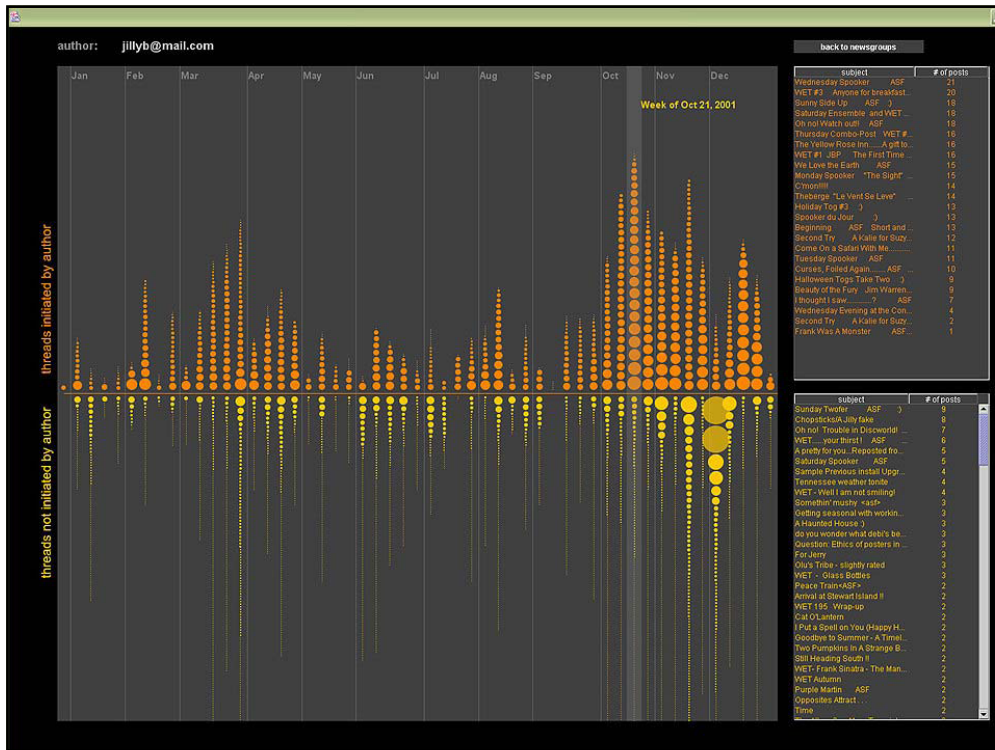
92



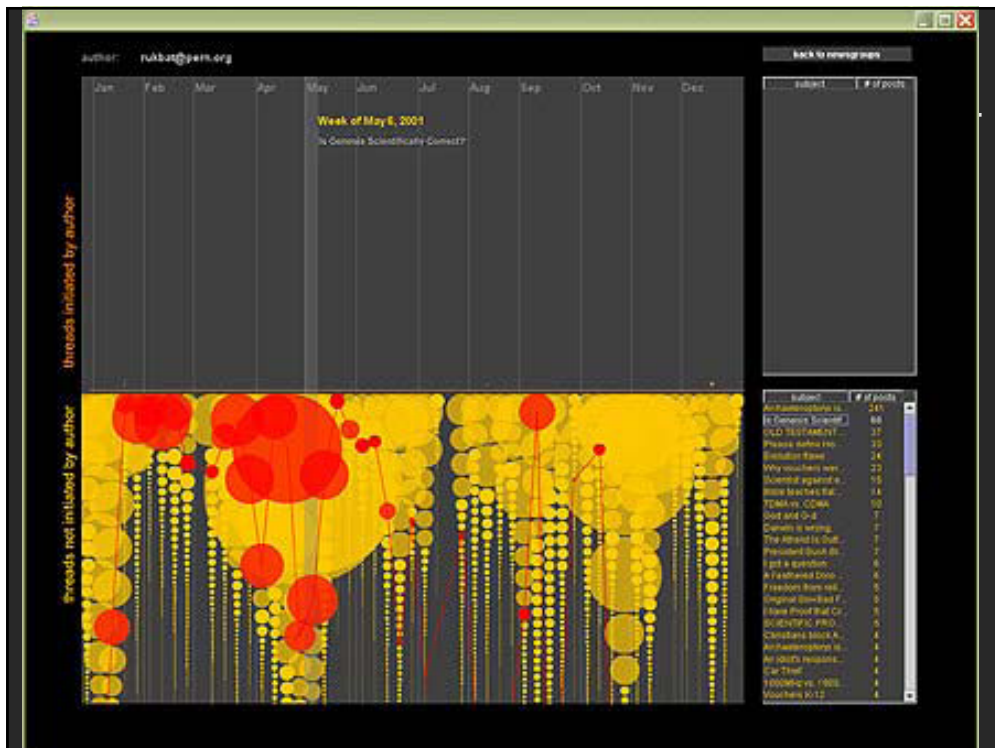
93



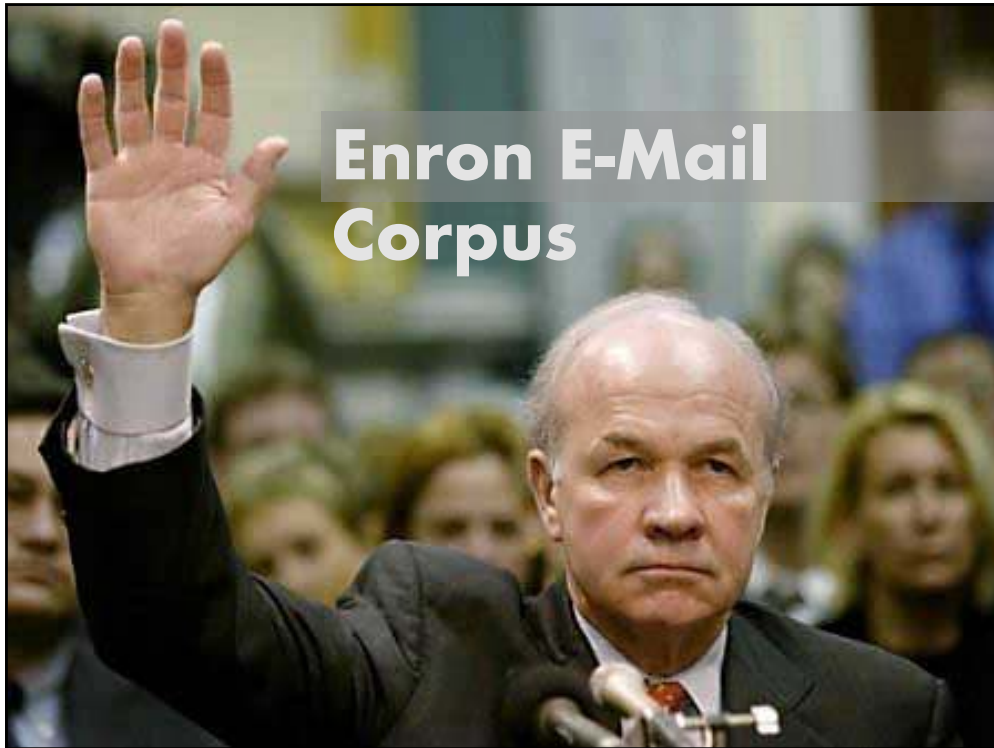
96



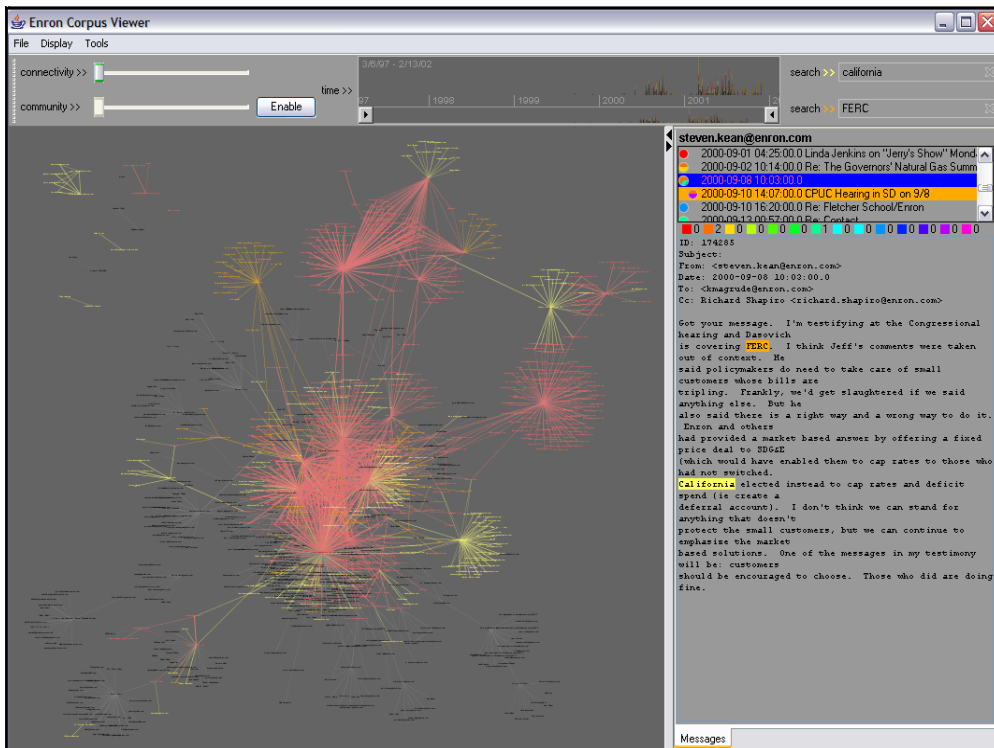
97



98



101



102

The screenshot displays the 'Enron Corpus Viewer' interface. On the left, a network graph titled 'Washington Lobby?' shows connections between various Enron employees, with nodes labeled with names and email addresses such as 'john.shelk@enron.com', 'Susan Mara', 'Tim Belden', and 'Dan Laff'. The graph is overlaid on a background of red lines. At the top, there are controls for 'connectivity' and 'community', a date range '1/20/01 - 6/27/01', and search fields containing 'california' and 'ferc'. On the right, a news article is displayed with the headline 'Enron 'Mastermind' Pleads Guilty' and a sub-headline 'SAN FRANCISCO, Oct. 17, 2002'. The article text includes: '(AP) A former top energy trader, considered the mastermind of Enron Corp.'s scheme to drive up California's energy prices, pleaded guilty Thursday to a federal conspiracy charge. Timothy Belden, the former head of trading in Enron's Portland, Ore., office, admitted to one count of conspiracy to commit wire fraud and promised to cooperate with state and federal prosecutors as well as any non-criminal effort to investigate the energy industry. "I did it because I was trying to maximize profit for Enron," Belden told U.S. District Judge Martin Jenkins.'

103

Visualizing Document Collections

104

Named Entity Recognition

Identify and classify named entities in text:

John Smith → PERSON

Soviet Union → COUNTRY

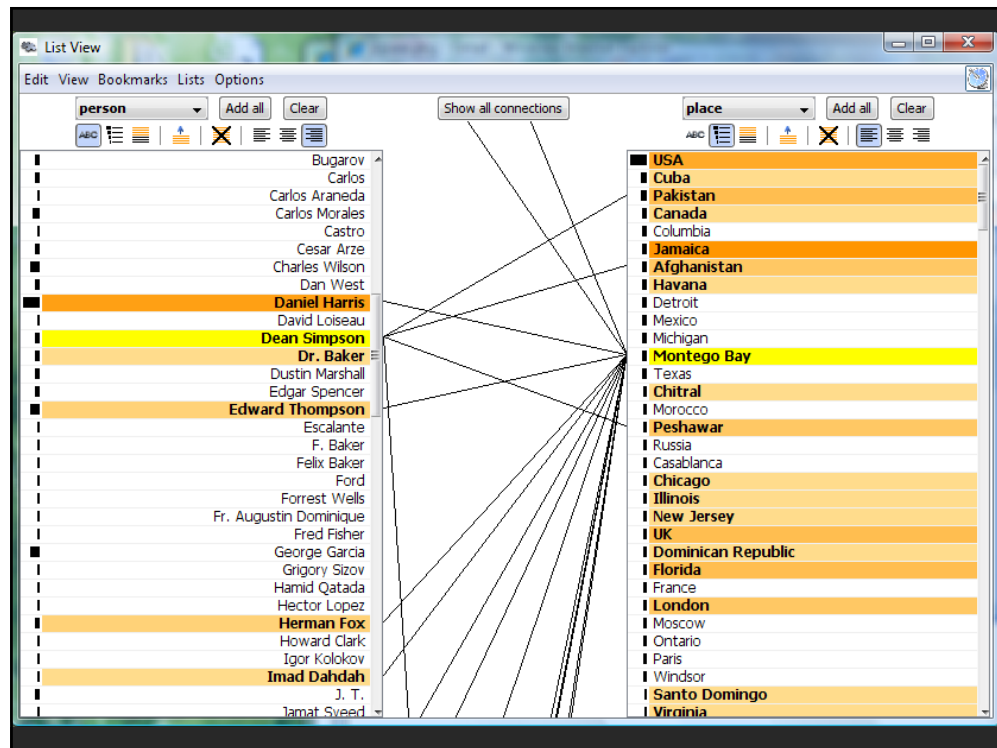
353 Serra St → ADDRESS

(555) 721-4312 → PHONE NUMBER

Entity relations: how do the entities relate?

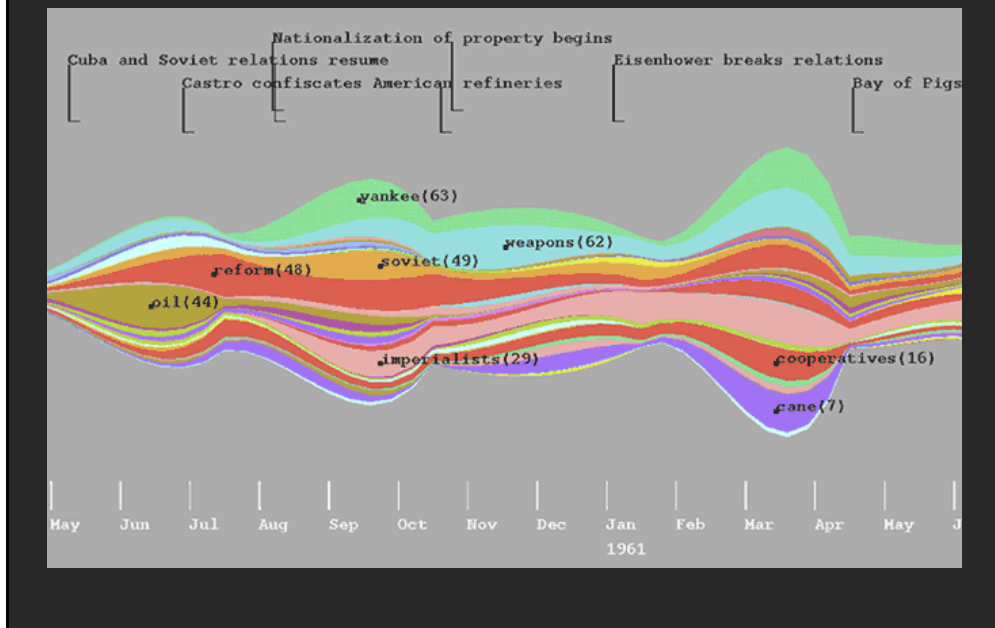
Simple approach: do they co-occur in small window of text?

