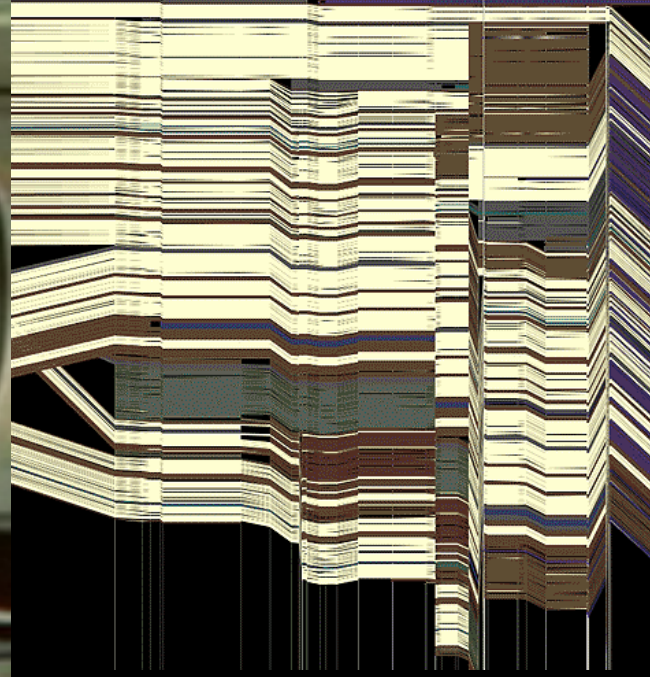
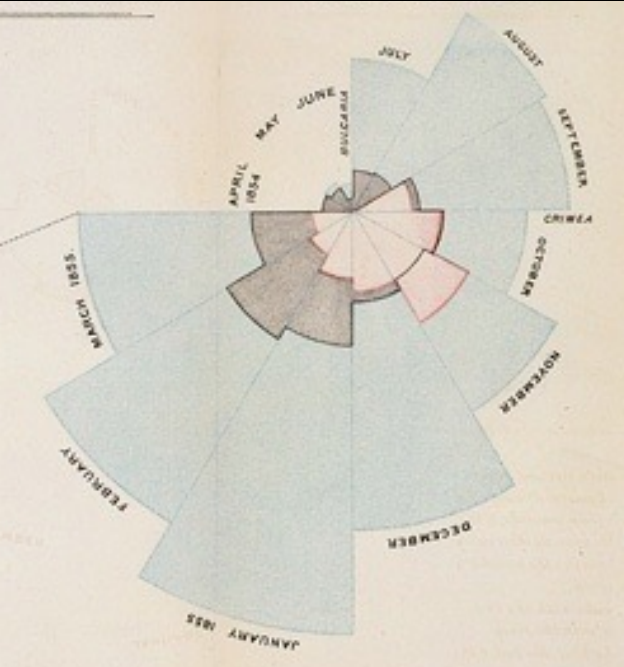


CSE 512 - Data Visualization

Text Visualization



Jeffrey Heer University of Washington

Why Visualize Text?

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

Text as Data

Documents

Articles, books and novels

E-mails, web pages, blogs

Tags, comments

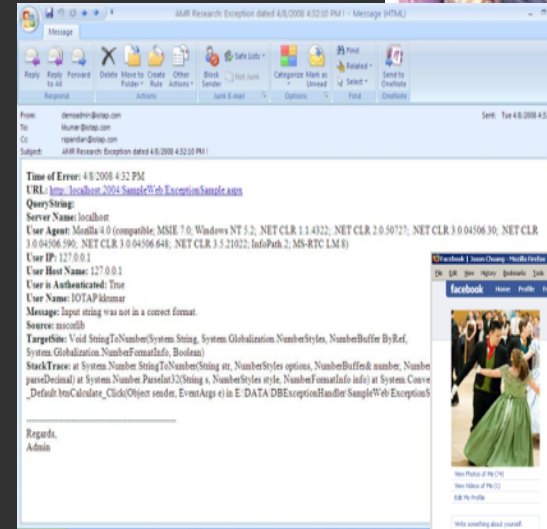
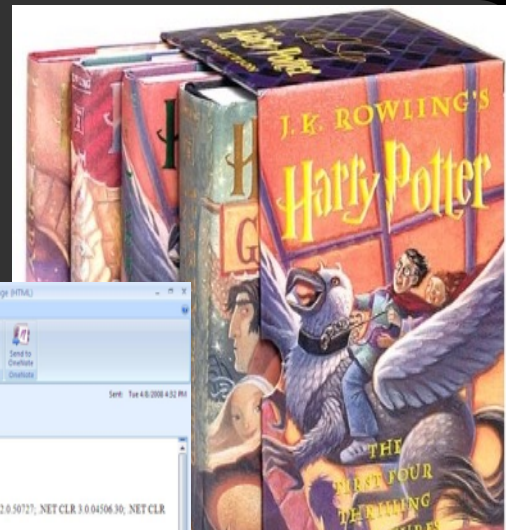
Computer programs, logs

Collections of Documents

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

Social networks (personal profiles)

Academic collaborations (publications)



Example:
Health Care Reform

Example: Health Care Reform

Recent History

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?

A Concrete Example

September 10, 2009

TEXT

Obama's Health Care Speech to Congress

Following is the prepared text of President Obama's speech to Congress on the need to overhaul health care in the United States, as released by the White House.

Madame Speaker, Vice President Biden, Members of Congress, and the American people:

When I spoke here last winter, this nation was facing the worst economic crisis since the Great Depression. We were losing an average of 700,000 jobs per month. Credit was frozen. And our financial system was on the verge of collapse.

As any American who is still looking for work or a way to pay their bills will tell you, we are by no means out of the woods. A full and vibrant recovery is many months away. And I will not let up until those Americans who seek jobs can find them; until those businesses that seek capital and credit can thrive; until all responsible homeowners can stay in their homes. That is our ultimate goal. But thanks to the bold and decisive action we have taken since January, I can stand here with confidence and say that we have pulled this economy back from the brink.

I want to thank the members of this body for your efforts and your support in these last several months, and especially those who have taken the difficult votes that have put us on a path to recovery. I also want to thank the American people for their patience and resolve during this trying time for our nation.

But we did not come here just to clean up crises. We came to build a future. So tonight, I return to speak to all of you

Word Tree: Word Sequences

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

Search

il

Back

Forward

Start

End

Occurrence Order

Clicks Will Zoom

52 hits



Search **i will**

Back

Forward

Start End

Occurrence Order

Clicks Will Zoom

12 hits

i will

not

let up until those americans who seek jobs can find them - - (applause) - - until those

back down on the basic principle that if americans can't find affordable coverage , w

sign

a plan that adds one dime to our deficits - - either now or in the future

it if it adds one dime to the deficit , now or in the future , period .

make that same mistake with health care .

waste time with those who have made the calculation that it's better politics to kill thi

- - and i will not accept the status quo as a solution .

accept the status quo as a solution .

make sure that no government bureaucrat or insurance company bureaucrat gets between you and t

protect medicare .

continue to seek common ground in the weeks ahead .

be there to listen .

still believe

we can

act even when it's hard .

replace economy with civility , and gridlock with progress .

do great things , and that here and now we will meet history's test .

- - i still believe that we can act when it's hard .

that we can act when it's hard .

Gulfs of Evaluation

Many text visualizations do not represent the 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?

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.

Topics

Text as Data

Visualizing Document Content

Evolving Documents

Visualizing Conversation

Document Collections

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 Processing Pipeline

1. Tokenization

Segment text into terms.

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

Numbers and symbols? *#gocard, @stanfordfball, 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

Tokenization & Stemming

Well-formed text to support stemming?

txt u l8r!

Word meaning or entities?

#berkeley -> #berkelei

Reverse stems for presentation.

Ha appl made programm cool?

Has Apple made programmers cool?

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

	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

WordCounts (Harris '04)

The screenshot displays the WordCounts website interface. At the top right, there is a button labeled "WORDCOUNT". Below this, a navigation bar contains "PREVIOUS WORD" with a left arrow and "NEXT WORD" with a right arrow. The main content area shows a horizontal bar chart representing word frequencies. The word "the" is the largest and is labeled with a red "1" below it. Other words are smaller and labeled with red numbers: "of" (2), "and" (3), "to" (4), "ain" (5), "that" (6), "is" (7), "was" (8), "i" (9), "for" (10), "on" (11), "you" (12), "he" (13), "be" (14), "with" (15), "by" (16), "have" (17), "the" (18), "but" (19), "ed" (20), "is" (21), "from" (22), "which" (23), "ow" (24), "are" (25), "em" (26), "a" (27), "an" (28), "in" (29), "at" (30), "on" (31), "the" (32), "of" (33), "to" (34), "and" (35), "for" (36), "with" (37), "by" (38), "have" (39), "the" (40), "but" (41), "ed" (42), "is" (43), "from" (44), "which" (45), "ow" (46), "are" (47), "em" (48), "a" (49), "an" (50). Below the chart, there is a "CURRENT WORD" label. At the bottom, there are search controls: "FIND WORD:" followed by an input field and a right arrow, "BY RANK:" followed by an input field and a right arrow, "REQUESTED WORD: THE", and "RANK: 1". On the right side, it says "86800 WORDS IN ARCHIVE" and "ABOUT WORDCOUNT".

<http://wordcount.org>

Tag 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

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: $tf_{td} / \sum_t tf_{td}$

Keyword Weighting

Term Frequency

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

TF.IDF: Term Freq by Inverse Document Freq

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

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

Keyword Weighting

Term Frequency

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

TF.IDF: Term Freq by Inverse Document Freq

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

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

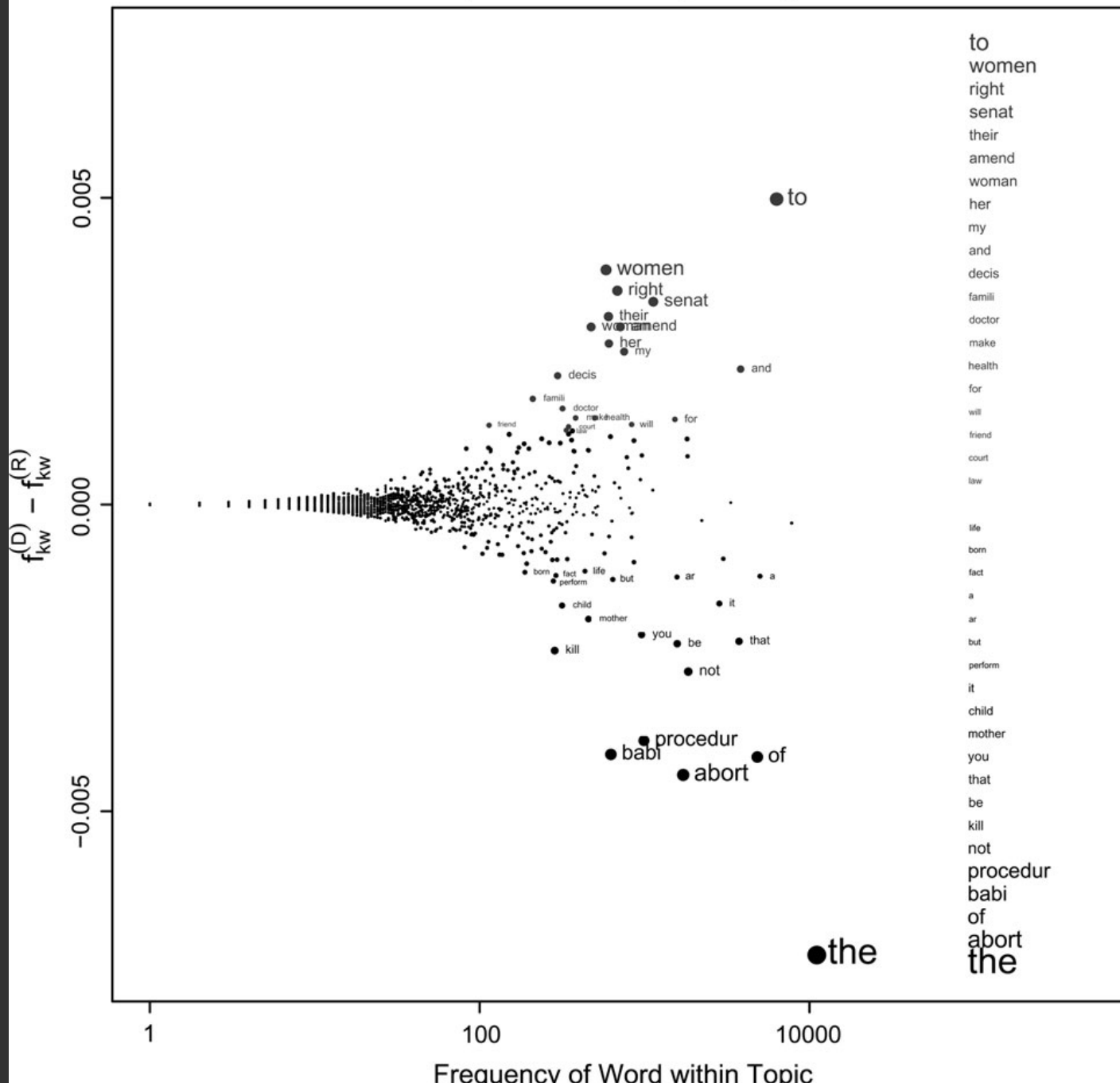
G^2 : Probability of different word frequency

