

# Text Visualization

Maneesh Agrawala

CS 448B: Visualization  
Fall 2021

1

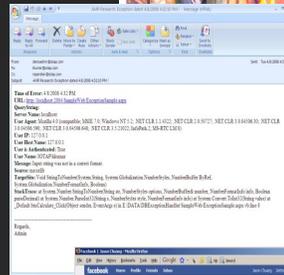
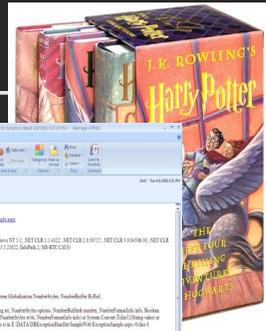
## Text as data

### 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

# **Text Visualization**

3

# **Why visualize text?**

4

## 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

5

## 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?**

6







# Topics

---

**Text as Data**

**Visualizing Document Content**

**Visualizing Conversation**

**Document Collections**

15

**Text as Data**

16

# 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

17

# 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.*

18

# 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*

19

# 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

20

# 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

21

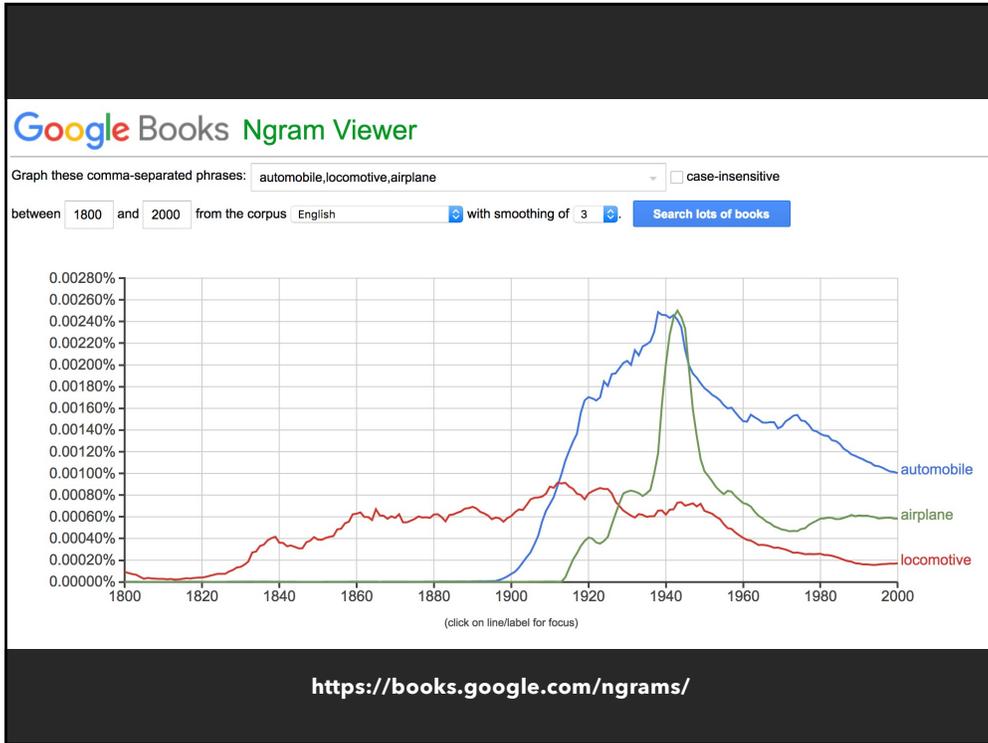
# Document $\times$ Term matrix

Each document is a vector of term weights

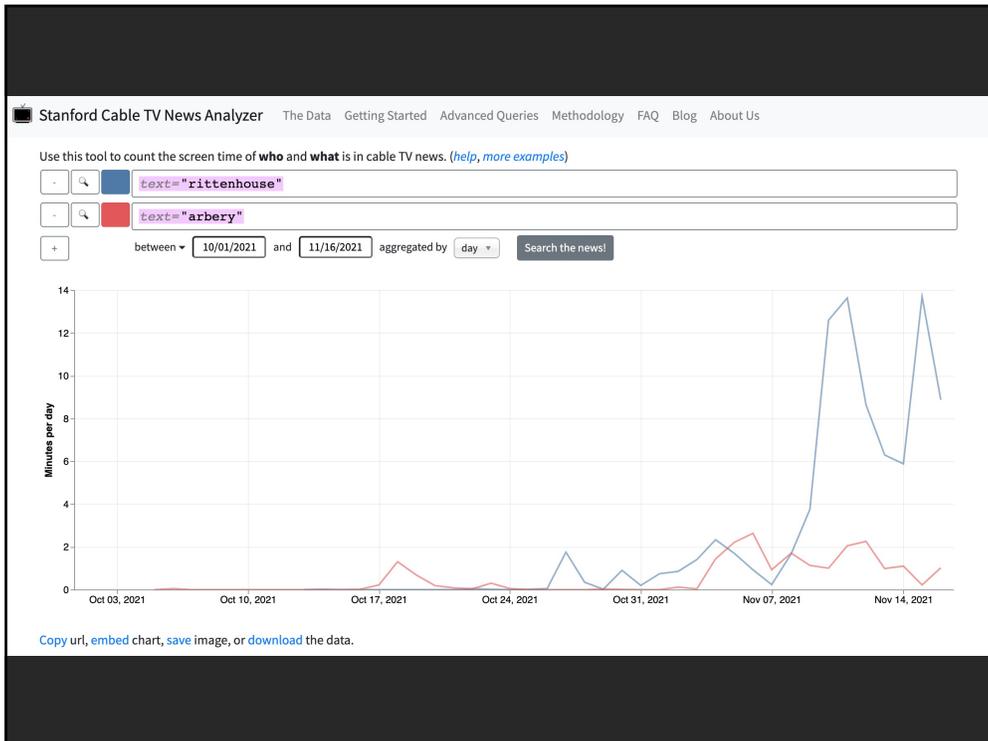
Simplest weighting is to just count occurrences

	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

22



24



26



# Announcements

29

## 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 different interactive visualizations
- **Research:** Implementation of solution and web-based demo if possible
- **Short video (2 min)** demoing and explaining the project

### Schedule

- Project proposal: **Wed 11/3**
- Design Review and Feedback: **10<sup>th</sup> week of quarter**
- Final code and video: **Fri 12/10 11:59pm**

### Grading

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

30

## Feedback (Week 10)

---

### Signup for a ~ 10 min slot

[https://docs.google.com/spreadsheets/d/1U-Q7DVvWTmTt\\_nYumJlqSDgySvAd2lq0EuUXGVHo4cE/edit?usp=sharing](https://docs.google.com/spreadsheets/d/1U-Q7DVvWTmTt_nYumJlqSDgySvAd2lq0EuUXGVHo4cE/edit?usp=sharing)