107

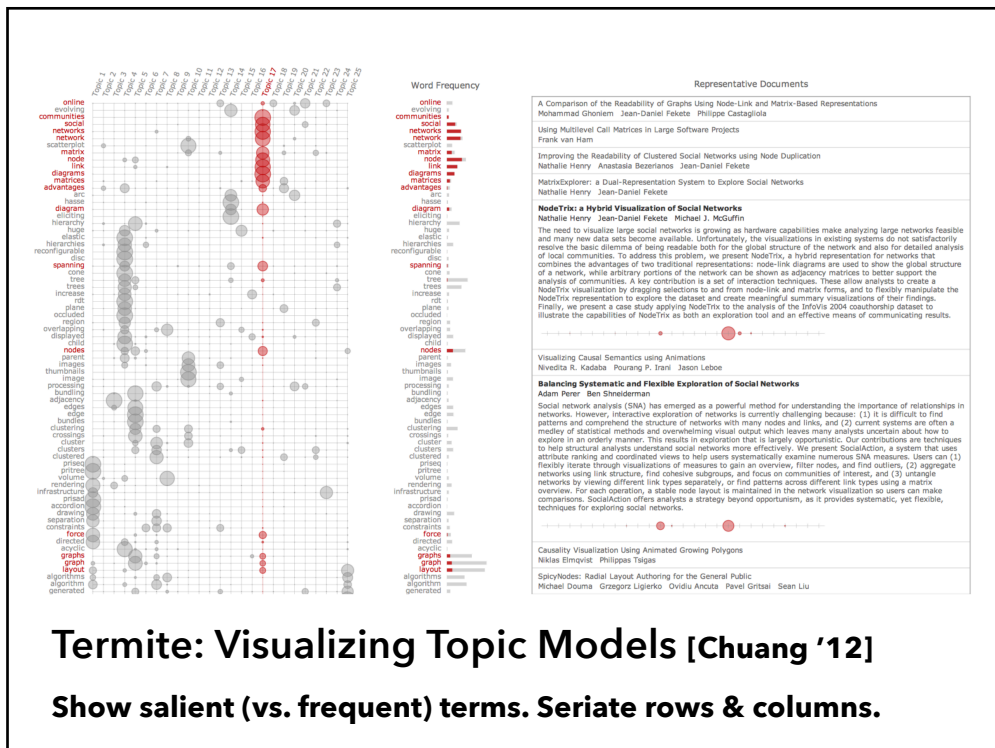


108

ThemeRiver (Havre et al 99)