$$E_1 = |d| \times (tf_{td} + tf_{t(C-d)}) / |C|$$

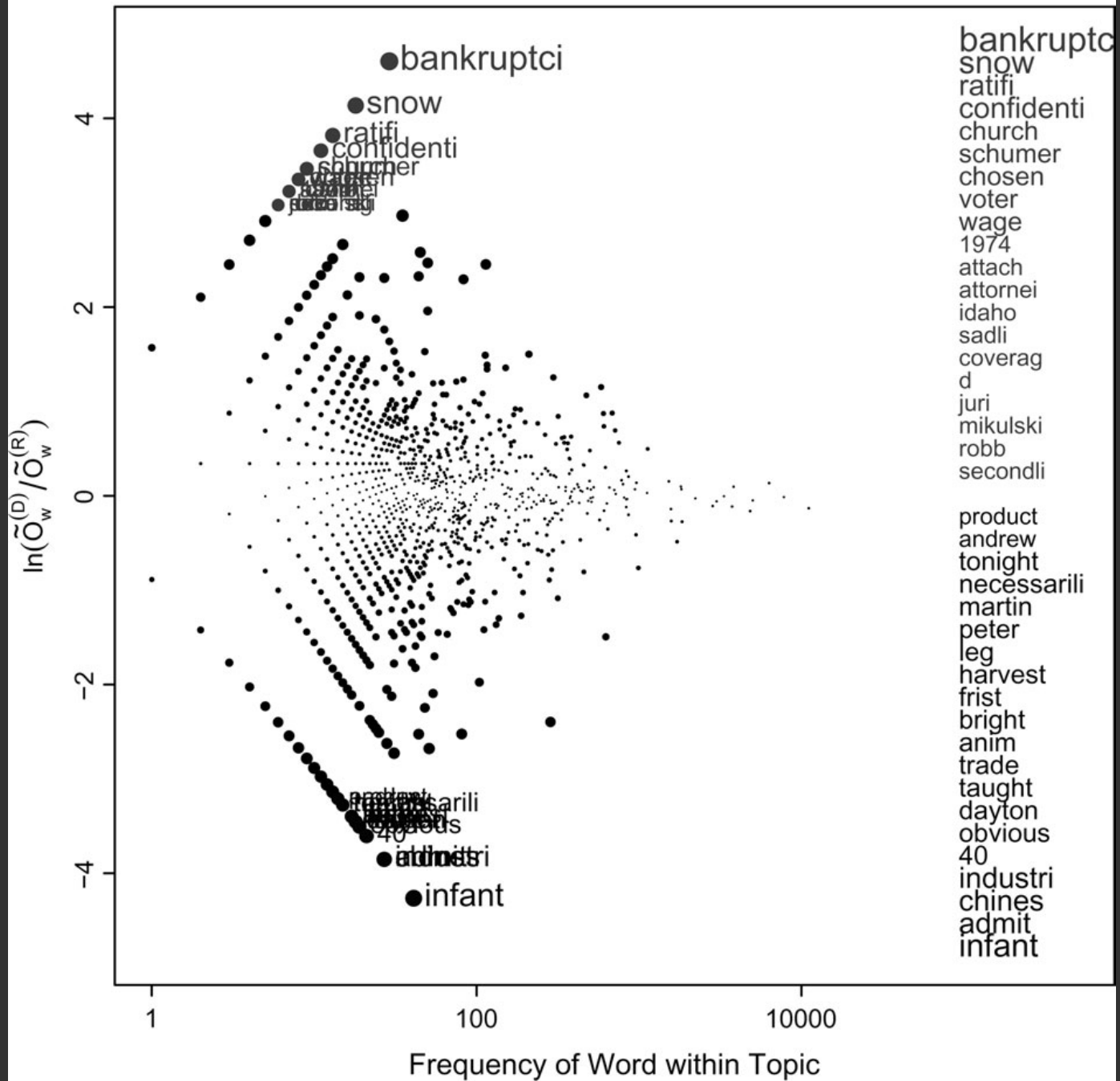
$$E_2 = |C-d| \times (tf_{td} + tf_{t(C-d)}) / |C|$$

$$G^2 = 2 \times (tf_{td} \log(tf_{td}/E_1) + tf_{t(C-d)} \log(tf_{t(C-d)}/E_2))$$

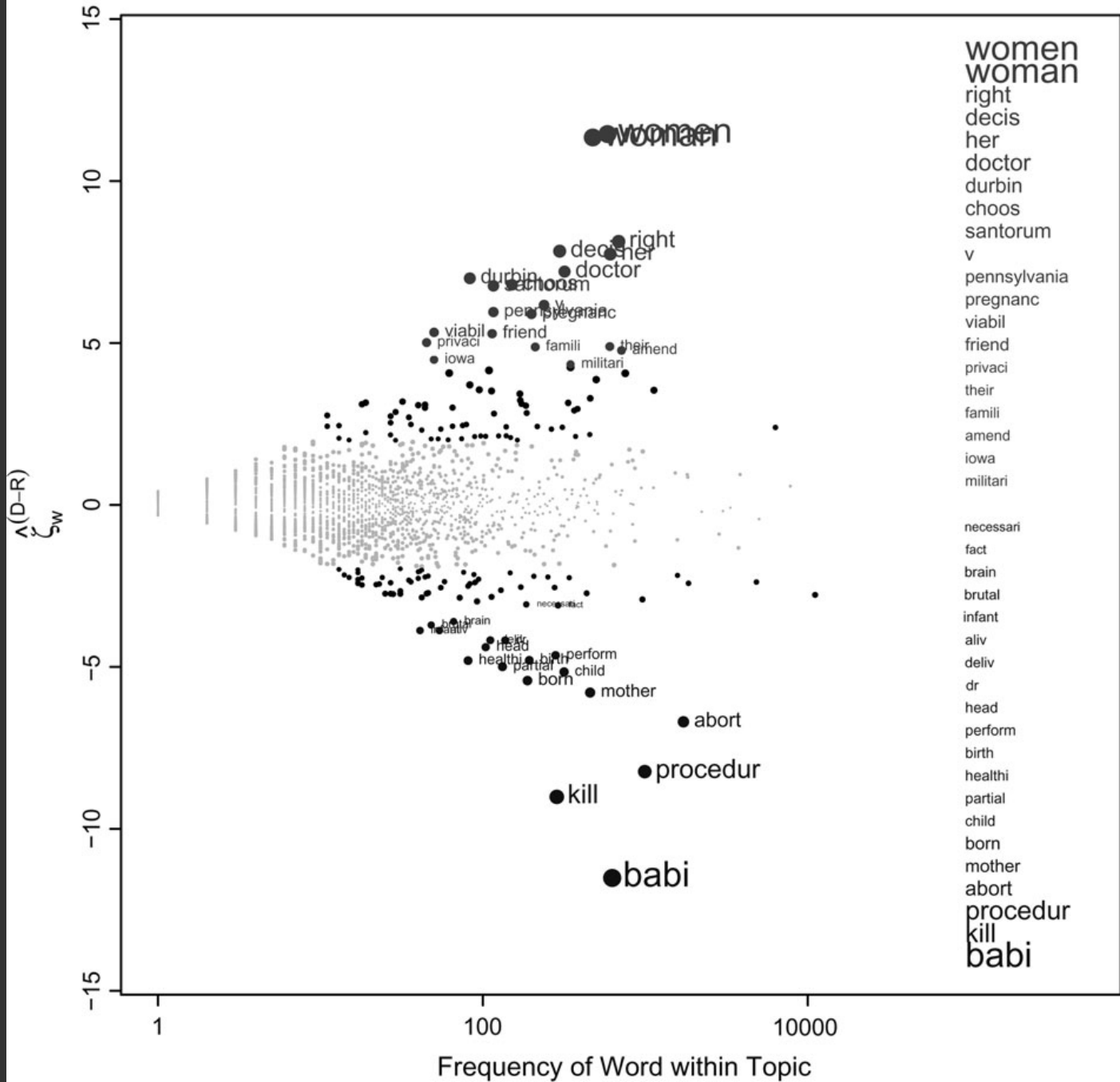
Partisan Words, 106th Congress, Abortion (Difference of Proportions)



Partisan Words, 106th Congress, Abortion
 (Log-Odds-Ratio, Smoothed Log-Odds-Ratio)



Partisan Words, 106th Congress, Abortion (Weighted Log-Odds-Ratio, Informative Dirichlet Prior)



Limitations of Freq. Statistics

Typically focus on unigrams (single terms)

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

Not clear that these provide best description

A “bag of words” ignores additional information

Grammar / part-of-speech

Position within document

Recognizable entities

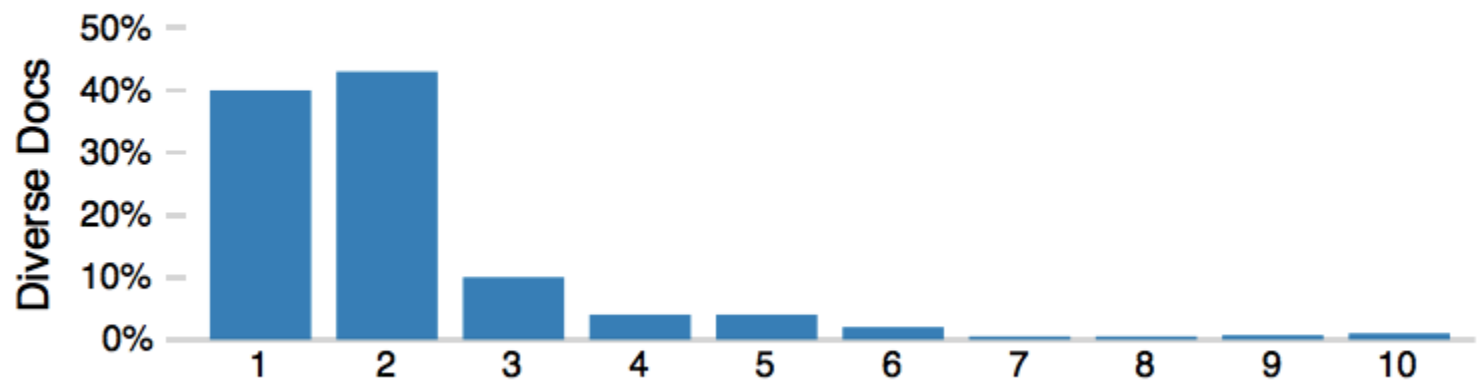
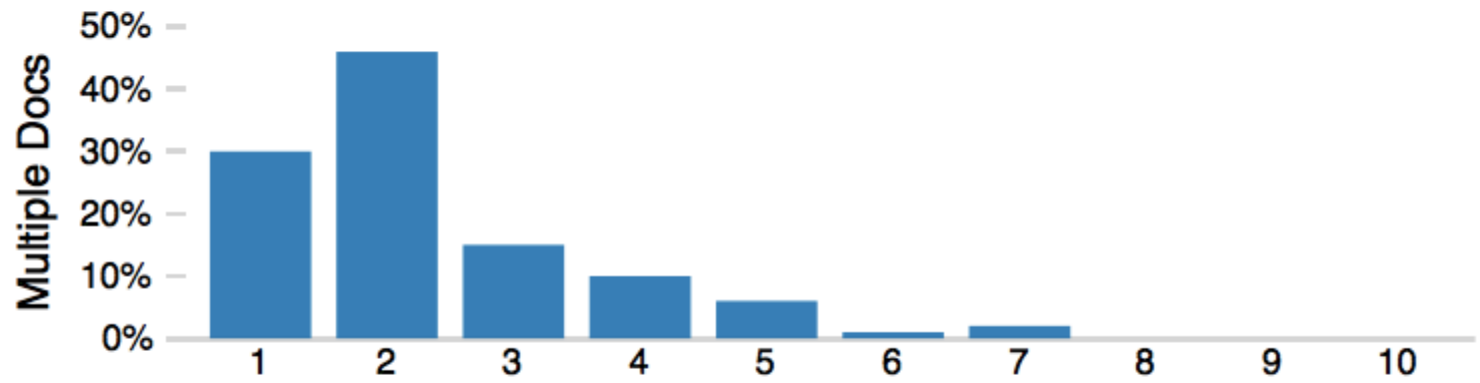
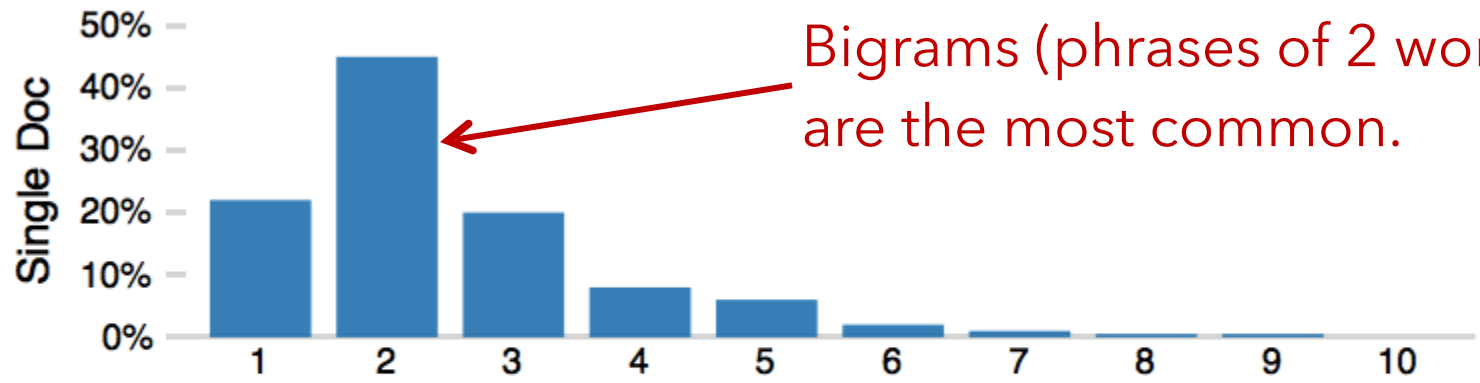
How do people describe text?

We asked 69 subjects (graduate students) to read and describe dissertation abstracts.