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

31

**Given a text, what are the best descriptive words?**

32

# Keyword Weighting

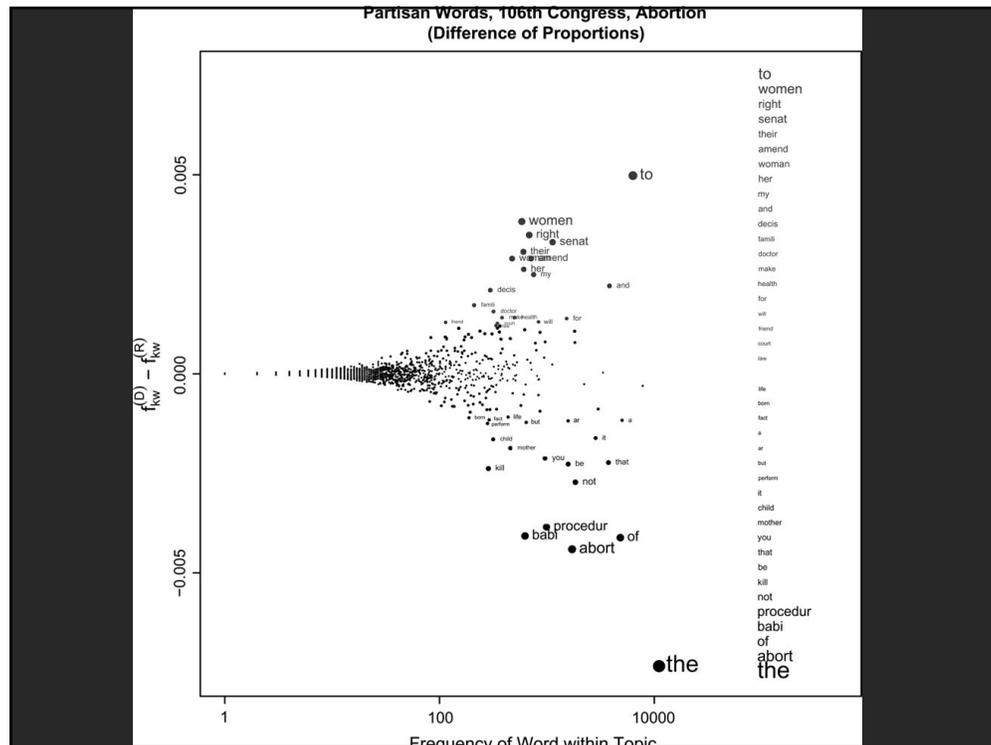
## Term Frequency

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

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

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

33



34

# Keyword Weighting

## Term Frequency

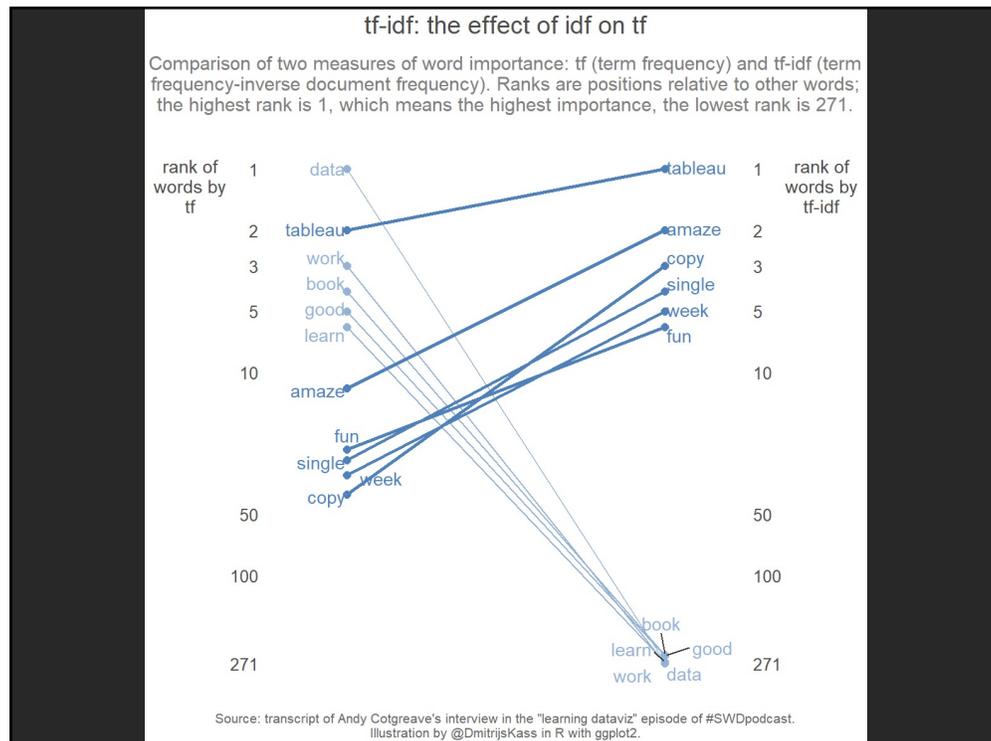
$$tf_{td} = \text{count}(t) \text{ in } d$$

## TF.IDF: Term Freq by Inverse Document Freq

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

$df_t$  = # docs containing  $t$ ;  $N$  = # of docs

35



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?

---

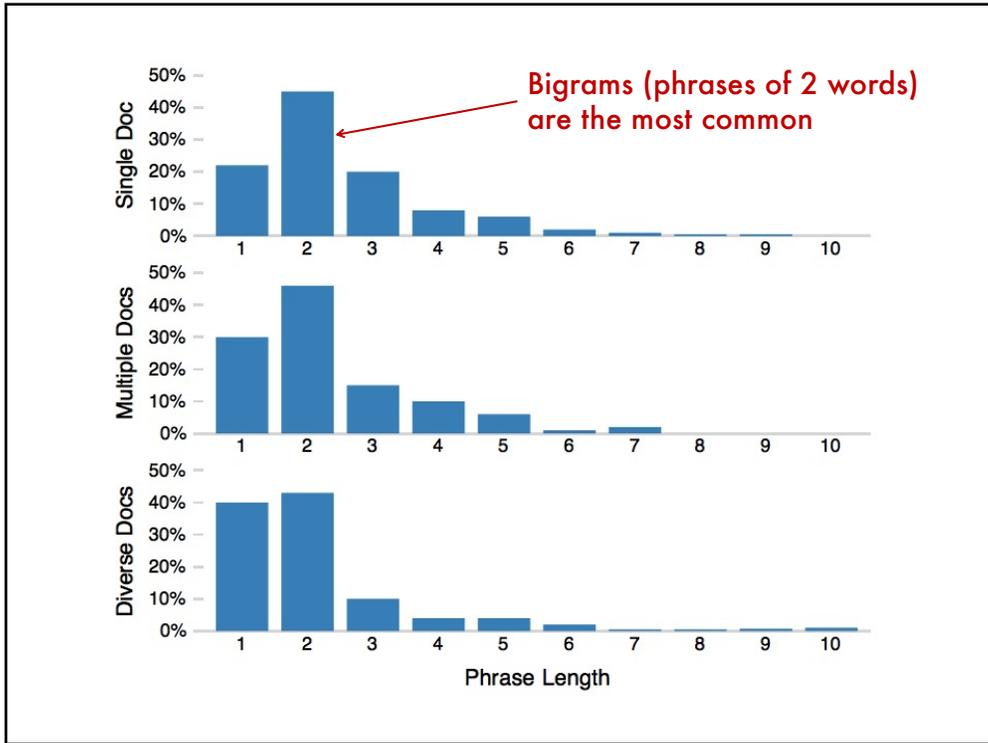
Asked 69 graduate students to read and describe dissertation abstracts

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

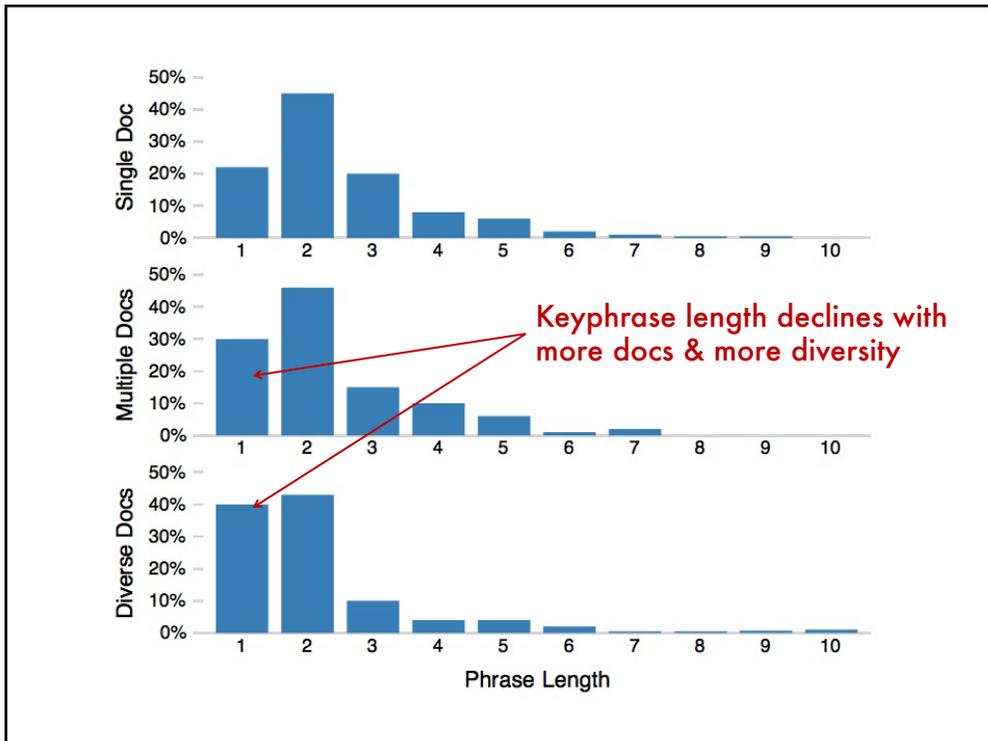
Were matched to both *familiar* and *unfamiliar* topics; *topical diversity* within a collection was varied systematically

[Chuang 2012]