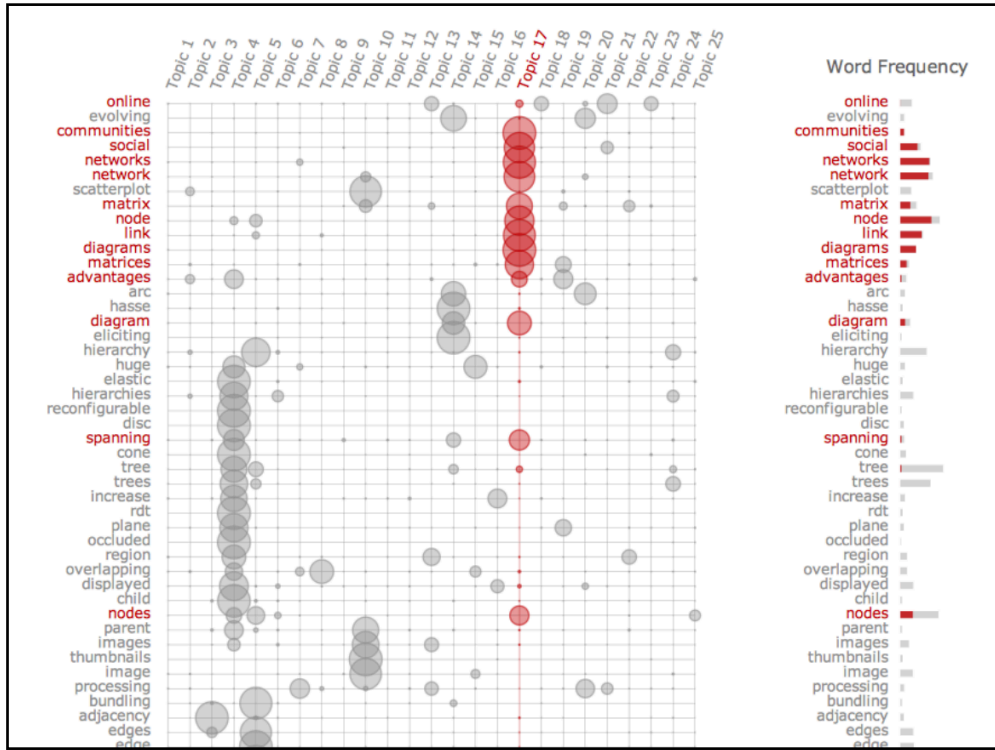
114



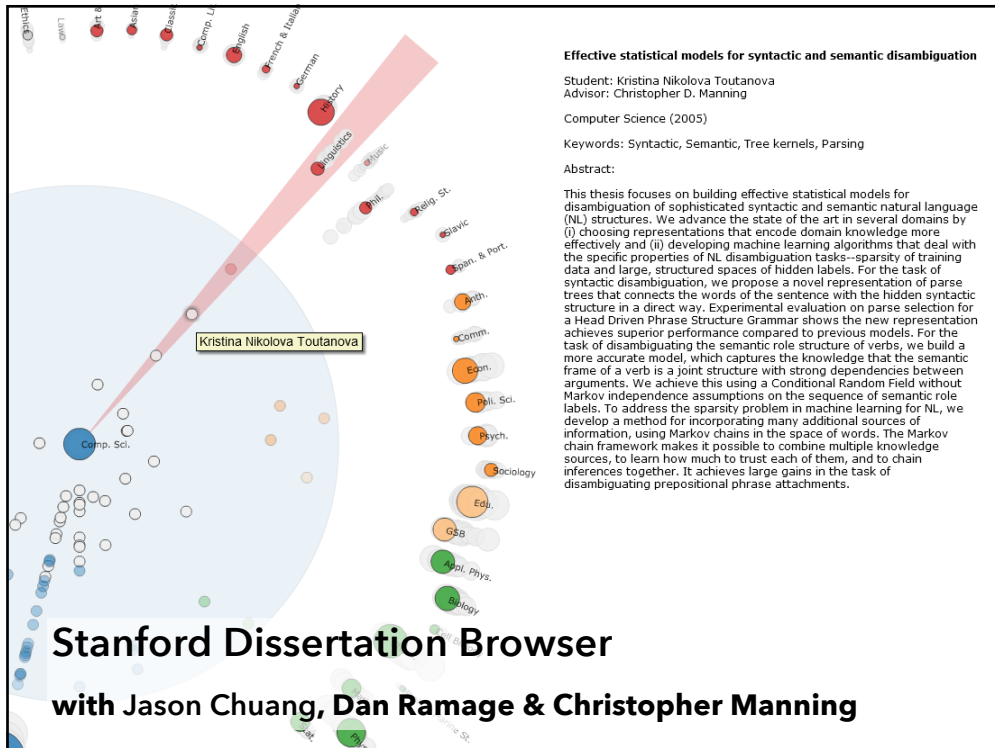
122

Termite: Visualizing Topic Models [Chuang '12]