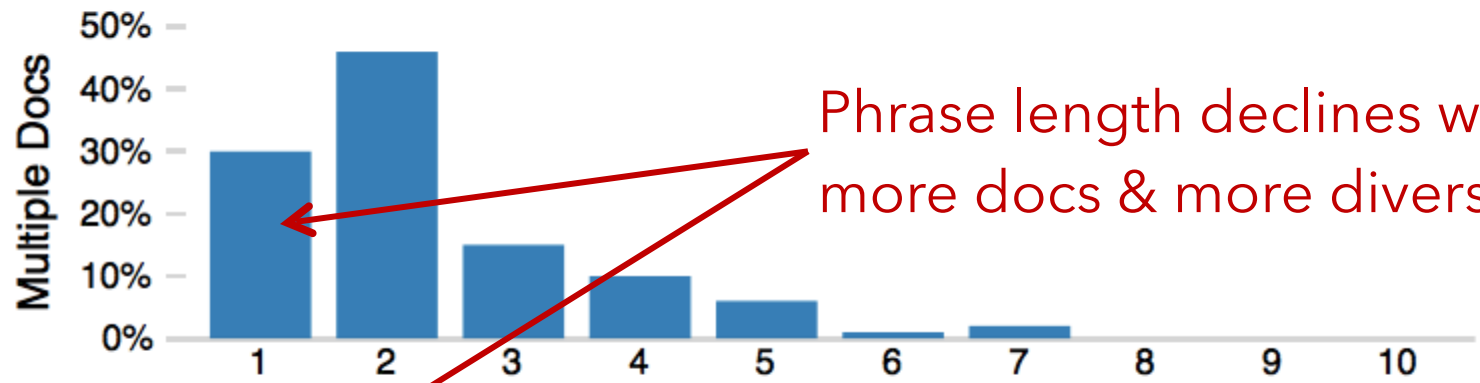
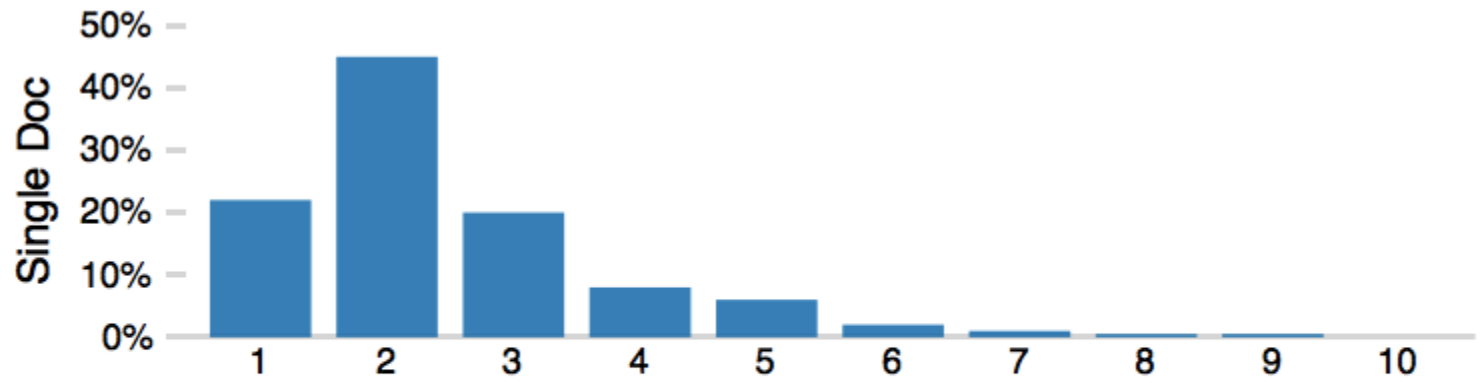
Students were given 3 documents in sequence; they then described the collection as a whole.

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

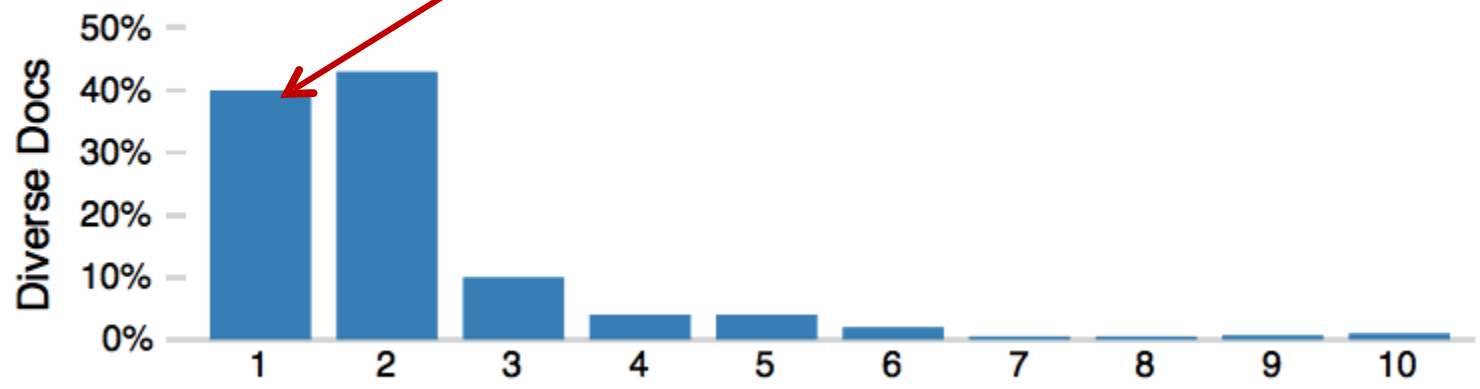
[Chuang, Manning & Heer, 2012]



Phrase Length



Phrase length declines with more docs & more diversity.



Phrase Length

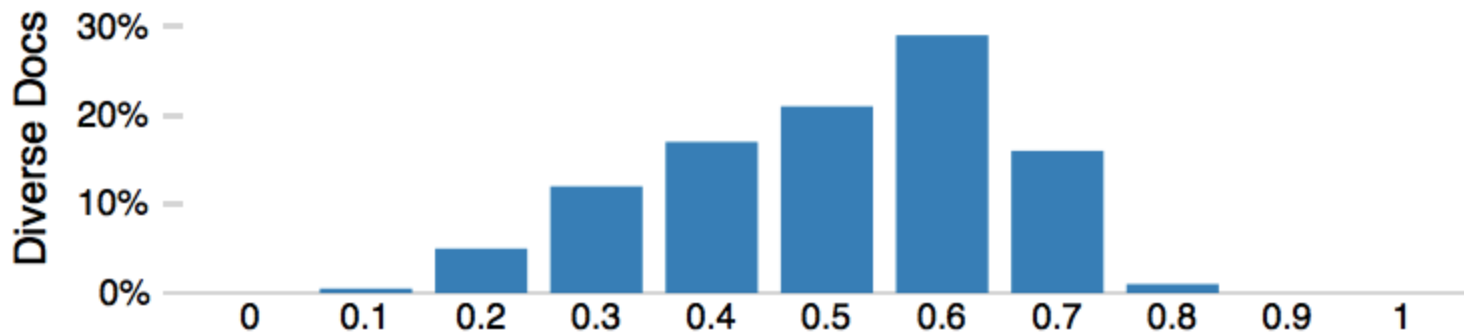
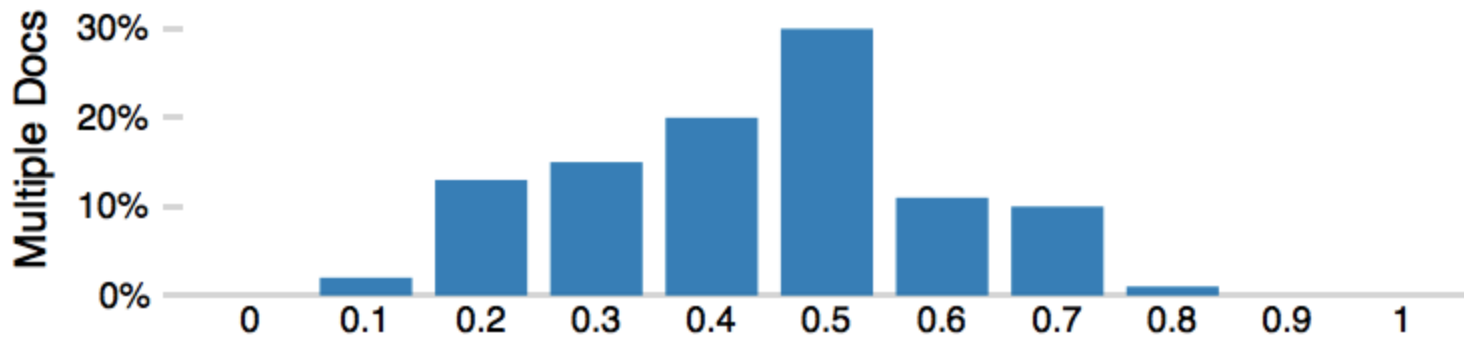
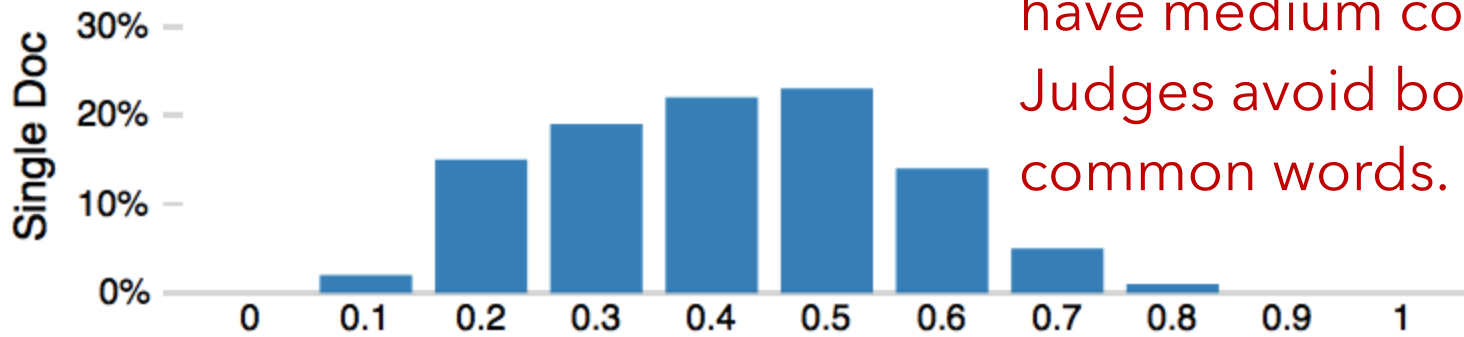
Term Commonness

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

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

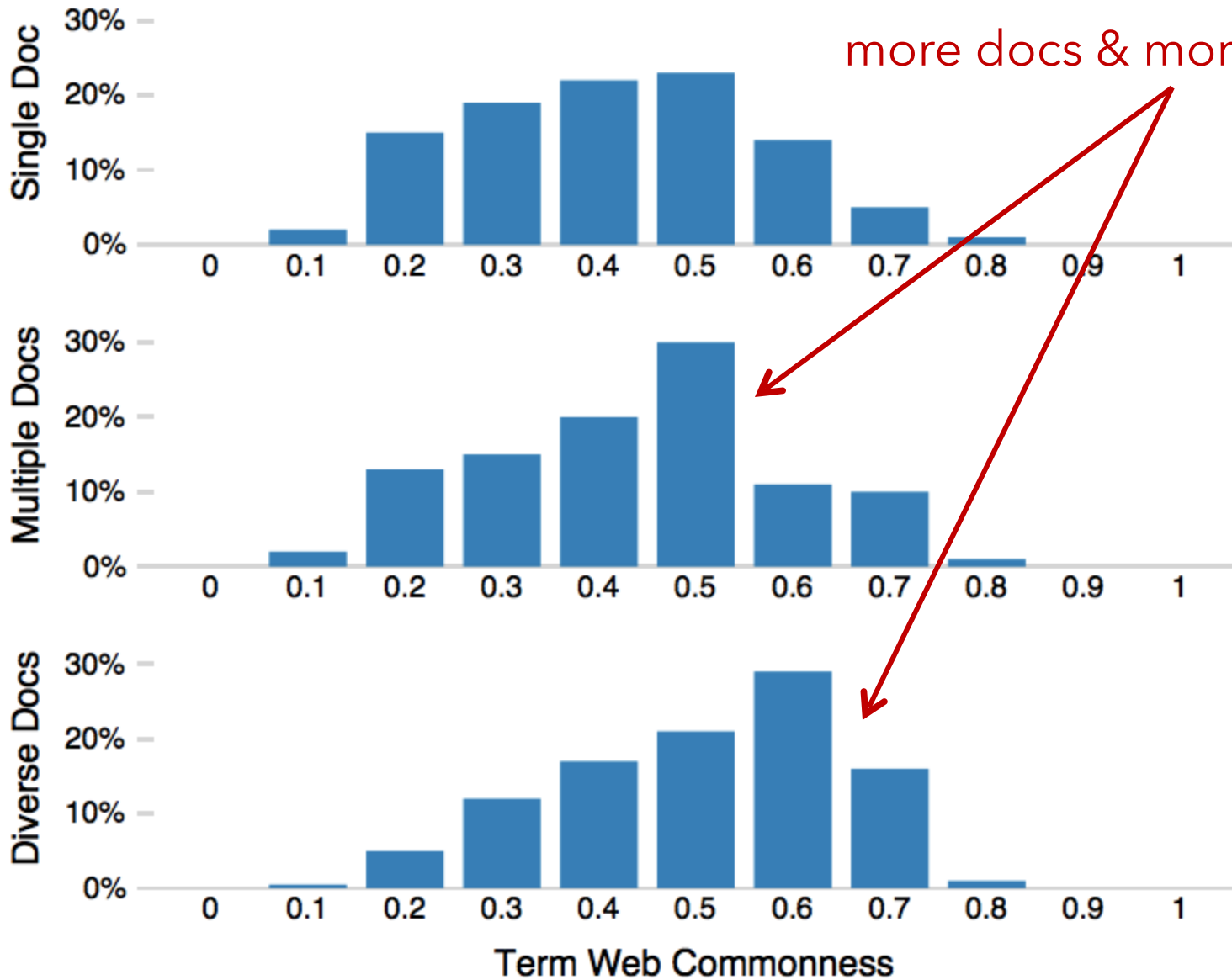
Measured across a corpus or across the entire English language (using Google n-grams)

Selected descriptive terms have medium commonness. Judges avoid both rare and common words.