43



44



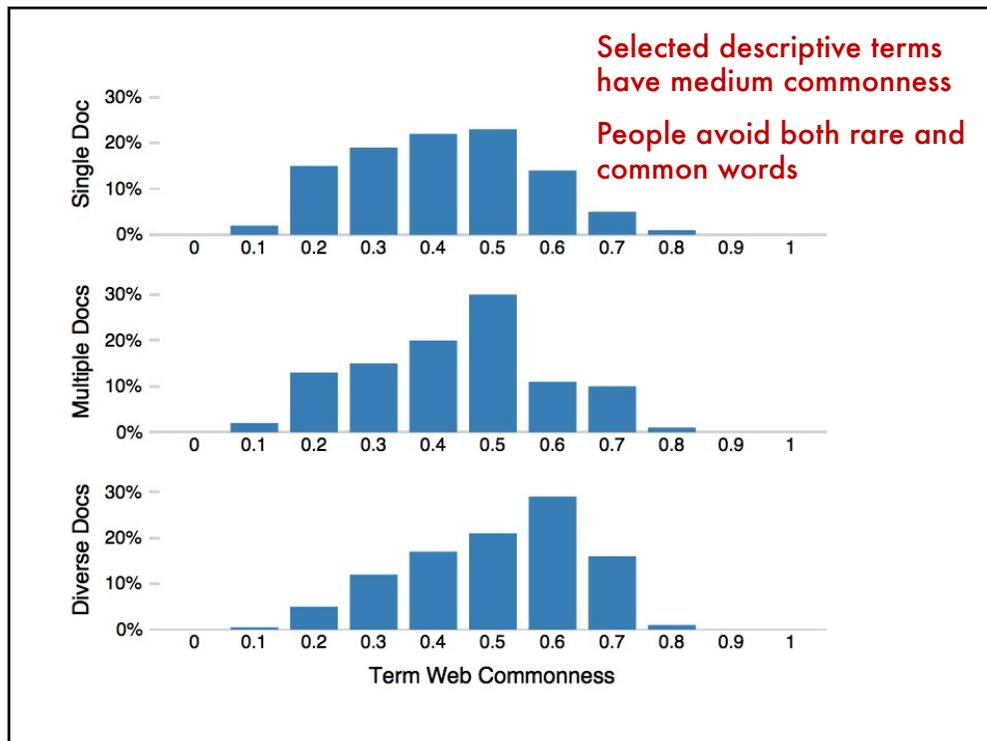
45

## Term Commonness

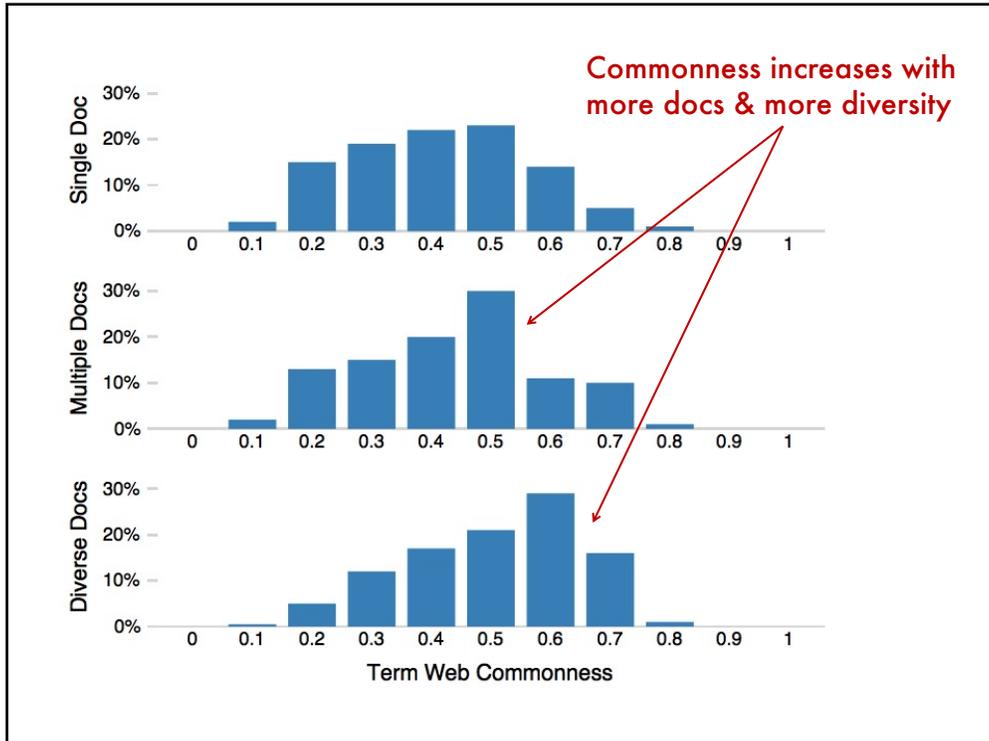
$$\log(\text{tf}_w) / \log(\text{tf}_{\text{the}})$$

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

46

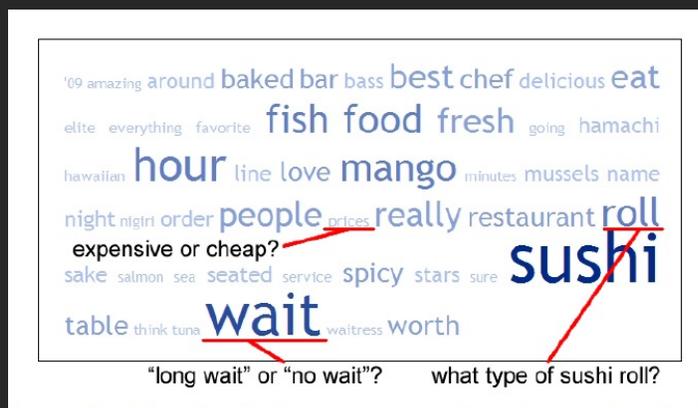


47



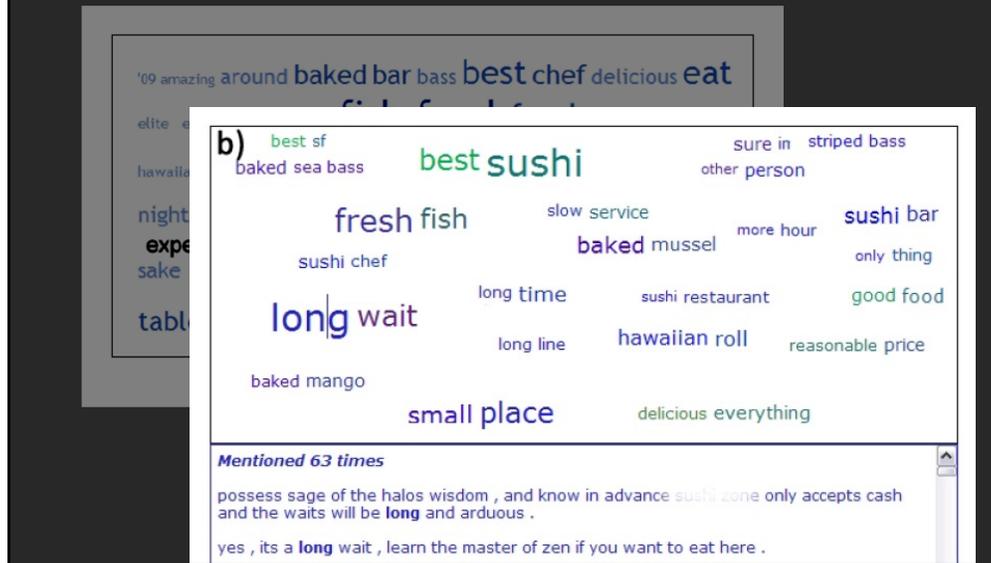
48

## 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)'. On the left side, there are filters for 'Articles', 'Any time', 'Since 2020', 'Since 2019', 'Since 2016', and 'Custom range...'. There are also options to 'Sort by relevance' or 'Sort by date', and checkboxes for 'include patents', 'include citations', and 'Create alert'. The main results area displays three entries:

- Entry 1:** A supervised learning approach to acronym identification. Authors: D. Nadeau, F.D. Turney. Conference of the Canadian Society for ... 2005 - Springer. Abstract: ... 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.
- Entry 2:** Leveraging PubMed to Create a Specialty-Based Sense Inventory for Spanish Acronym Resolution. Authors: A. Pomares-Quimbaya, P. López-Úbeda. Studies in health ... 2020 - researchgate.net. Abstract: 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.
- Entry 3:** Using word embeddings for unsupervised acronym disambiguation. Authors: J. Charbonnier, C. Wartena. servwis.bib.hs-hannover.de. Abstract: ... 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.

Additional results include 'SLD: a folk acronym?' by GA Ringwood - ACM Sigplan Notices, 1989 - dl.acm.org and 'Find it@Stanford'.