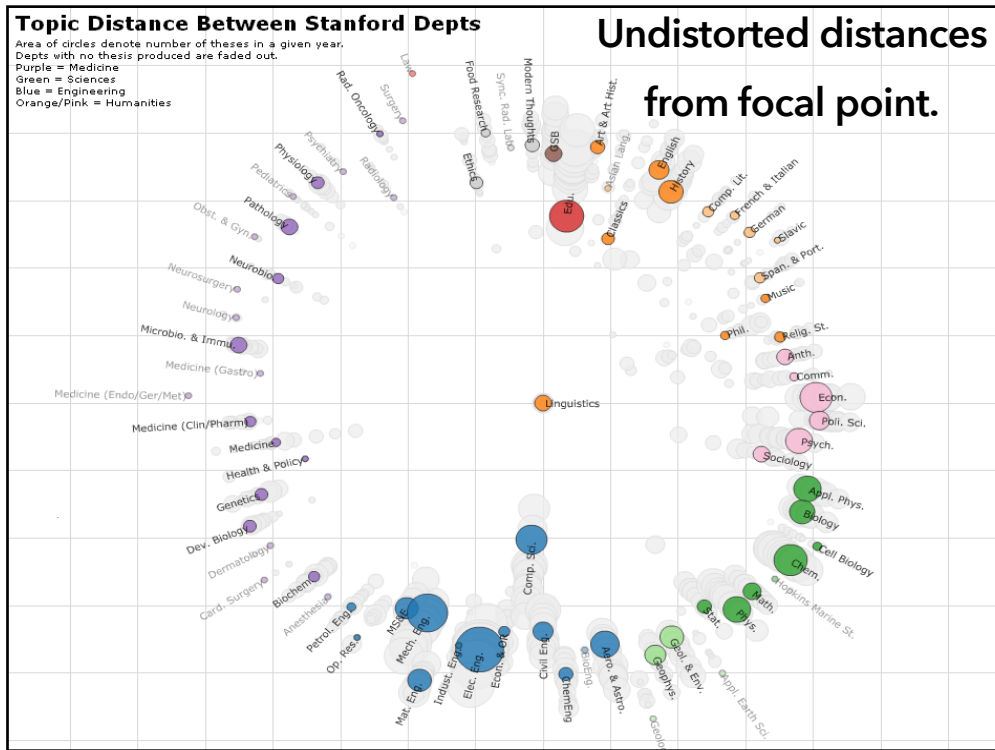
Show salient (vs. frequent) terms. Seriate rows & columns.



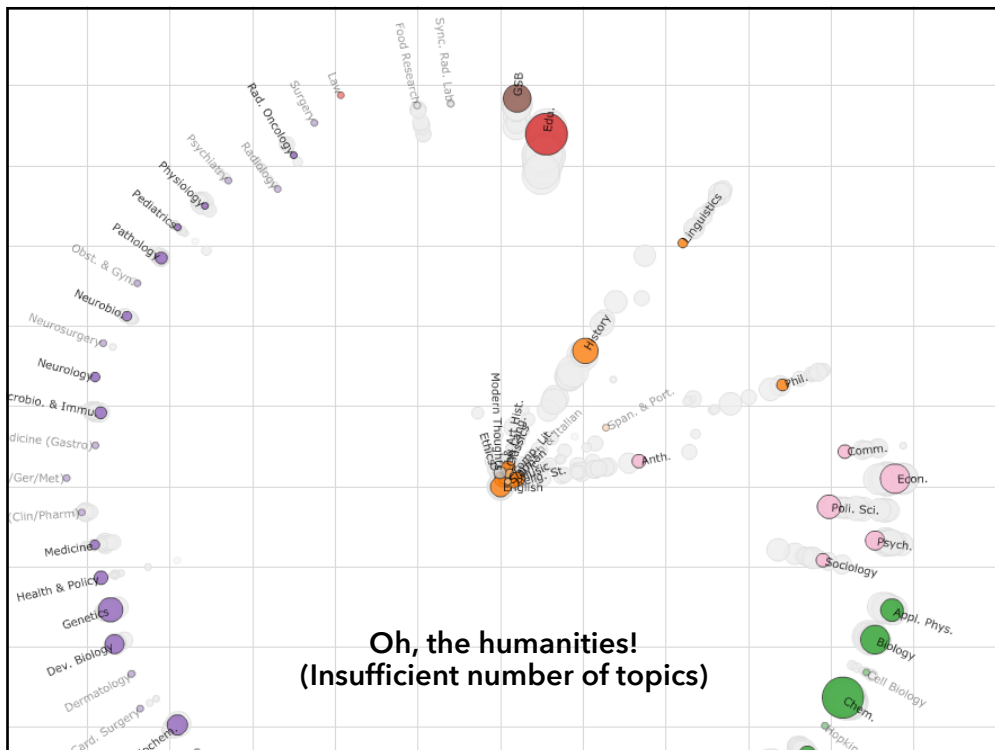
123



124



126



127

Summary

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

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

Currently: from bag-of-words to *vector space embeddings*