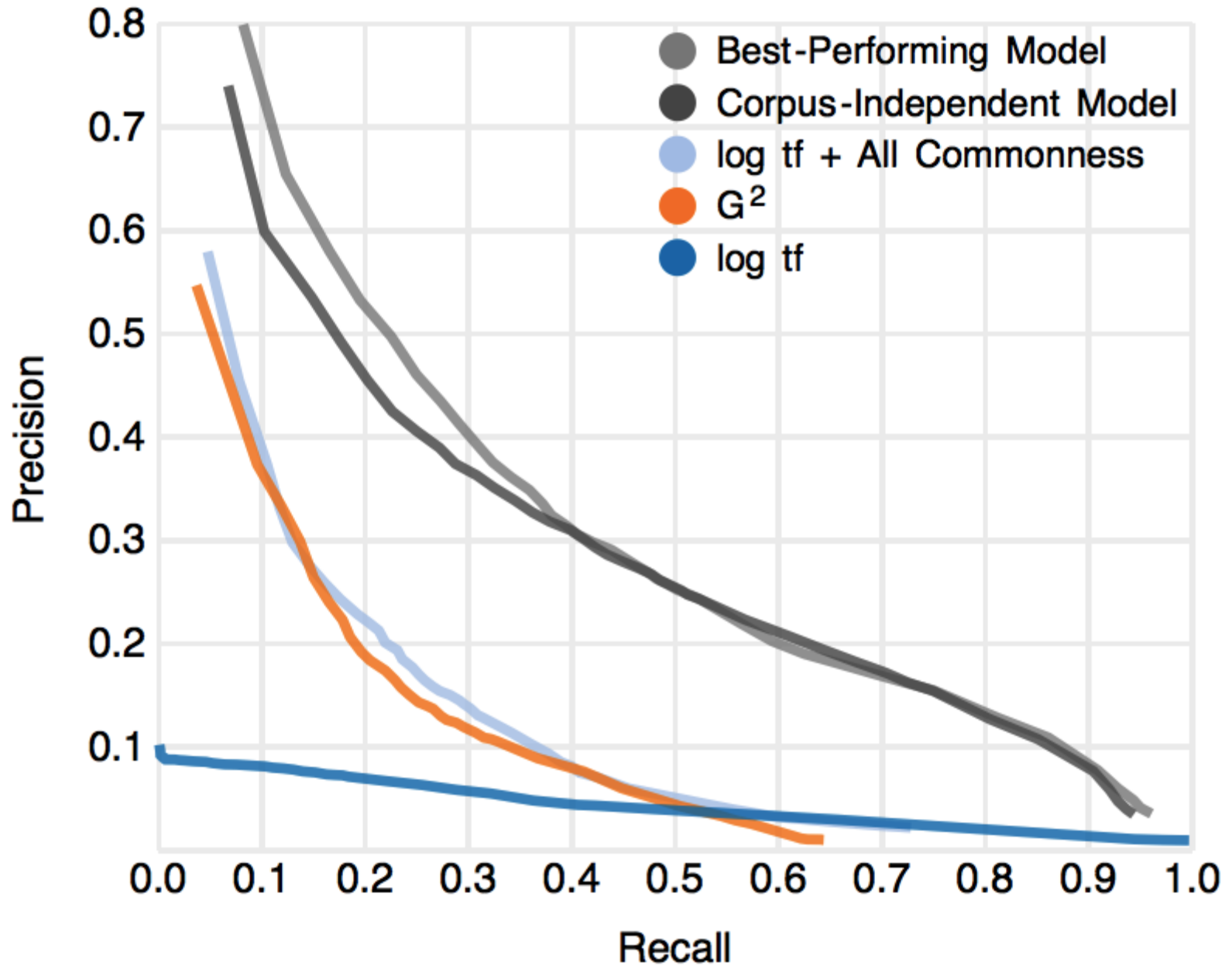


Term Web Commonness

Commonness increases with more docs & more diversity.



Scoring Terms with Freq, Grammar & Position





A fighter jet rain check

Story and video by [Chamila Jayaweera](#)

Have you ever thought about what it takes to make sure that sea-based fighter jets stay dry?

When it comes to the F/A-18 Super Hornet, Boeing engineers in St. Louis use a special process called the Water Check Test to rule out areas where moisture could seep into the aircraft and its electronics suite.

Program experts douse the jet with simulated rain at a 15-inch-per-hour rate for about 20 minutes inside an enormous hangar in St. Louis.

"Our ultimate customers are U.S. Navy fighter pilots, and we want to ensure their safety in flight and on the ground, and water-tight integrity of the aircraft also helps increase their effectiveness," said Boeing's Rich Baxter, F/A-18 Super Hornet final assembly manager.

To find out more about how the process works and watch the action unfold, click above to see the video story.



CHAMILA JAYAWEERA/BOEING

The Water Check team rolls in a large metal frame, which they affectionately call their "spray tree," over a Super Hornet inside a St. Louis hangar.



G²

Regression Model

fighter

F/A

Hornet

Super

Boeing

-18

rain

St.

jet

Louis

15-inch-per-hour

douse

hangar

water-tight

Check

Baxter

sea-based

aircraft

Rich

seep

click

Navy

sure

Water

moisture

watch

enormous

stay

want

Super Hornet

F/A -18

fighter jet

Boeing engineers

special process

rain check

electronics suite

Program experts

simulated rain

ultimate customers

enormous hangar

water-tight integrity

Rich Baxter

15-inch-per-hour rate

video story

aircraft

U.S. Navy fighter pilots

Super Hornet final assembly manager

U.S.
Navy fighter
fighter pilot
sea-based fighter

Yelp Review Spotlight (Yatani 2011)

'09 amazing around baked bar bass best chef delicious eat

elite e

hawaii

night

expe

sake

tabl

b) best sf
baked sea bass best sushi sure in striped bass
other person
fresh fish slow service sushi bar
sushi chef baked mussel more hour
only thing
long wait long time sushi restaurant good food
long line hawaiian roll reasonable price
baked mango
small place delicious everything

Mentioned 63 times

possess sage of the halos wisdom , and know in advance sushi zone only accepts cash and the waits will be **long** and arduous .

yes , its a **long** wait , learn the master of zen if you want to eat here .

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 interface. The search bar contains the text 'acronym resolution'. Below the search bar, there are filters for 'Articles and patents', 'anytime', and 'include citations'. The search results are displayed in a list format, with each entry including a title, authors, publication details, a brief abstract, and citation information. The first result is 'A supervised learning approach to acronym identification' by D Nadeau and P Turney, published in 2005. The second result is 'Biomedical term mapping databases' by JD Wren, JT Chang, and J Pustejovsky, published in 2005. The third result is 'Anthropogenic climate change over the Mediterranean region simulated by a global variable resolution model' by AL Gibelin, published in 2003. The fourth result is 'Metaphrase: an aid to the clinical conceptualization and formalization of patient problems in healthcare enterprises' by MS Tuttle, NE Olson, KD Keck, and WG Cole, published in 1998.

Google scholar [Advanced Scholar Search](#)

Scholar

[A supervised learning approach to acronym identification](#) [\[PDF\] from nrc-cnrc.ca](#)
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](#) [\[HTML\] from nih.gov](#)
JD Wren, JT Chang, J Pustejovsky... - *Nucleic acids ...*, 2005 - Oxford Univ Press
... the prevalence of polynoms, 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](#) [Find it@Stanford](#)
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

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

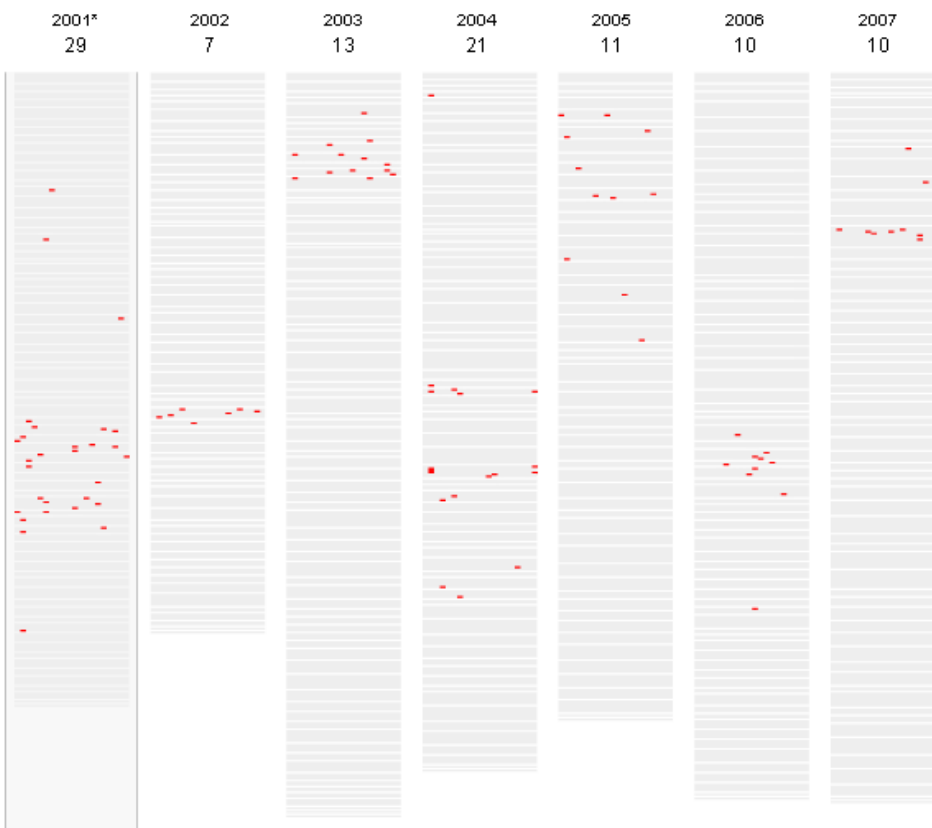
The screenshot shows the TileBars software interface. At the top, there is a 'User Query' section with three lines of input: 'osteoporosis', 'prevention', and 'research'. To the right of the query input are three buttons: 'Run Search', 'New Query', and 'Quit'. Below these are two rows of search controls: 'Search Limit' with a range from 50 to 1000 (250 is selected) and 'Number of Clusters' with a range from 3 to 10 (5 is selected). Below the search controls is a 'Mode' dropdown set to 'TileBars'. Underneath are three buttons: 'Cluster', 'Titles', and 'Backup'. The main display area is divided into two panes. The left pane, labeled 'Cluster', shows a vertical list of 12 horizontal bars representing different clusters. Each bar has a small colored square on the left (pink, yellow, or purple) and a grayscale heatmap. The right pane, labeled 'Titles', shows a list of search results corresponding to the clusters. The results include document IDs (e.g., FR88513-0157, FR88120-0046) and titles (e.g., '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', '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...').

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...

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.

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



The word in context

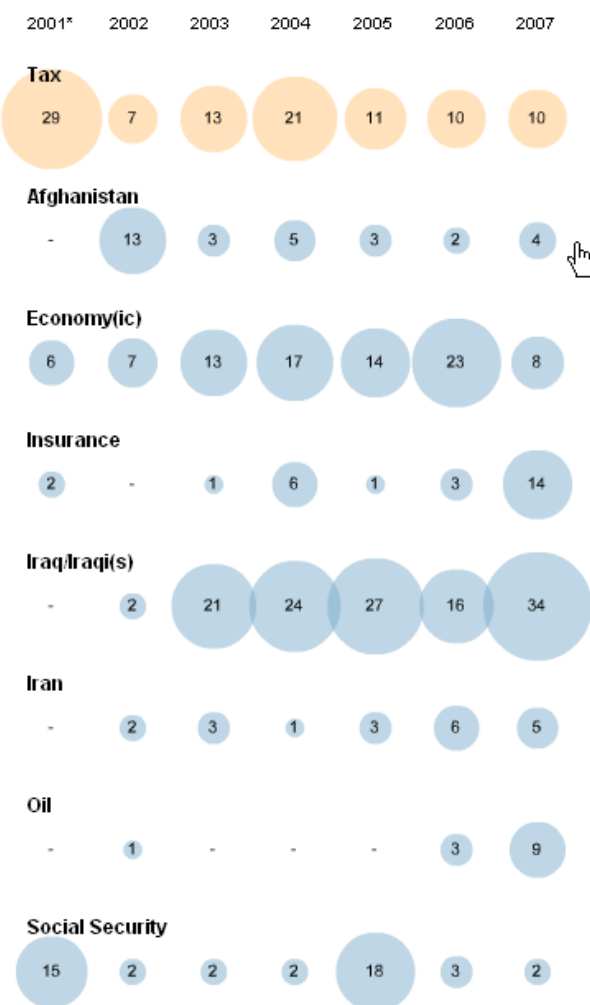
[Next Instance of 'Tax'](#)

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)

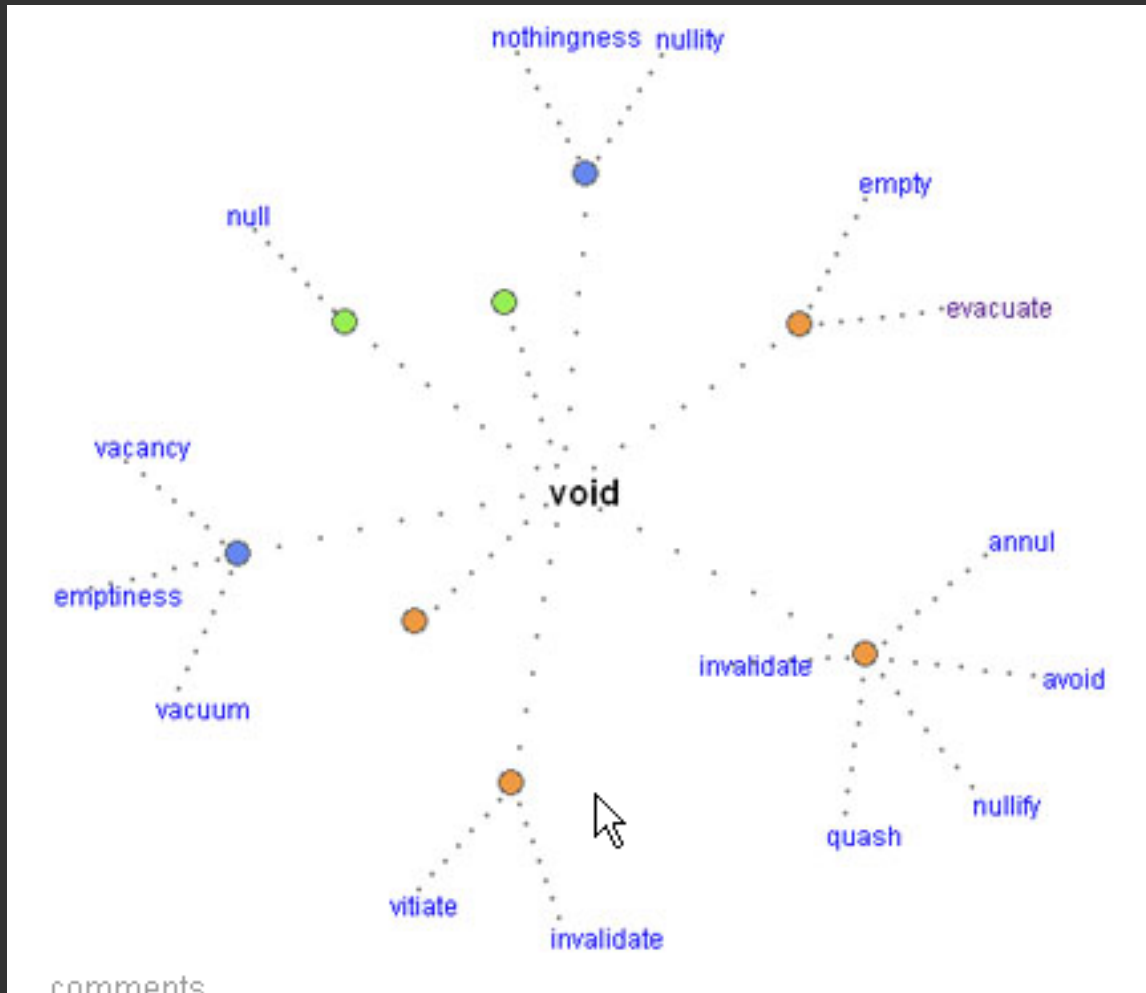
New York Times

Compared with other words



* 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.