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

#mp/words22058
0

conscience

angel

adultery

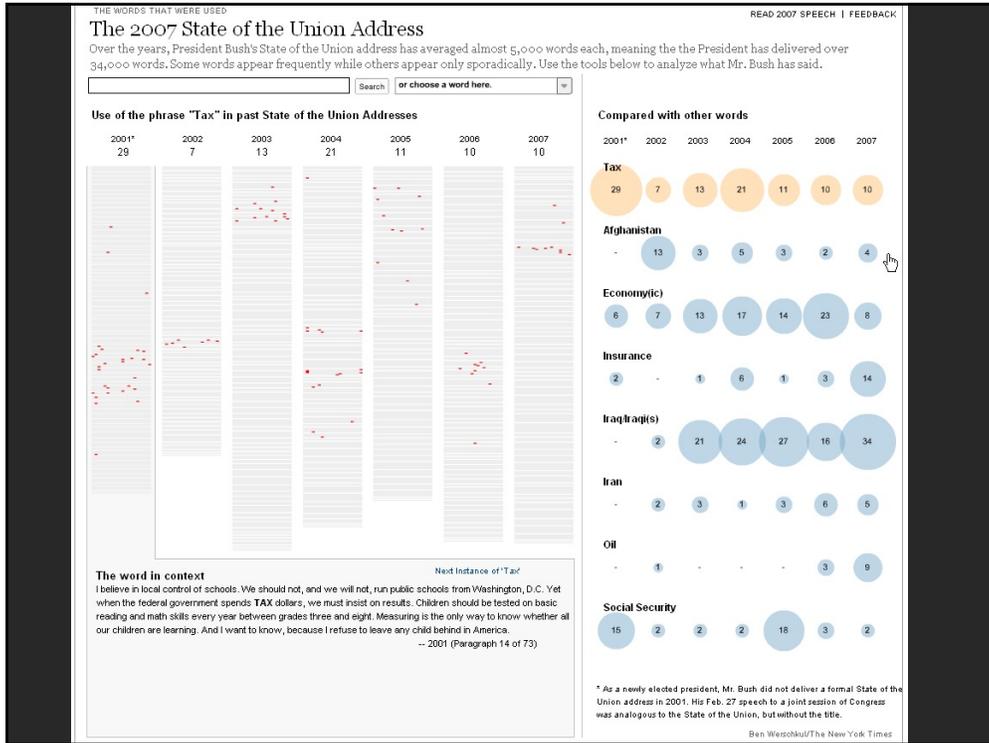
Lines: 7957 / 7957    Fast    0.50

Indent    Animate    Slow

Browser    Gray

feat: RDM 9.5. Whose are the fellows, and of whom as concealing the flesh Christ came, who is over all, God blessed for ever. Amen.  
#mp/words22058

58



59

# Concordance

## What is the common local context of a term?

Concordance - Larkin Concordance

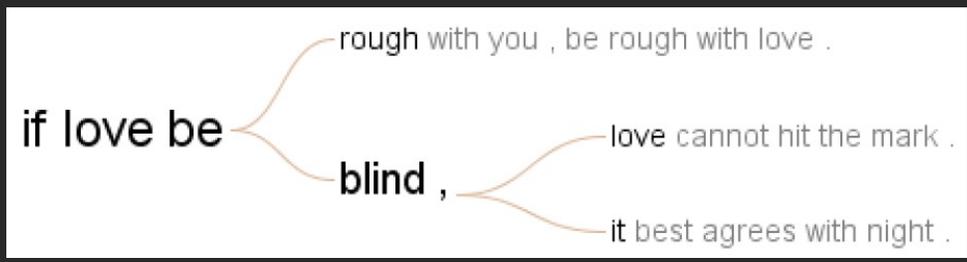
File Text Search Edit Headwords Contexts View Tools Help

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

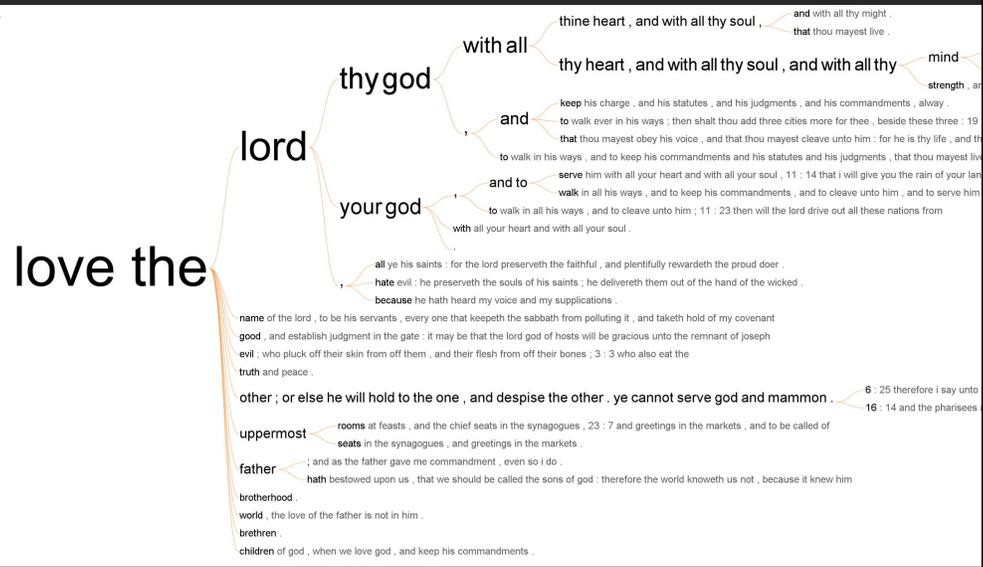
62

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 .



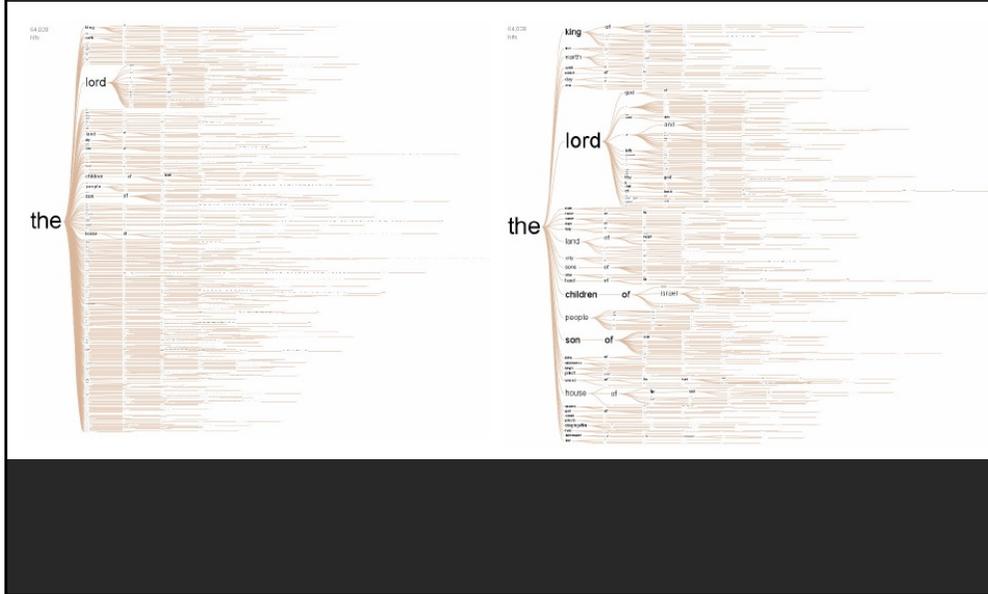
64

# WordTree



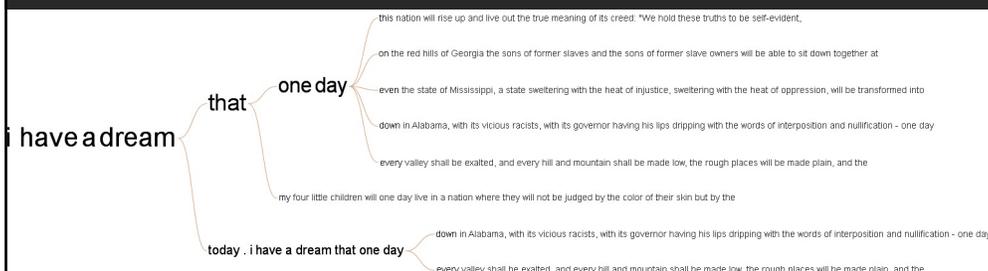
65

# Filter infrequent runs



66

# Recurrent themes in speech



67

**word tree** reverse tree one phrase per line

act

we must act knowing that our work will be imperfect. We must act, knowing that today's victories will be only partial, and that it will be up to those who stand here in four years, and forty years, and four hundred years hence to advance the timeless spirit once conferred to us in a spare Philadelphia hall.

knowing that today's victories will be only partial, and that it will be up to those who stand here in four years, and forty years, and four hundred years hence to advance the timeless spirit once conferred to us in a spare Philadelphia hall.

knowing that our work will be imperfect. We must act, knowing that today's victories will be only partial, and that it will be up to those who stand here in four years, and forty years, and four hundred years hence to advance the timeless spirit once conferred to us in a spare Philadelphia hall.