Visual Thesaurus [ThinkMap]



Concordance

What is the common local context of a term?

The screenshot shows the 'Concordance - Larkin.Concordance' window. The interface includes a menu bar (File, Text, Search, Edit, Headwords, Contexts, View, Tools, Help), a toolbar with icons for file operations and text formatting, and a main display area with a table of concordance results. The table has columns for 'Headword', 'No.', 'Context...', 'Word', '...Context', and 'Reference'. The 'HEART' headword is highlighted in blue, with a count of 25. The table shows various contexts for 'heart', such as 'That my own heart drifts and cries, having no...', 'By the shout of the heart continually at work', and 'Nothing to adapt the skill of the heart to, skill'. The interface also features a vertical sidebar on the right with alignment options (Centered, Left-aligned, Index, None) and a status bar at the bottom showing statistics: Words (7318), Tokens (37070), At word (2990), Deleted lines (1 [24]), Word sort (Asc alpha (string)), and Context sort (Asc occurrence order).

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 u
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 Spac
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 flc
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 the

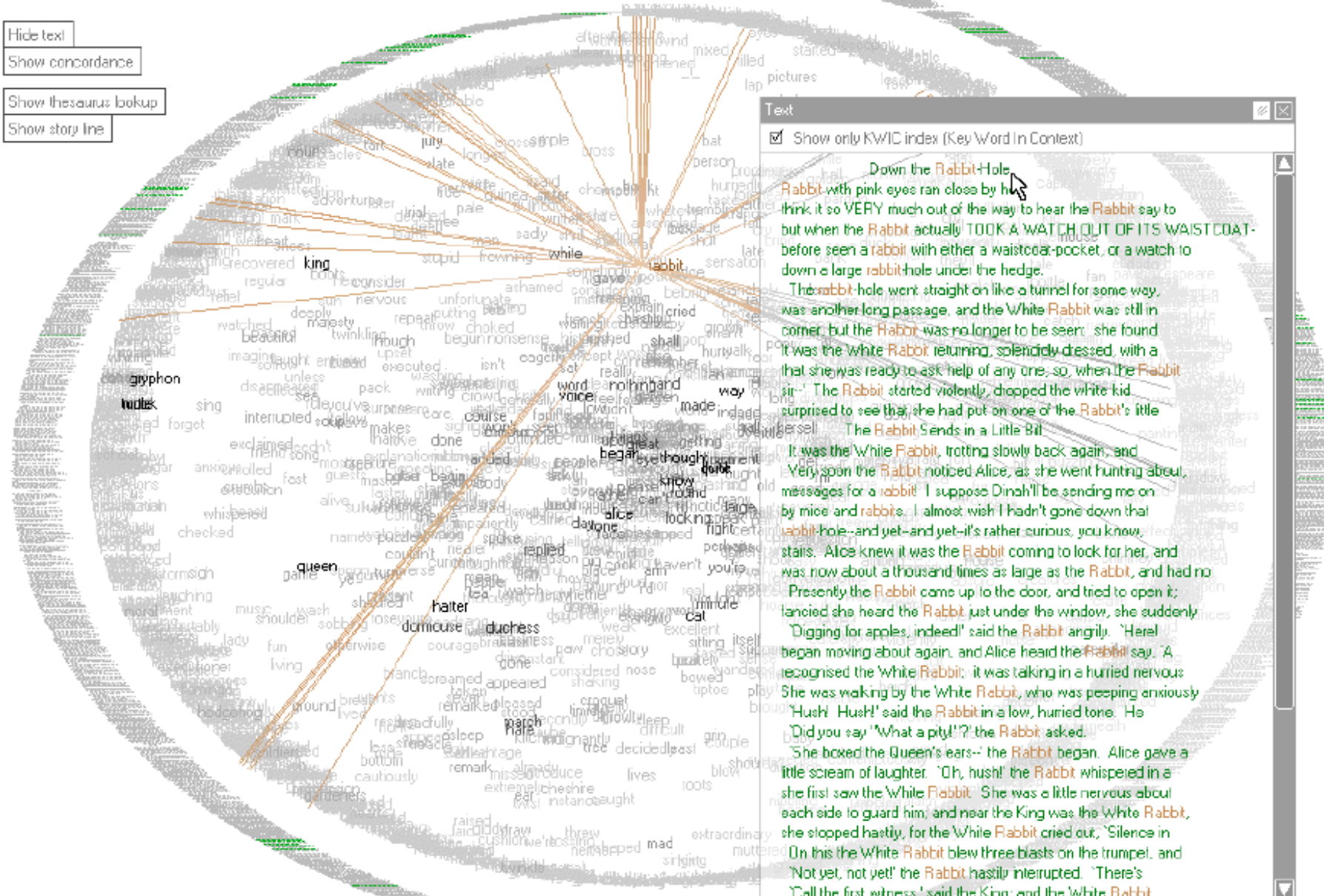
Down the Rabbit-Hole

Hide text

Show concordance

Show thesaurus lookup

Show story line



Text

Show only KWIC index (Key Word In Context)

Down the Rabbit-Hole
 Rabbit with pink eyes ran close by her.
 I think so VERY much out of the way to hear the Rabbit say to
 but when the Rabbit actually TOOK A WATCH OUT OF ITS WAIST COAT-
 before seen a rabbit with either a waistcoat-pocket, or a watch to
 down a large rabbit-hole under the hedge.
 The rabbit-hole went straight on like a tunnel for some way,
 was another long passage, and the White Rabbit was still in
 corner, but the Rabbit was no longer to be seen: she found
 it was the white Rabbit returning, splendidly dressed, with a
 that she was ready to ask help of any one: so, when the Rabbit
 sir'. The Rabbit started violently, dropped the white kid,
 surprised to see that she had put on one of the Rabbit's little
 bersel. The Rabbit Sends in a Little Bill
 It was the White Rabbit, trotting slowly back again, and
 Very soon the Rabbit noticed Alice, as she went hunting about,
 messages for a rabbit: I suppose Dinah'll be sending me on
 by mice and rabbits. I almost wish I hadn't gone down that
 rabbit-hole-and yet-and yet-it's rather curious, you know,
 stars. Alice knew it was the Rabbit coming to look for her, and
 was now about a thousand times as large as the Rabbit, and had no
 Presently the Rabbit came up to the door, and tried to open it;
 fancied she heard the Rabbit just under the window, she suddenly
 'Digging for apples, indeed!' said the Rabbit angrily. 'Here!
 began moving about again, and Alice heard the Rabbit say, 'A
 recognised the White Rabbit: it was talking in a hurried nervous
 She was waking by the White Rabbit, who was peeping anxiously
 'Hush! Hush!' said the Rabbit in a low, hurried tone. He
 'Did you say "What a pity!"?' the Rabbit asked.
 'She boxed the Queen's ears--' the Rabbit began. Alice gave a
 little scream of laughter. 'Oh, hush!' the Rabbit whispered in a
 she first saw the White Rabbit. She was a little nervous about
 each side to guard him; and near the King was the White Rabbit,
 she stopped hastily, for the White Rabbit cried out, 'Silence in
 On this the White Rabbit blew three blasts on the trumpet, and
 'Not yet, not yet!' the Rabbit hastily interrupted. 'There's
 'Call the first witness,' said the King, and the White Rabbit

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

rough with you , be rough with love .

blind ,

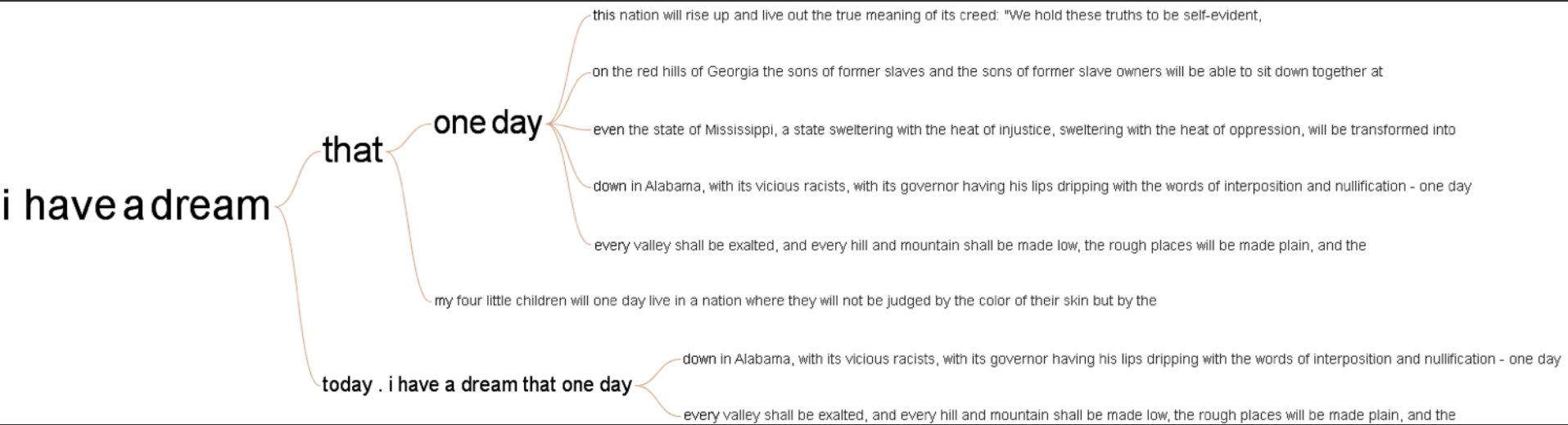
love cannot hit the mark .

it best agrees with night .

Filter Infrequent Runs



Recurrent Themes in Speeches



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

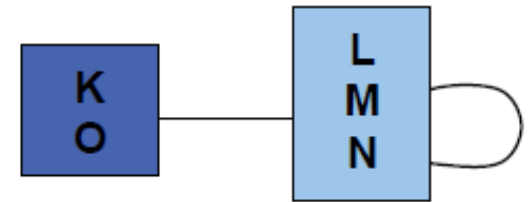
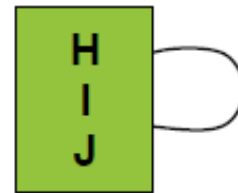
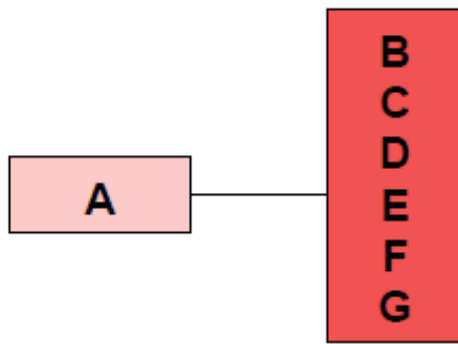
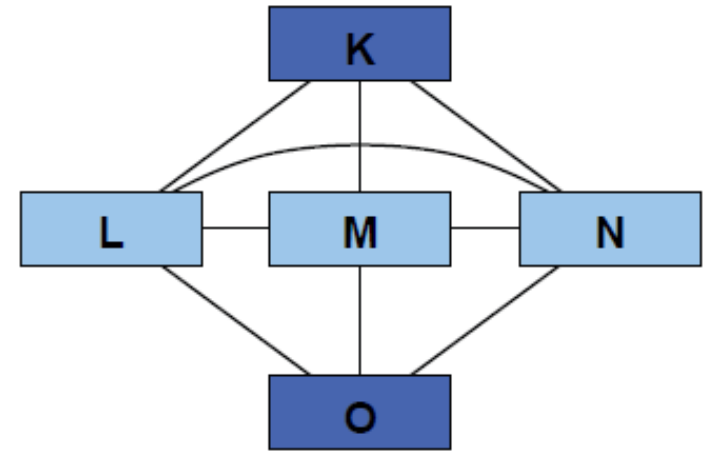
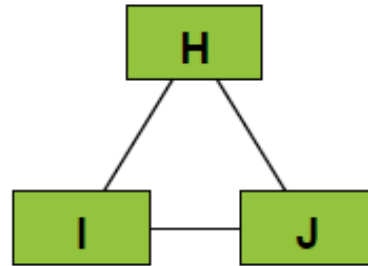
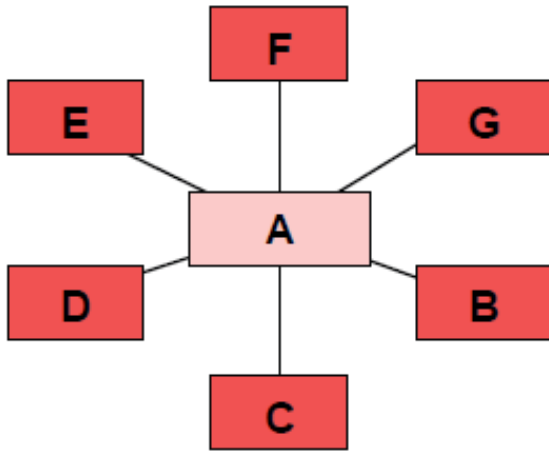
Could be output of regexp or parser.

Visualize patterns in a node-link view

Occurrences -> Node size

Pattern position -> Edge direction

Node Grouping

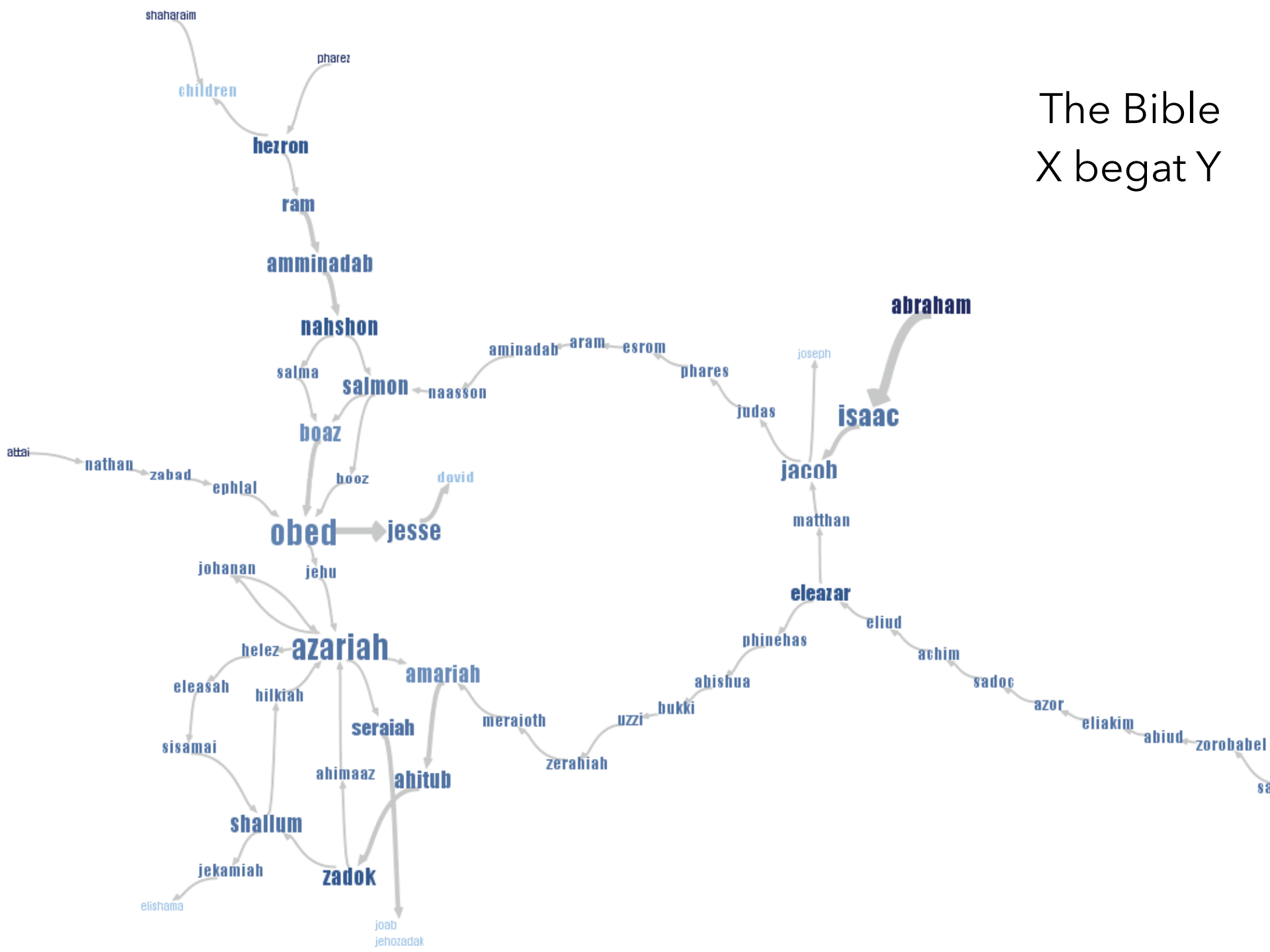


(a)

(b)

(c)

The Bible X begat Y



Document Content

Understand Your Analysis Task

Visually: Word position, browsing, brush & link

Semantically: Word sequence, hierarchy, clustering

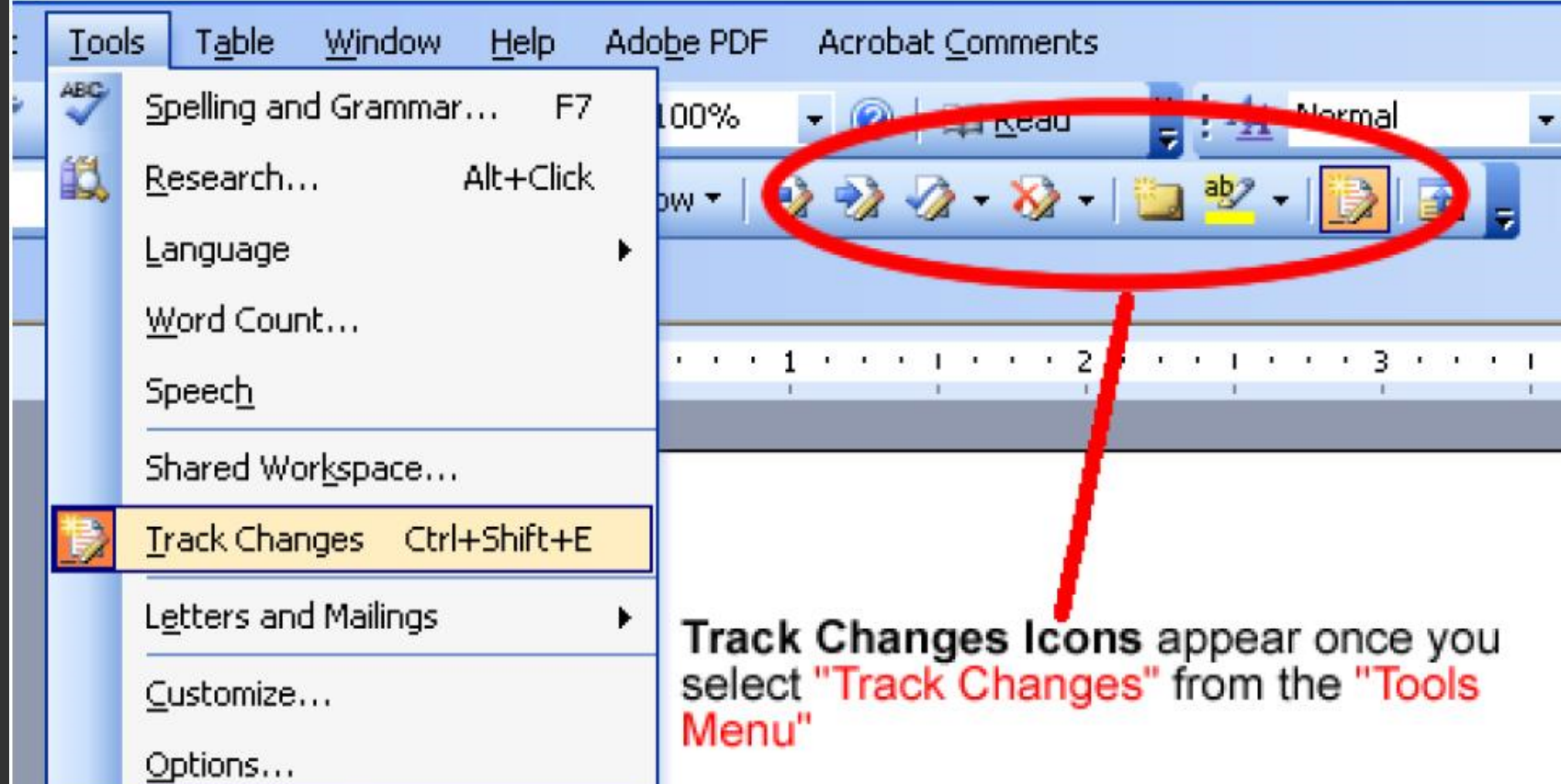
Both: Spatial layout reflects semantic relationships

The Role of Interaction

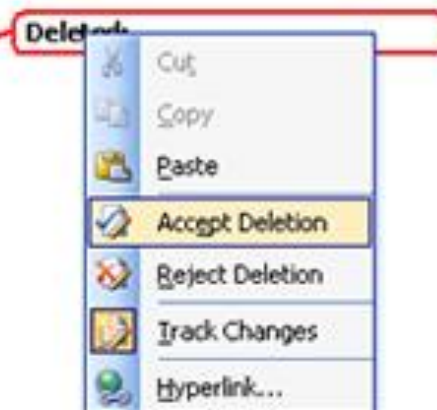
Language model supports visual analysis cycles

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

Evolving Documents



This is a test document to demonstrate the use of tracking changes. The characters in black font represent the original document while the characters in red font represent the changes which are being tracked.



Visualizing Revision History

How to depict contributions over time?

Example: Wikipedia history log

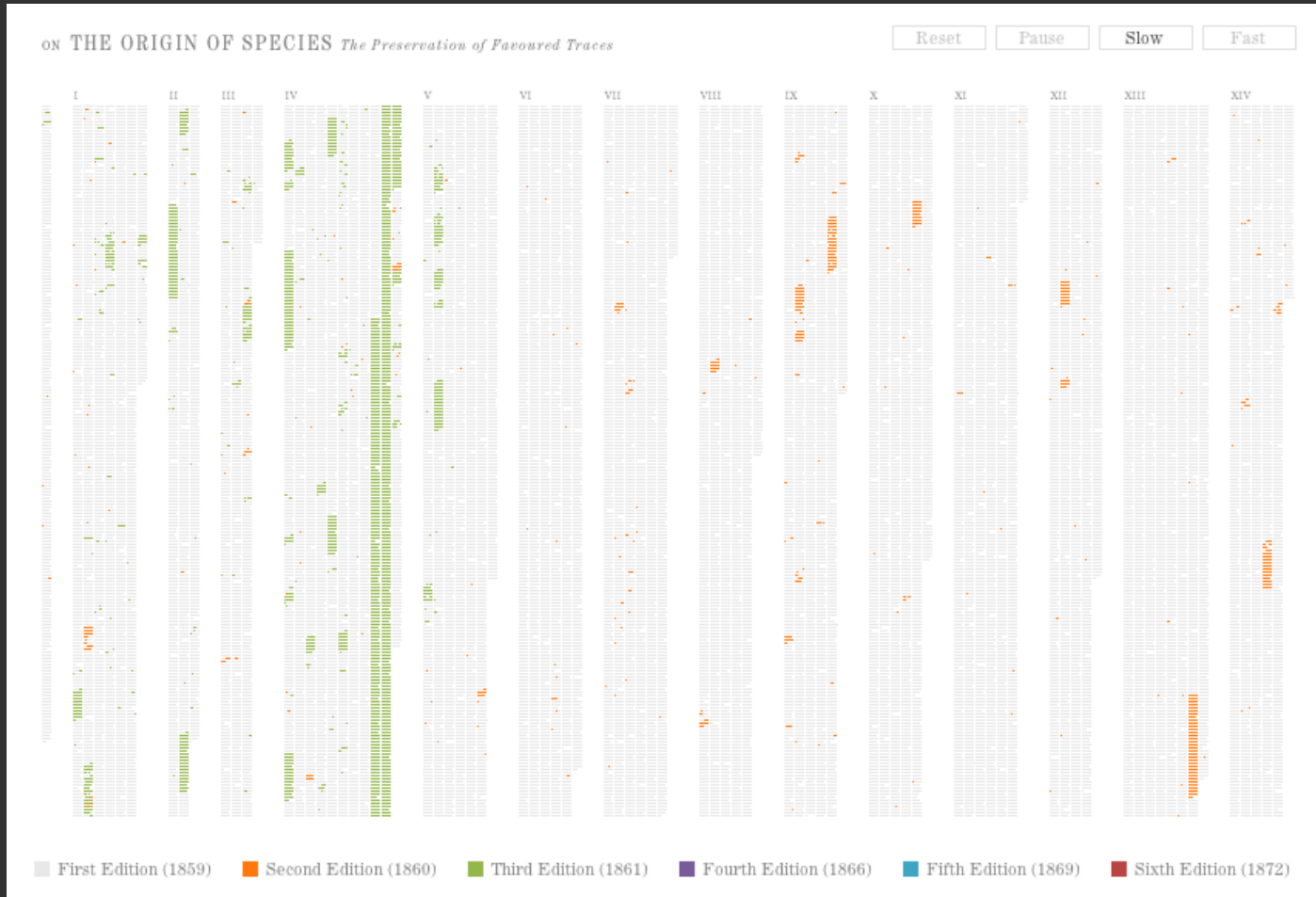
Chocolate

Revision history

Legend: (cur) = difference with current version, (last) = difference with preceding version, M = minor edit

- [\(cur\)](#) [\(last\)](#) . . [12:01, 20 Aug 2003](#) . . [Dysprosia](#) (*neaten to do, rearrange see also*)
- [\(cur\)](#) [\(last\)](#) . . [11:59, 20 Aug 2003](#) . . [Patrick](#)
- [\(cur\)](#) [\(last\)](#) . . [11:52, 20 Aug 2003](#) . . [81.203.98.109](#)
- [\(cur\)](#) [\(last\)](#) . . **M** [18:36, 6 Aug 2003](#) . . [Manika](#) (*corrected spelling*)
- [\(cur\)](#) [\(last\)](#) . . [18:32, 6 Aug 2003](#) . . [Daniel Quinlan](#) (*removing obscure heraldry information, belongs on [[heraldry]] if anywhere*)
- [\(cur\)](#) [\(last\)](#) . . [15:21, 6 Aug 2003](#) . . [Rmhermen](#)
- [\(cur\)](#) [\(last\)](#) . . [15:08, 6 Aug 2003](#) . . [Cyp](#) (*Chocolate often has odd shapes.*)
- [\(cur\)](#) [\(last\)](#) . . [19:14, 3 Aug 2003](#) . . [Daniel C. Boyer](#) (*"chocolate" as shade of gules in heraldry*)
- [\(cur\)](#) [\(last\)](#) . . **M** [02:00, 30 Jul 2003](#) . . [Evercat](#) (*fmt*)

Animated Traces [Ben Fry]



<http://benfry.com/traces/>



Package Explorer

- Struts demo
 - src
 - WEB-INF/src
 - com.mia_software.booster.page_flow_example
 - CollectInformationDispatchAction.java
 - CollectInformationForm.java
 - ResultDispatchAction.java
 - ResultForm.java
 - Test1DispatchAction.java
 - Test1Form.java
 - WelcomePageDispatchAction.java
 - WelcomePageForm.java
 - com.mia_software.struts.generic.back
 - com.mia_software.struts.generic.form

Generation Results

- Results (30)
 - ResultDispatchAction.java
 - Test1DispatchAction.java
 - WelcomePageDispatchAction.java
 - Generated Code
 - Manual Code
 - Generated Code
 - Context menu:
 - Sort
 - Group by Status
 - File Name
 - Relative Path
 - Full Path
 - Compare with Previous Version
 - Show In

Fragment Comparison

Previous Version	Current Version
<pre>// End of user imports public class WelcomePageDispatchAction // associated forward definitions public final static String COLLECTI public final static String LOGOUTAM public final static String TEST1_TO // inherited forward definitions // dispatch action methods declarat public ActionForward enter(ActionMe ActionForward actionForward = r if (form != null) { WelcomePageForm currentForm // execute code on exit of currentForm.onExit(); int i_ENTERXXX=0; CollectInformationForm coll // Start of user code : ret</pre>	<pre>// End of user imports public class WelcomePageDispatchAct // associated forward definitio public final static String COLL public final static String LOGO public final static String TEST // inherited forward definition // dispatch action methods decl public ActionForward enter(Acti ActionForward actionForward if (form != null) { WelcomePageForm current // execute code on exit currentForm.onExit(); CollectInformationForm // Start of user code :</pre>