claim to possess. That's how we will maintain our economic vitality and our national treasure - our forests and waterways; our croplands and snowcapped peaks.

harness new ideas and technology to remake our government, rework our tax code, reform our schools, and empower our citizens with the skills they need to work harder, learn more, reach higher.

make the hard choices to reduce the cost of health care and the size of our deficit. But we reject the belief that America must choose between caring for the generation that built this country and investing in the generation that will build its future.

lead it. We cannot cede to other nations the technology that will power new jobs and new industries - we must claim its promise.

claim to possess. That's how we will maintain our economic vitality and our national treasure - our forests and waterways; our croplands and snowcapped peaks.

carry those lessons into this time as well. We will defend our people and uphold our values through strength of arms and rule of law.

be a source of hope to the poor, the sick, the marginalized, the victims of prejudice - not out of mere charity, but because peace in our time requires the constant advance of those principles that our common creed describes: tolerance and opportunity; human dignity; respect for the rights of all.

size it - so long as we set it together. For we, the people, understand that our country cannot succeed when a shrinking few do very well and a growing many barely make it.

respond to the threat of climate change, knowing that the failure to do so would betray our children and future generations.

maintain our economic vitality and our national treasure - our forests and waterways; our croplands and snowcapped peaks. That is how we will preserve our planet, commanded to our care by God.

preserve our planet, commanded to our care by God. That's what will lead meaning to the creed our fathers once declared.

defend our people and uphold our values through strength of arms and rule of law. We will show the courage to try and resolve our differences with other nations peacefully - not because we are naive about the dangers we face, but because engagement can more durably lift suspicion and fear.

show the courage to try and resolve our differences with other nations peacefully - not because we are naive about the dangers we face, but because engagement can more durably lift suspicion and fear.

rense those institutions that extend our capacity to manage crisis ahead, for no one has a greater stake in a peaceful world than its most powerful nation.

support democracy from Asia to Africa; from the Americas to the Middle East, because our interests and our conscience compel us to act on behalf of those who long for freedom.

made for this moment, and we will seize it - so long as we set it together. For we, the people, understand that our country cannot succeed when a shrinking few do very well and a growing many barely make it.

truth to our creed when a little girl born into the Meadest poverty knows that she has the same chance to succeed as anybody else, because she is an American, she is free, and she is equal, not just in the eyes of God but also

also help to those who won the peace and not just the war, who turned sworn enemies into the surest of friends, and we must carry those lessons into this time as well.

naive about the dangers we face, but because engagement can more durably lift suspicion and fear. America will remain the anchor of strong alliances in every corner of the globe; and we will renew those institutions that extend our capacity to manage crisis ahead, for no one has a greater stake in a peaceful world than its most powerful nation.

truly created equal, then surely the laws we commit to one another must be equal as well. Our journey is not complete until no citizen is forced to wait for hours to exercise the right to vote.

every citizen deserves a basic measure of security and dignity. We must make the hard choices to reduce the cost of health care and the size of our deficit.

understand that our country cannot succeed when a shrinking few do very well and a growing many barely make it.

our obligations as Americans are not just to ourselves, but to all posterity. We will respond to the threat of climate change, knowing that the failure to do so would betray our children and future generation

enduring security and lasting peace do not require perpetual war. Our brave men and women in uniform, tempered by the flames of battle, are unmatched in skill and courage.

declare today that the most evident of truths - that all of us are created equal - is the star that guides us still; just as it guided our forebears through Seneca Falls, and Selma, and Stone Mountain; just as it guided all those men and women

never relinquished our skepticism of central authority, nor have we succumbed to the fiction that all society's ills can be cured through government alone.

Through blood drawn by lash and blood drawn by sword, we learned that no nation founded on the principles of liberty and equality could survive half-slave and half-free.

always understood that when times change, so must we that fidelity to our founding principles requires new responses to new challenges; that preserving our individual freedoms ultimately requires collective action.

code to other nations the technology that will power new jobs and new industries - we must claim its promise. That's how we will maintain our economic vitality and our national treasure - our forests and waterways; our croplands and snowcapped peaks.

walk alone, to bear a King proclaim that our individual freedom is inextricably bound to the freedom of every soul on Earth.

afford delay. We cannot mistake idealism for principle, or substitute spectacle for politics, or treat name-calling as reasoned debate. We must act, knowing that our work will be imperfect.

mislike absolutism for principle, or substitute spectacle for politics, or treat name-calling as reasoned debate. We must act, knowing that our work will be imperfect.

define liberty exactly the same way, or follow the same precise path to happiness. Progress does not compel us to settle centuries-long debates about the role of government for all time - but it does require us to act in our time.

make to the flag that waves above and that fills our hearts with pride. They are the words of citizens, and they represent our greatest hope.

gather to inaugurate a president, we bear witness to the enduring strength of our Constitution. We affirm the promise of our democracy.

bear witness to the enduring strength of our Constitution. We affirm the promise of our democracy. We recall that what binds this nation together is not the colors of our skin or the tenets of our faith or the origins of our names.

affirm the promise of our democracy. We recall that what binds this nation together is not the colors of our skin or the tenets of our faith or the origins of our names.

recall that what binds this nation together is not the colors of our skin or the tenets of our faith or the origins of our names.

hold these truths to be self-evident, that all men are created equal, that they are endowed by their Creator with certain unalienable rights, that among these are Life, Liberty, and the pursuit of Happiness.

continue a never-ending journey, to bridge the meaning of those words with the realities of our time. For history tells us that while these truths may be self-evident, they have never been self-executing; that freedom is a gift from

learned that no nation founded on the principles of liberty and equality could survive half-slave and half-free. We made ourselves anew, and vowed to move forward together.

made ourselves anew, and vowed to move forward together. Together, we determined that a modern economy requires railroads and highways to speed travel and commerce; schools and colleges to train our workers.

determined that a modern economy requires railroads and highways to speed travel and commerce; schools and colleges to train our workers.

discovered that a free market only thrives when there are rules to ensure competition and fair play. Together, we resolved that a great nation must care for the vulnerable, and protect its people from life's worst hazards and misfortune.

resolved that a great nation must care for the vulnerable, and protect its people from life's worst hazards and misfortune.

succumbed to the fiction that all society's ills can be cured through government alone. Our celebration of initiative and enterprise; our insistence on hard work and personal responsibility; these are constant in our character.

that fidelity to our founding principles requires new responses to new challenges; that preserving our individual freedoms ultimately requires collective action.

It need to equip our children for the future, or build the roads and networks and research labs that will bring new jobs and businesses to our shores.

possess all the qualities that this world without boundaries demands: youth and drive; diversity and openness; an endless capacity for risk and a gift for reinvention.

size it together. For we, the people, understand that our country cannot succeed when a shrinking few do very well and a growing many barely make it.

believe that America's prosperity must rest upon the broad shoulders of a rising middle class. We know that America thrives when every person can find independence and pride in their work, when the wages of honest labor liberate families from the brink of hardship.

know that America thrives when every person can find independence and pride in their work, when the wages of honest labor liberate families from the brink of hardship.

understand that outworn programs are inadequate to the needs of our time. We must harness new ideas and technology to remake our government, rework our tax code, reform our schools, and empower our citizens with the skills they need to work harder, learn more, reach higher.

reject the belief that America must choose between caring for the generation that built this country and investing in the generation that will build its future.

remember the lessons of our past, when twilight years were spent in poverty, and parents of a child with a disability had nowhere to turn.

do not believe that in this country, freedom is reserved for the lucky, or happiness for the few. We recognize that no matter how responsibly we live our lives, any one of us, at any time, may face a job loss, or a recognition that no matter how responsibly we live our lives, any one of us, at any time, may face a job loss, or a sudden illness, or a home swept away in a terrible storm.

live our lives, any one of us, at any time, may face a job loss, or a sudden illness, or a home swept away in a terrible storm.

make to each other - through Medicare, and Medicaid, and Social Security - these things do not sap our initiative; they strengthen us.

live our lives, any one of us, at any time, may face a job loss, or a sudden illness, or a home swept away in a terrible storm.

face, but because engagement can more durably lift suspicion and fear. America will remain the anchor of strong alliances in every corner of the globe; and we will renew those institutions that extend our capacity to manage crisis ahead, for no one has a greater stake

68

## Quantifying the Race War in America

October 23, 2020 by Will Orichton

Most people see hate filtered through the individual events that reach our daily lives or news feeds - a chance encounter, a video, a story from a friend. But hate rarely emerges at random. There are patterns, sources, and cultures that give rise to hateful trends.

When our team at the podcast *Sounds Like Hate* was approached with a trove of secret recordings from the Base's vetting room for domestic terrorists in training, our goal was to understand the patterns in this data. You can hear the full story of our investigation by listening to our podcast series.