Diff style: Side-by-side Enable syntax coloring

Files Changed:

1. [sshconsole.js](#): 1 change [1]

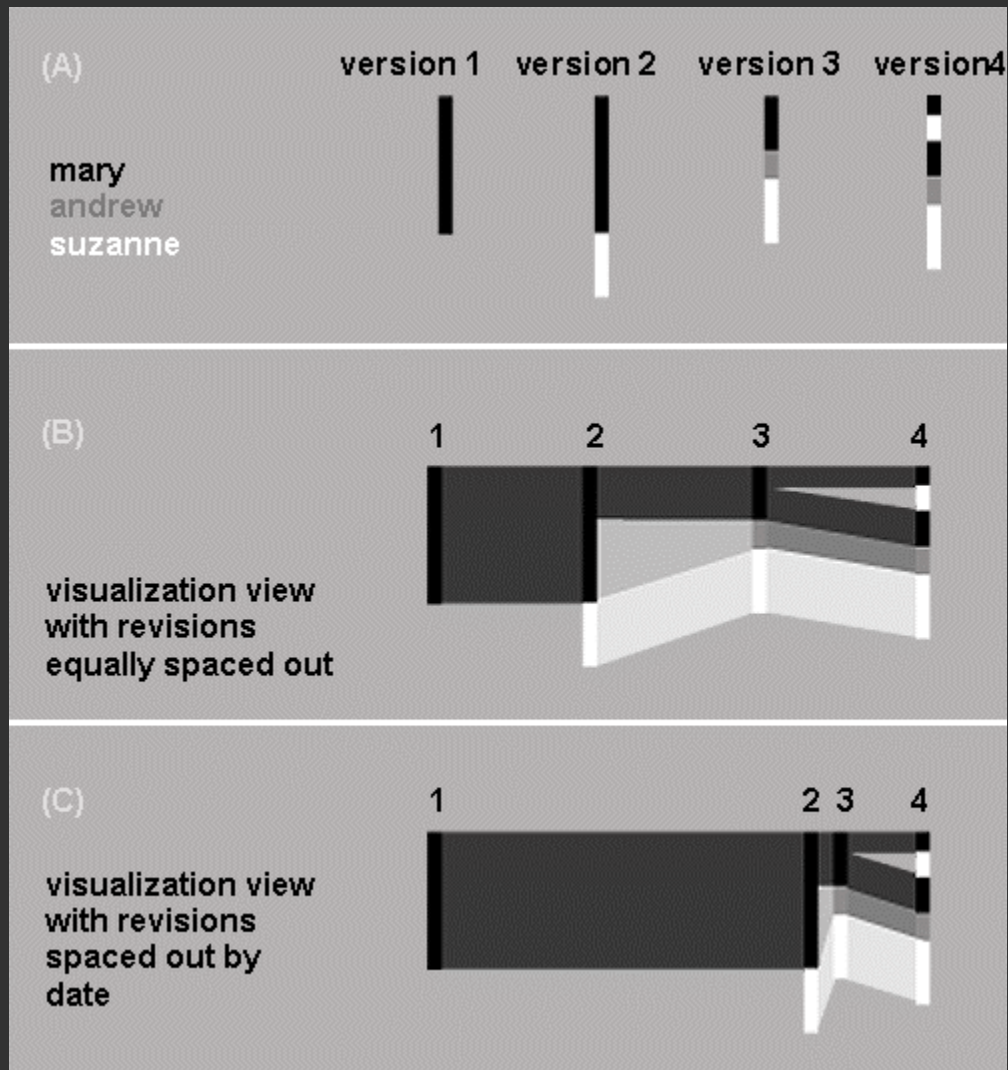
/home/toddw/src/sshconsole-read-only/content/sshconsole.js

50 lines hidden [\[Expand\]](#)

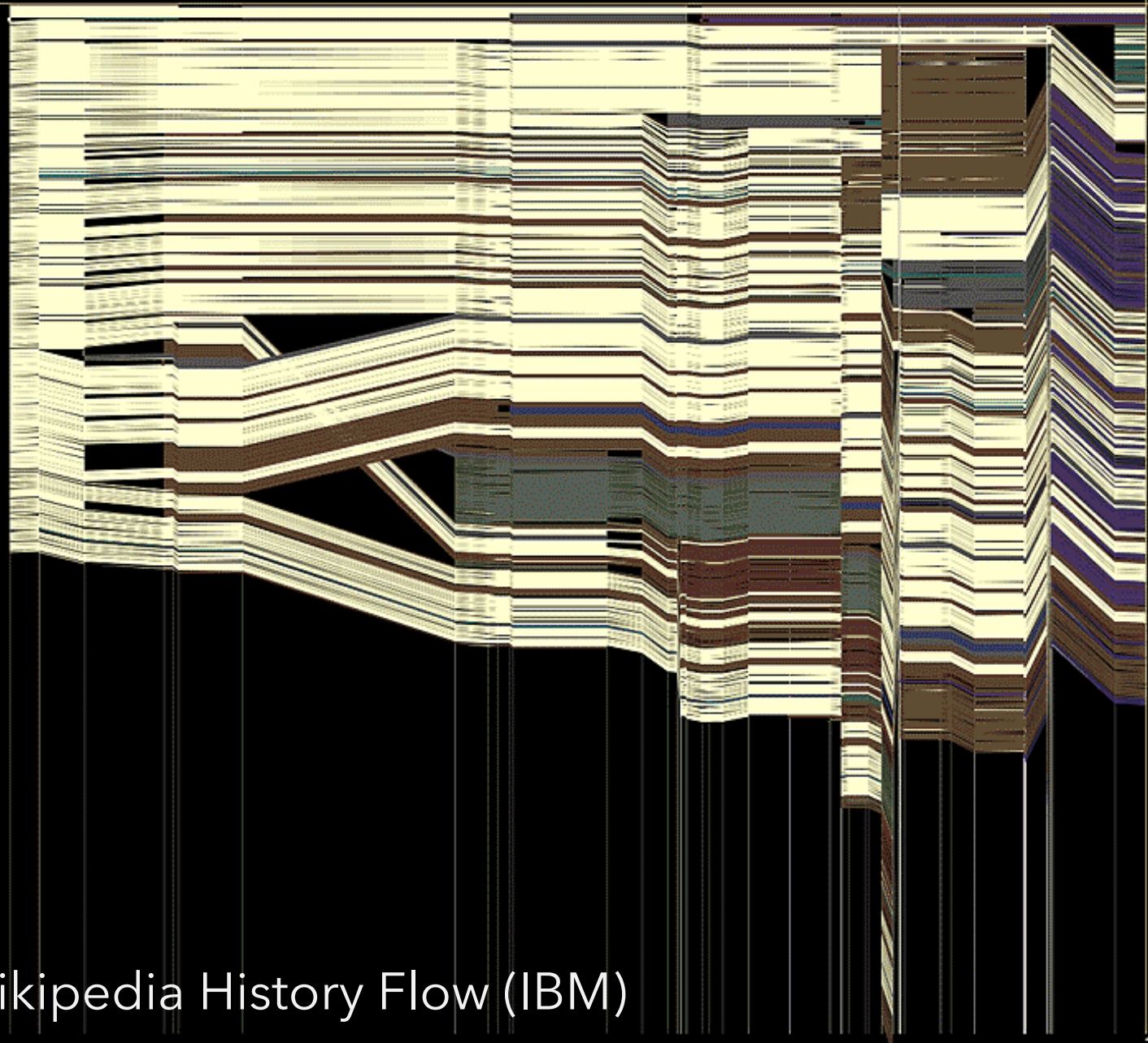
<pre> 51 _term = new VT100(80, 24, "term"); 52 // _term.debug_ = 1; 53 _term.curs_set(true, true, _term_box_element); 54 _term.noecho(); 55 56 // Replace the go_getch_ function with our own, this is called 57 // for every keypress that is passed through the terminal to the 58 // remote server. The character is already converted into the 59 // required VT100 character sequence(s). 60 VT100.go_getch_ = function() { 61 var vt = VT100.the_vt_; 62 if (vt === undefined) { 63 return; 64 } 65 var ch = vt.key_buf_.shift(); 66 //dump("go_getch_:: ch: '" + ch + "'\n"); 67 if (ch === undefined) { 68 return; 69 } 70 if (vt.echo_ && ch.length == 1) { 71 vt.addch(ch); 72 } 73 if (_ssh_channel) { 74 _ssh_channel.sendStdin(ch); 75 } 76 } 77 78 var serverTextbox = document.getElementById("sshconsole_server_textbox"); 79 var connectionText; 80 if ('connectionText' in window.arguments[0]) { 81 connectionText = window.arguments[0].connectionText; 82 } else { </pre>	<pre> 51 _term = new VT100(80, 24, "term"); 52 // _term.debug_ = 1; 53 _term.curs_set(true, true, _term_box_element); 54 _term.noecho(); 55 56 // Replace the go_getch_ function with our own, this is called 57 // for every keypress that is passed through the terminal to the 58 // remote server. The character is already converted into the 59 // required VT100 character sequence(s). 60 VT100.go_getch_ = function() { 61 var vt = VT100.the_vt_; 62 if (vt === somevalue) { 63 return; 64 } 65 var ch = vt.key_buf_.shift(); 66 67 if (ch === undefined) { 68 return; 69 } 70 if (vt.echo_ && ch.length == 1) { 71 vt.refres(); 72 } 73 if (_ssh_channel) { 74 _ssh_channel.sendStdin(ch); 75 } 76 } 77 78 var serverTextbox = document.getElementById("sshconsole_server_textbox"); 79 var connectionText; 80 if ('connectionText' in window.arguments[0]) { 81 connectionText = window.arguments[0].connectionText; 82 } else { </pre>
--	--

174 lines hidden [\[Expand\]](#)

History Flow [Viegas et al.]



authors	posts
Zundark	1
The Cunctator	1
The Epost	1
Conversion script	1
RK	1
Freob	1
B4hand	1
KarakizeArchon	1
Stephen Gilbert	1
Sraubenstein	8
Mimccorn	5
Iels	1
Derek Ross	1
Dante Alighieri	2
Maveric149	3
Jzzbug	2
Jdirl	8
Theanthrope	1
Wesley	2
Dreamword	1
Stevetigo	4
Camembert	1
Hephaestos	2
Zoe	1
MyRedDice	1
G-Man	2
Kingturtle	1
Montrealais	1
770	1



Abortion

(Revision as of 22:56 4 Jun 2003)

"**Abortion**," in its most commonly used sense, refers to the deliberate early termination of a pregnancy, resulting in the death of the embryo or fetus. [1] Medically, the term also refers to the early termination of a pregnancy by natural means ("spontaneous abortion" or *miscarriage*), or to the cessation of normal growth of the embryo or fetus (1 in 5 of all pregnancies, usually within the first 12 weeks) or to the cessation of normal growth of the body part or organ. What follows is a discussion of the issues related to deliberate or "induced" abortion.

Methods

Depending on the stage of pregnancy an abortion is performed by a number of different methods. The earliest terminations (before nine weeks) are usually performed by a chemical abortion. A chemical abortion is the usual method, though *mifepristone* is usually the only legal method. Although research has uncovered similar effects from *methotrexate* and *misoprostol*. Concern with chemical abortion and extending up until around the fifteenth week, suction-aspiration vacuum abortion is the most common approach, replacing the more risky dilation and curettage (D & C). From the fifteenth week up until around the eighteenth week a surgical dilation and extraction (D & E) is used.

As the fetus size increases other techniques may be used to secure abortion in the third trimester. premature expulsion of the fetus can be induced with prostaglandin, this can be coupled with injecting the amniotic fluid with saline or urea solution. Very late abortions can be brought about by the controversial intact dilation and extraction (D & X) or a hysterotomy abortion, similar to a caesarian section.

The controversy

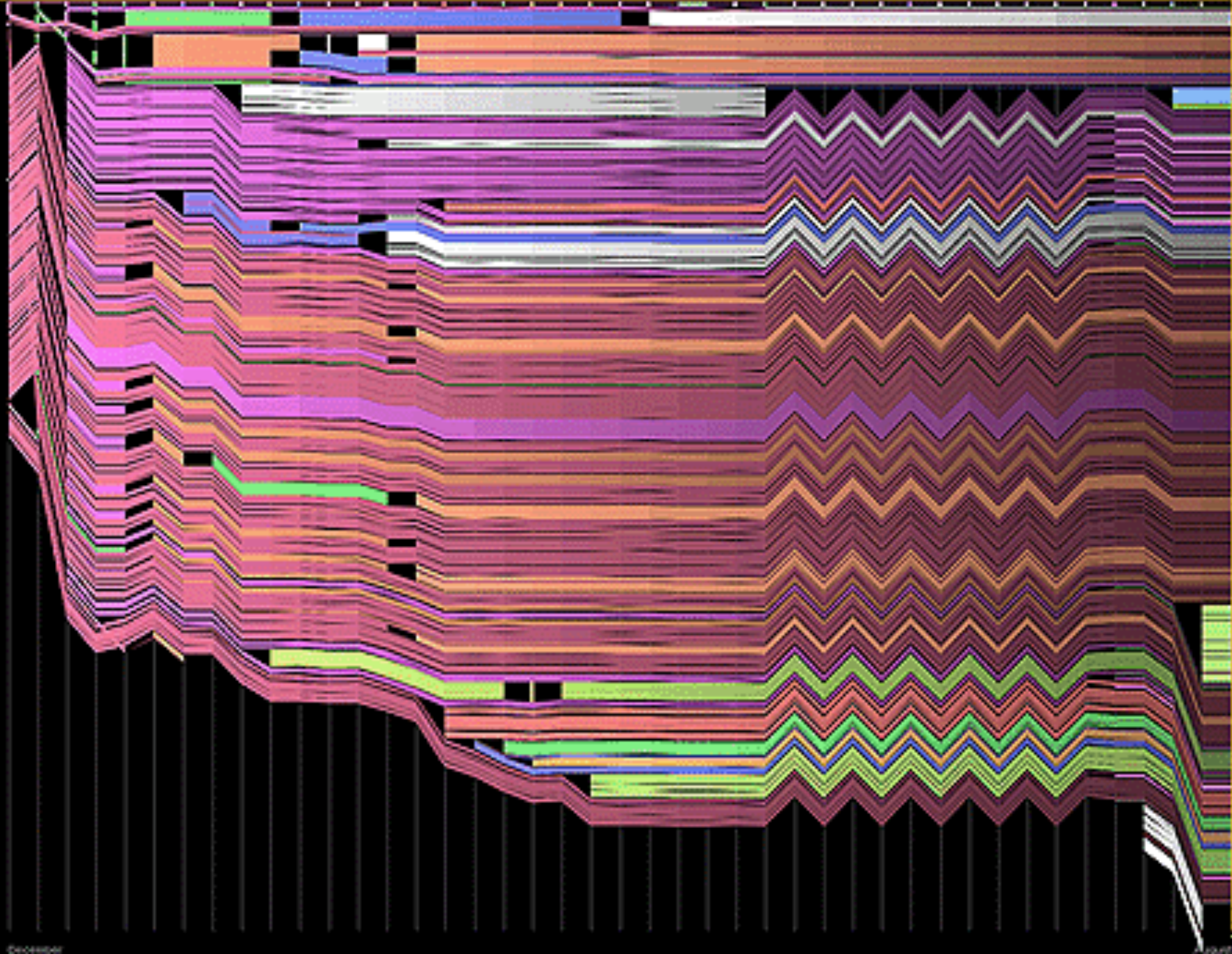
The morality and legality of abortion is a highly important topic in applied ethics, and is also discussed by legal scholars, and religious philosophers. Important facts about abortion are also reported by sociologists and historians.