The Base is a terrorist organization that began in 2018 to advance a white supremacist agenda of the collapse of America, an impending race war, and preparation for violence. Each recording contains a vetting call where members of the Base talk to potential recruits over the messaging app Wire. We were faced with the task of analyzing a significant amount of data, 83 hours in total. In this post, we will describe how we applied statistical analysis, data visualization, and machine learning to understand the trends beneath the hate.

we → are

definitely → being

very →

to → a → certain → extent

remain → anonymous

targeted → by → the → system

antifa → media → feds

for → infiltration



69

## Glimpses of structure

---

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

**For example**

Lexical: <A> at <B>

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

---

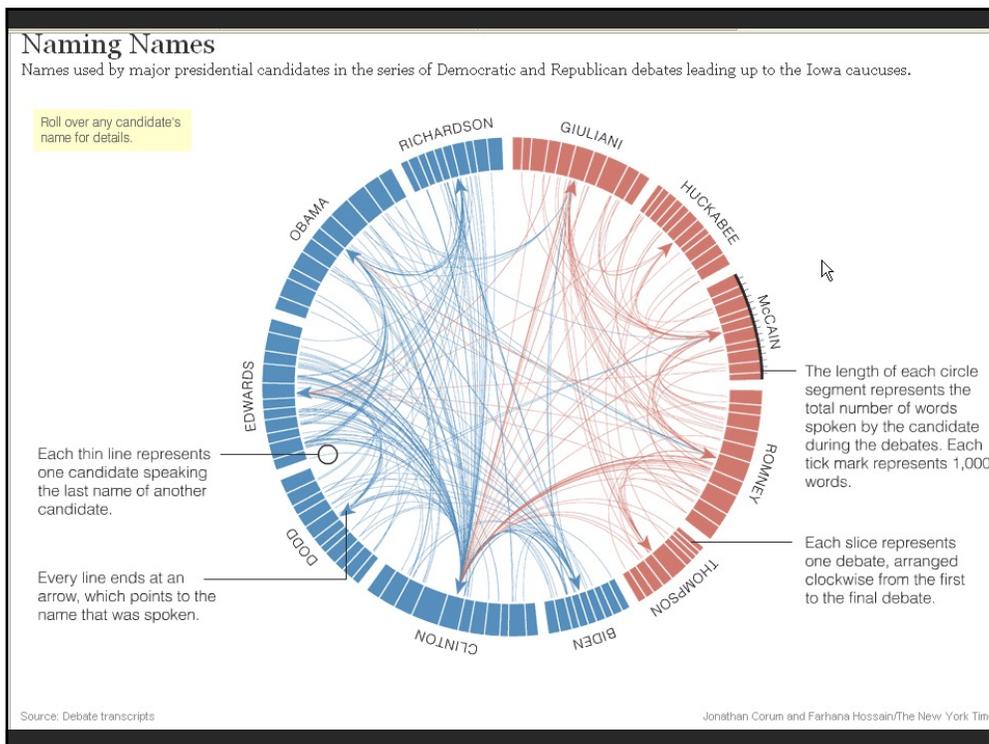
### Many dimensions to consider:

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

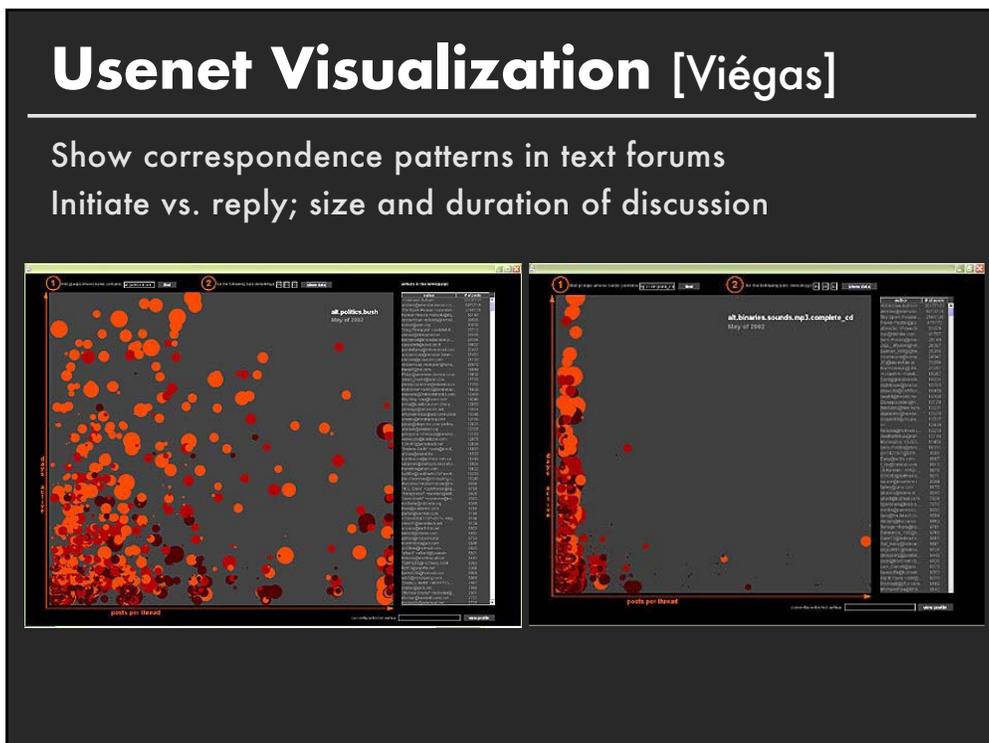
### Interesting cross-products:

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

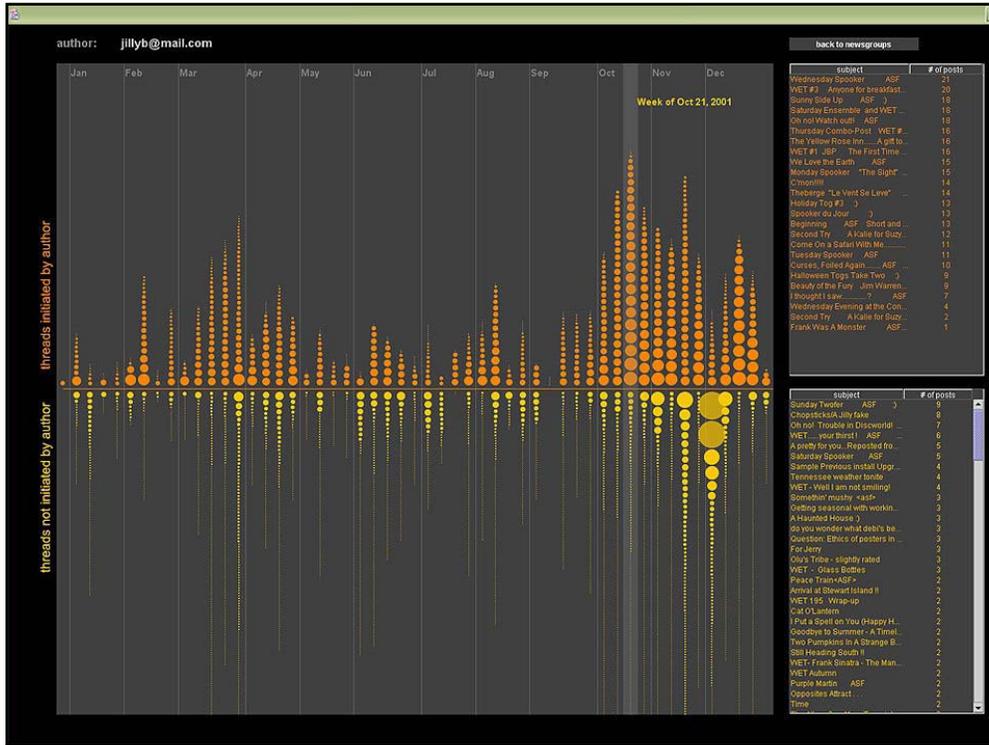
92



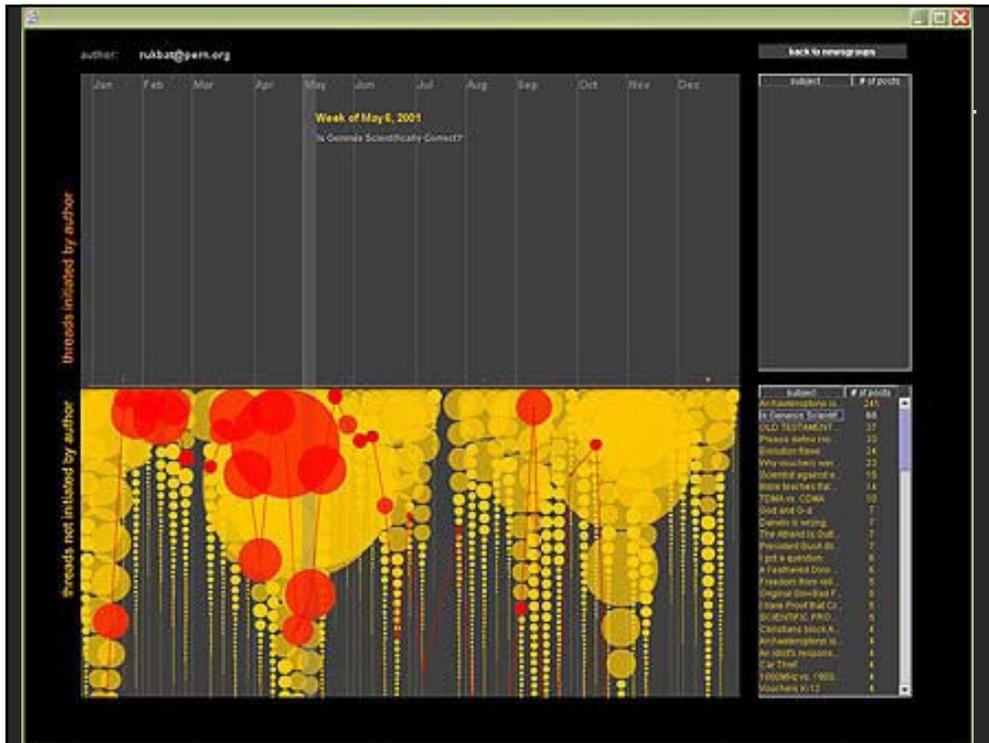
93



96

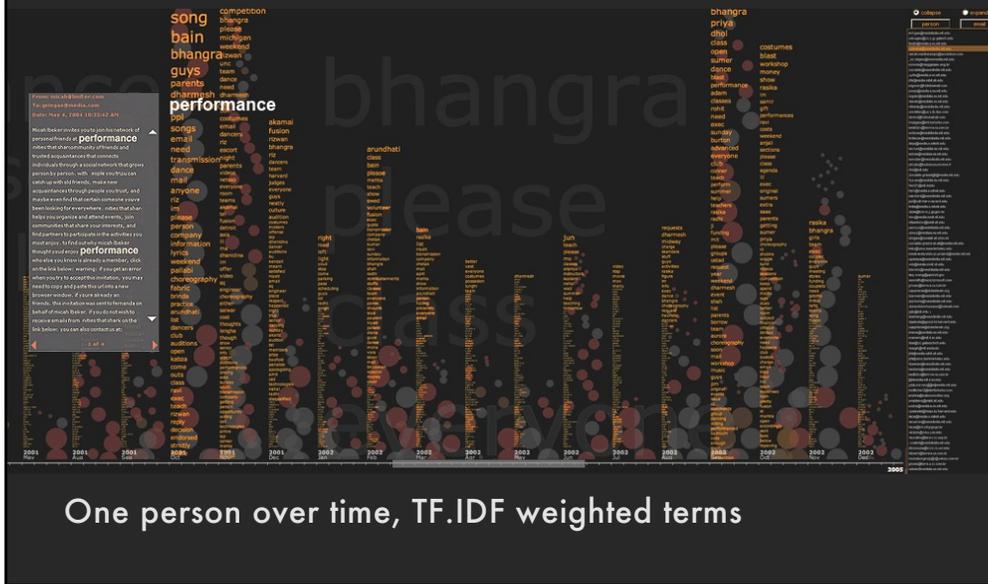


97



98

# Themail (Viégas)



One person over time, TF.IDF weighted terms

100

# Visualizing Document Collections

104

# Topic modeling

## Topic modeling approaches

Assume documents are a mixture of topics

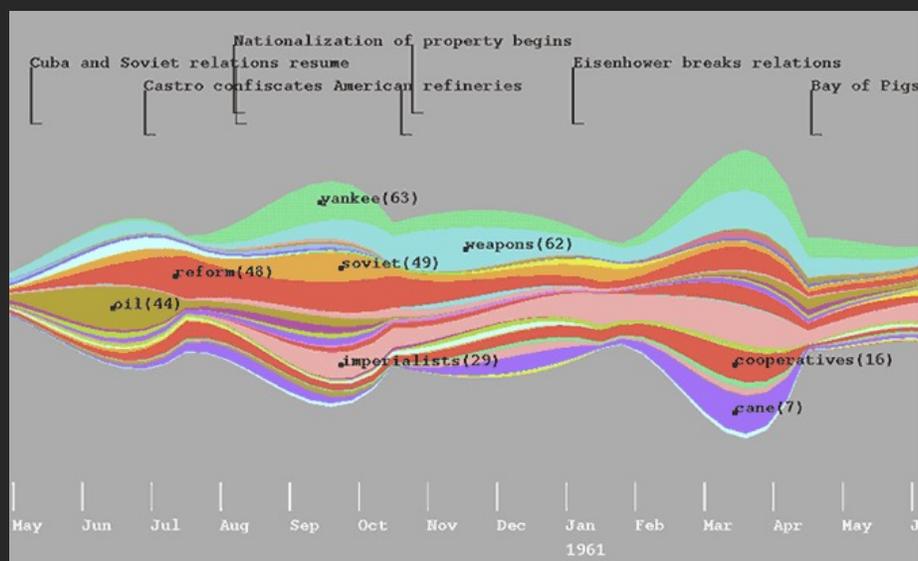
Topics are (roughly) a set of co-occurring terms

Latent Semantic Analysis (LSA): reduce term matrix

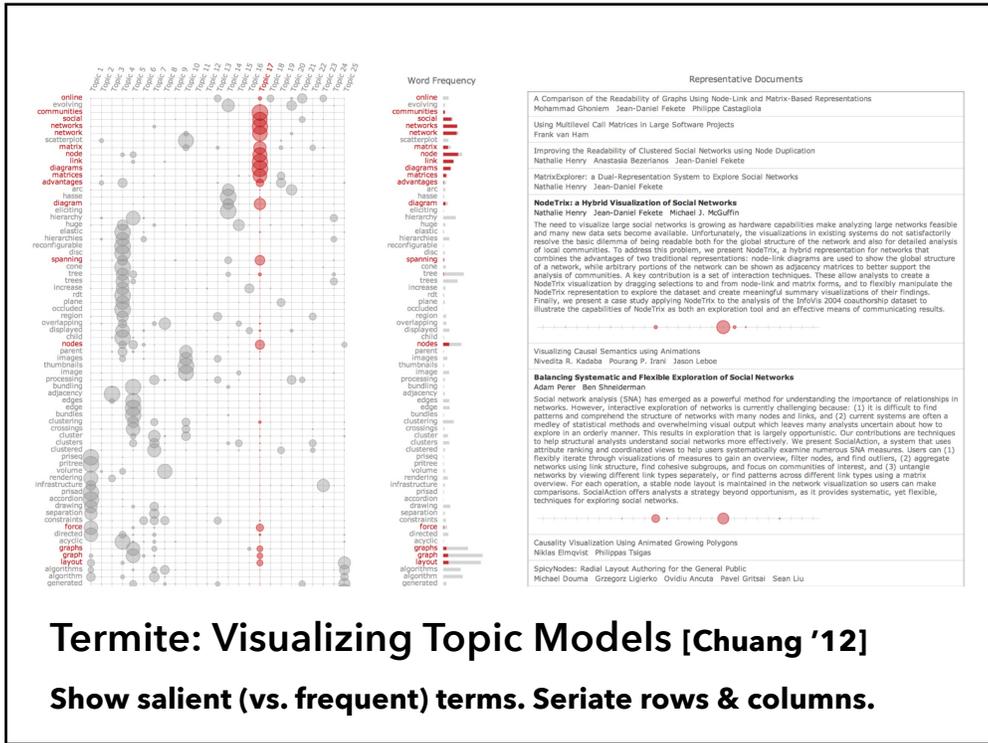
Latent Dirichlet Allocation (LDA): statistical model

113

# ThemeRiver (Havre et al 99)



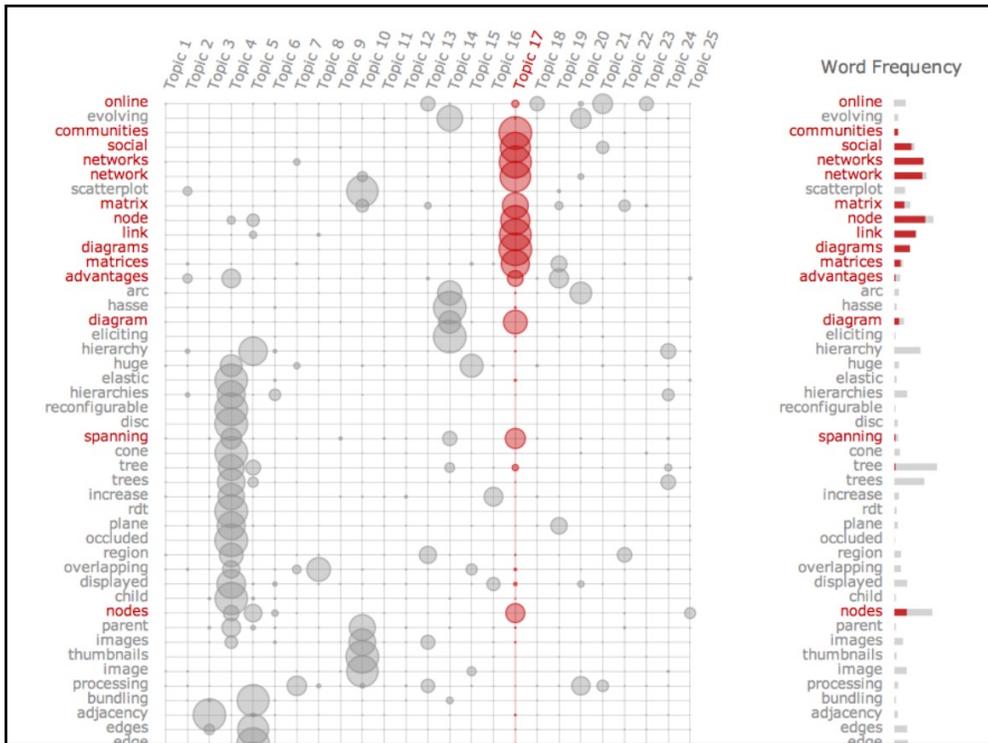
114



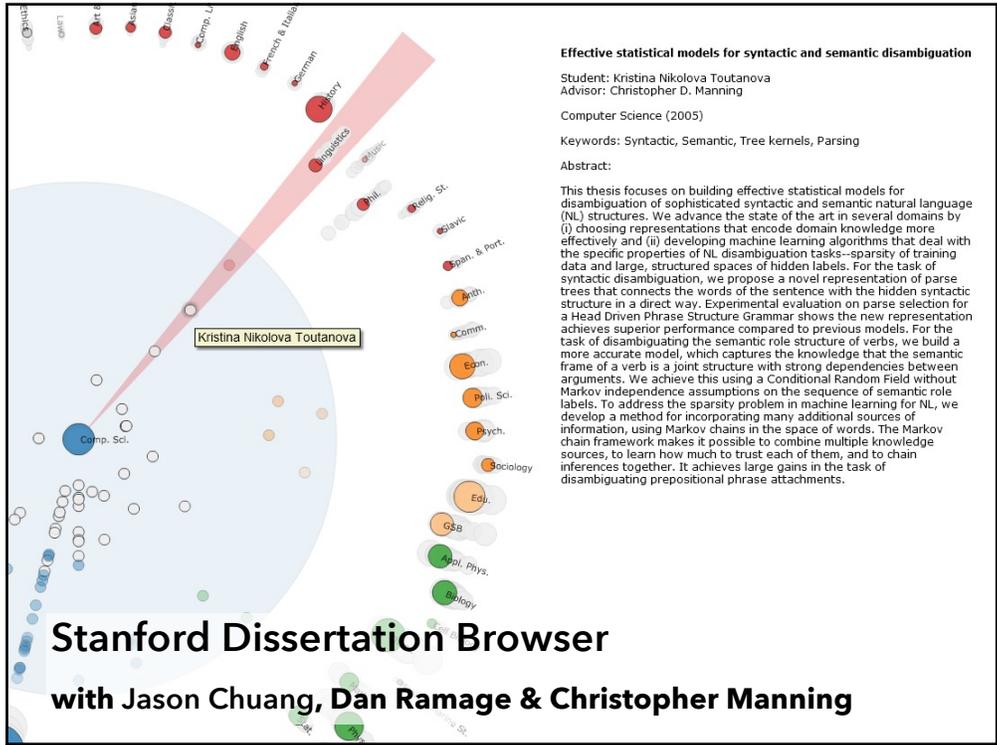
## Termite: Visualizing Topic Models [Chuang '12]

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

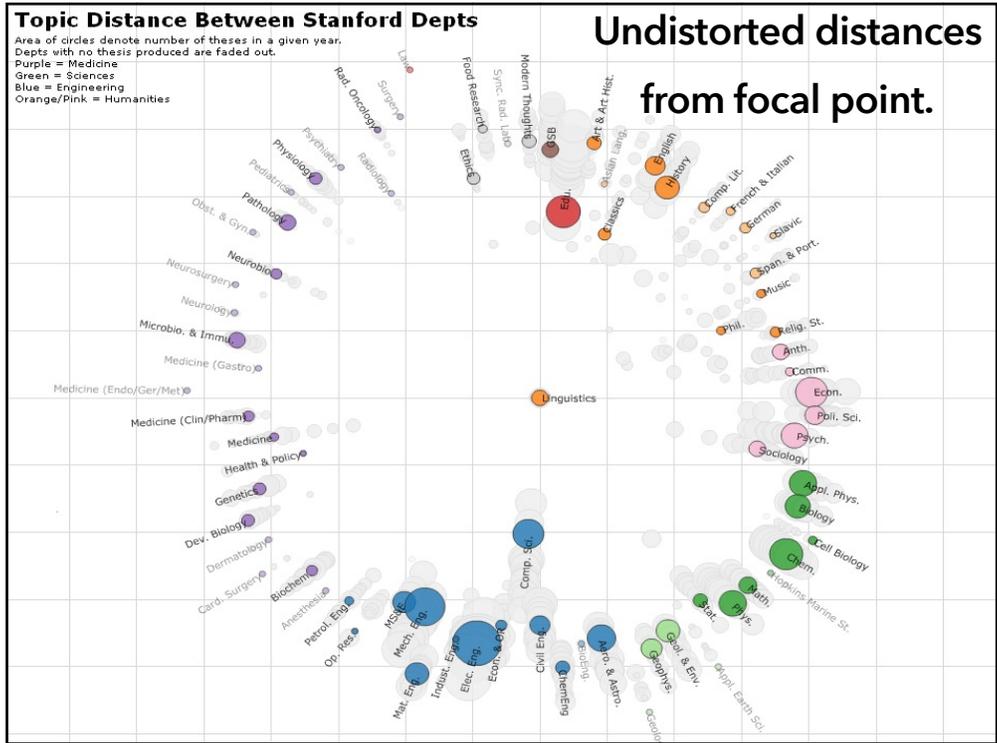
122



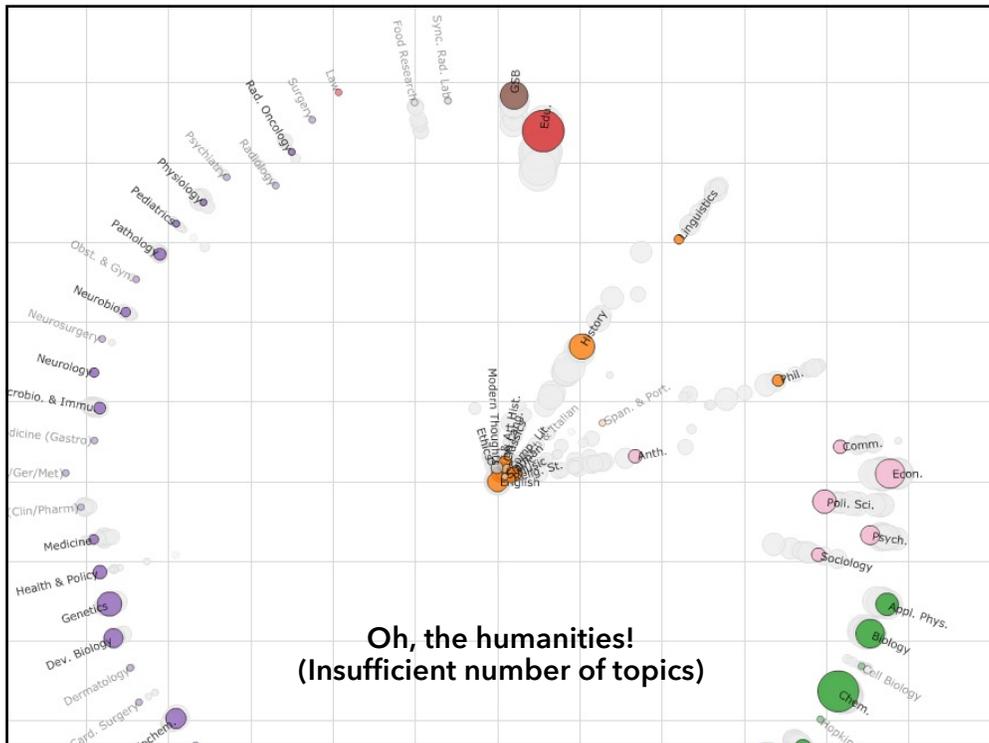
123



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.

134