Abortion has been common in most societies, although it has often been opposed by some institutionalized religions and governments. In the 19th century, politics in the United States and Europe, abortion became commonly accepted by the 20th century. Additionally, abortion is accepted in China, India and other populous countries. The Catholic Church remains opposed to the procedure, however, and in other countries, notably the United States and the (predominantly Catholic) Republic of Ireland, the controversy is extremely active, to the extent that even the respective positions are subject to heated debate. While those on both sides of the debate are generally peaceful, if heated, in their defense of their positions, the debate is sometimes characterized by violence. Though true of both sides, this is more marked on the side of those opposed to abortion, because of what they see as the gravity and urgency of their views.

The central question

The central question in the abortion debate is a clash of presumed or perceived rights. On the one hand, is a fetus (sometimes called the "unborn" by pro-life/anti-abortion advocates) a human with a right to life, and if so, at what point in pregnancy does the fetus become human? On the other hand, is a fetus part of a woman's

Wikipedia History Flow (IBM)



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"

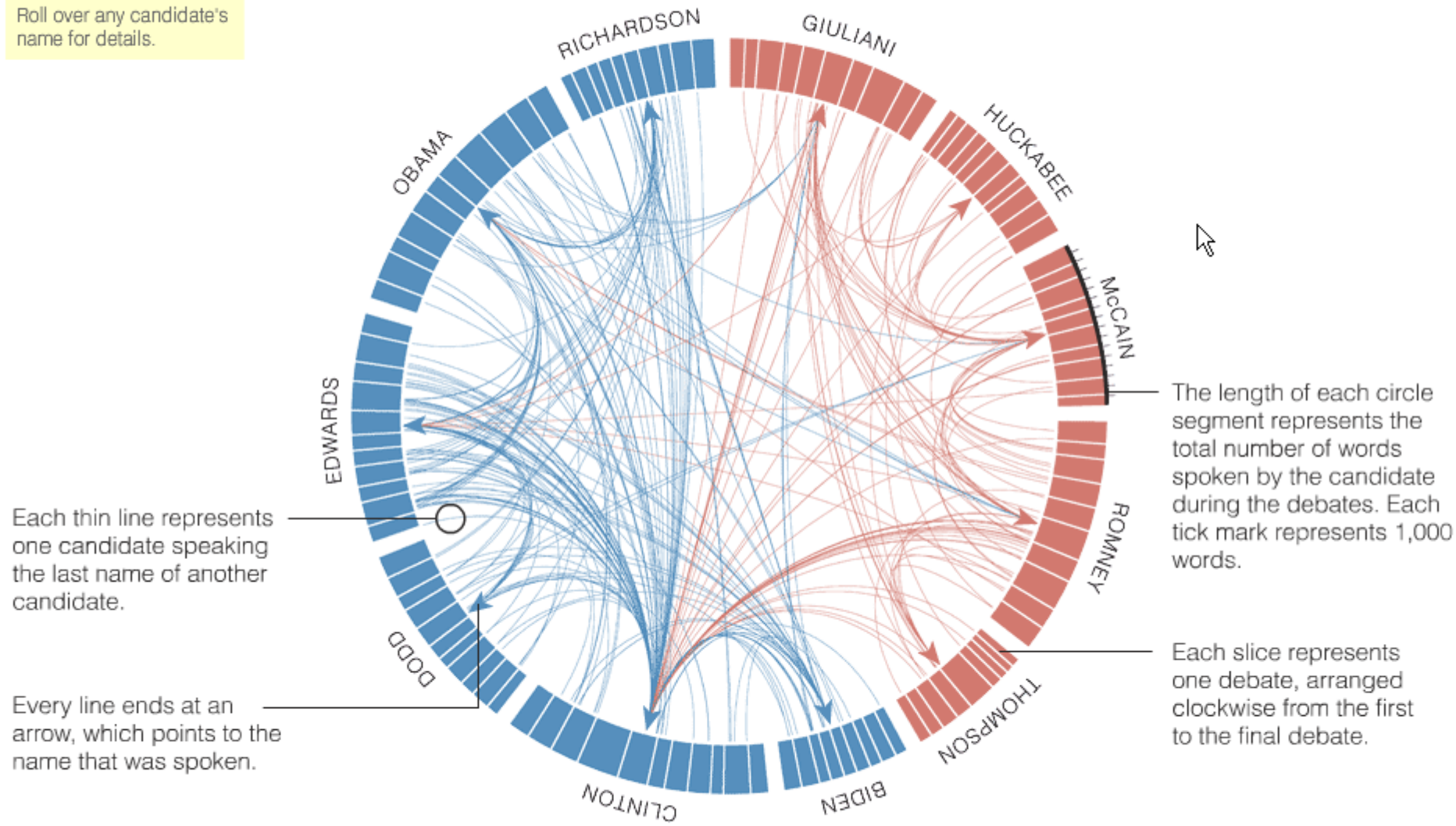
Who x Who -> Social network

Who x Who x What x When -> Information flow

Naming Names

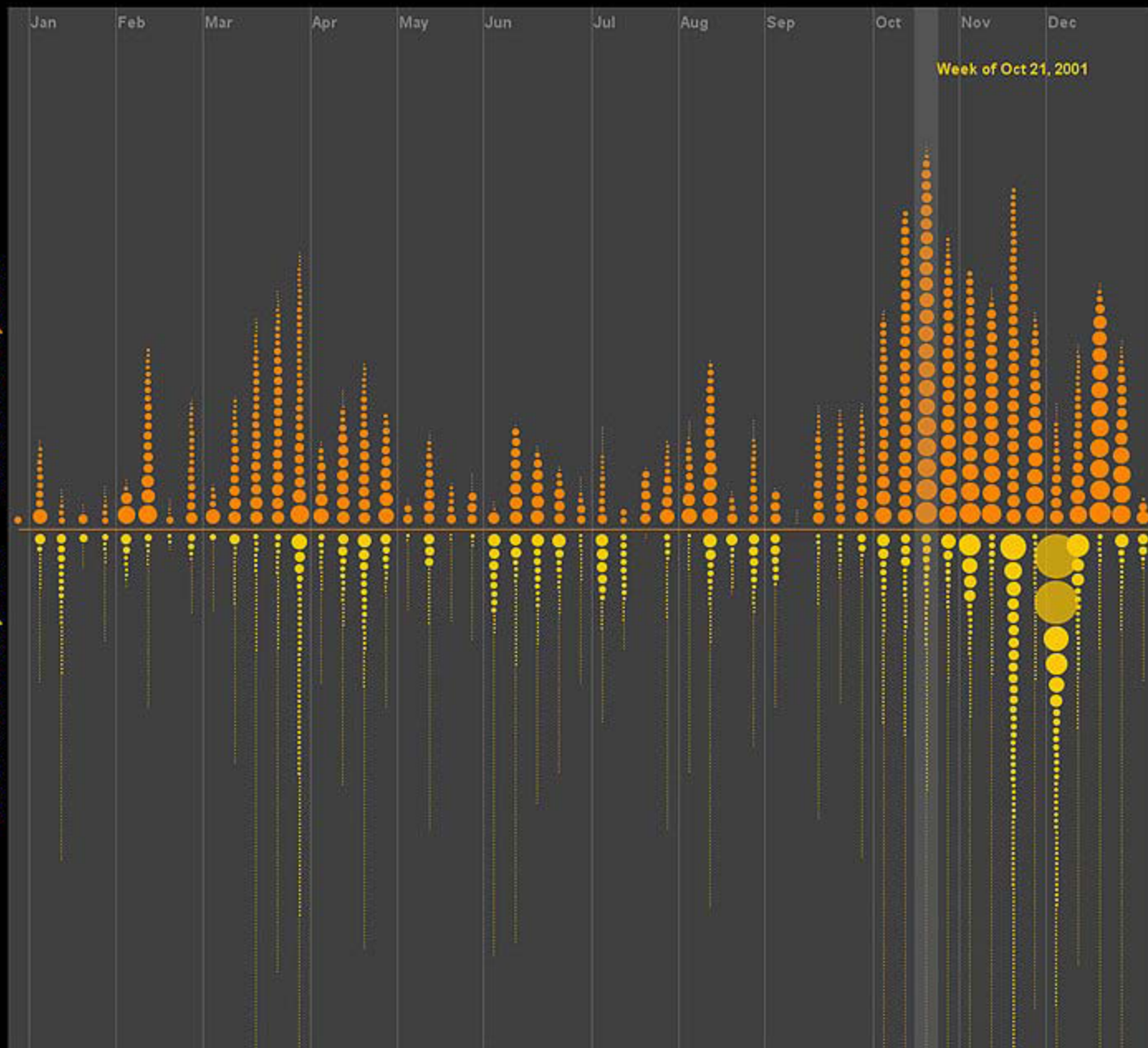
Names used by major presidential candidates in the series of Democratic and Republican debates leading up to the Iowa caucuses.

Roll over any candidate's name for details.



threads initiated by author

threads not initiated by author



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	16
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 Kalie 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 Kalie 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 debi'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 !!	2
WET 195 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

author: rubat@pam.org

[back to message](#)

Jan Feb Mar Apr May Jun Jul Aug Sep Oct Nov Dec

Week of May 6, 2001

In Genes Scientifically Correct?

subject | # of posts

threads initiated by author

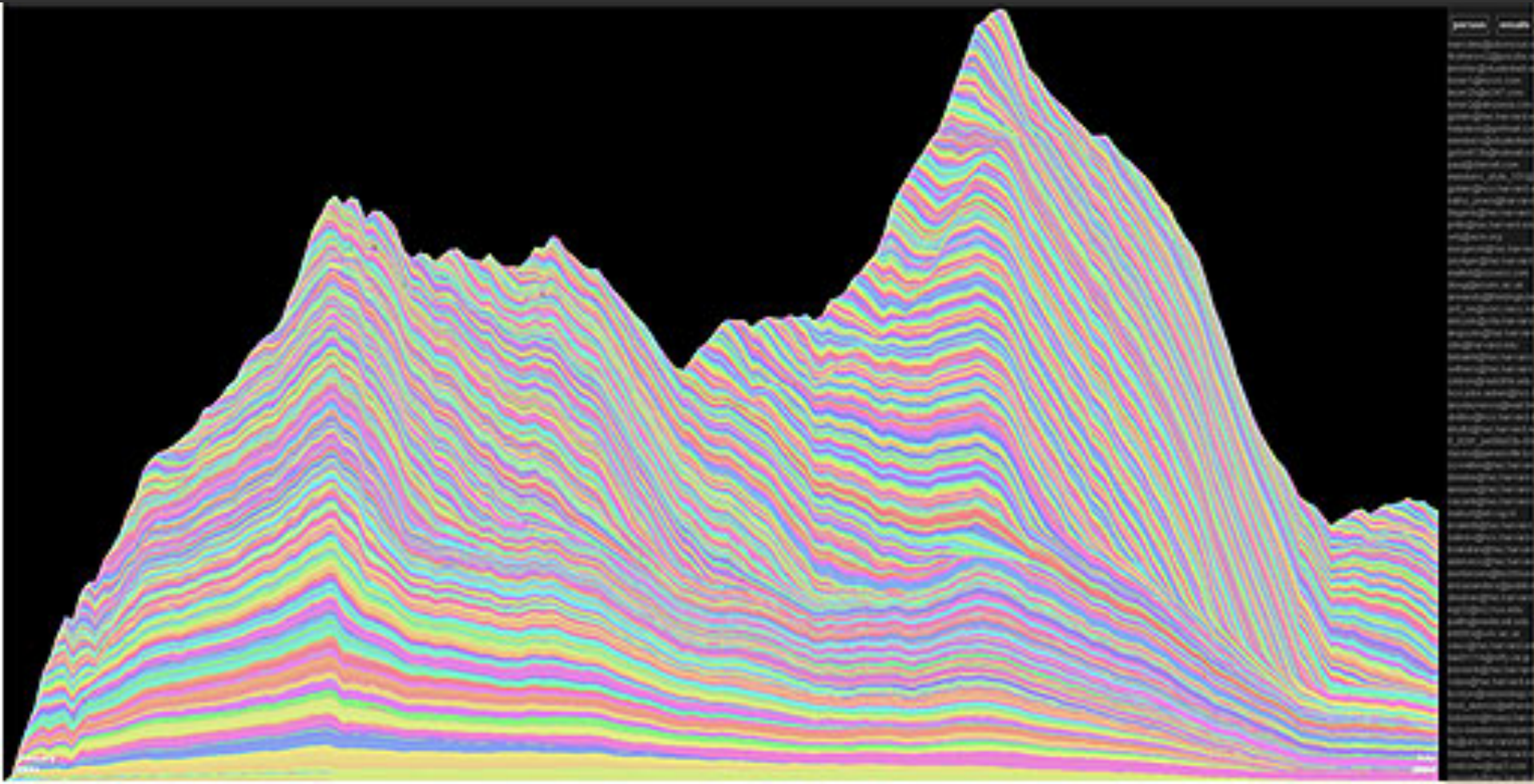
threads not initiated by author



subject	# of posts
Antimicrobials is	21
Etymology: Scientist	16
OLD TESTAMENT	17
Please define Ho	10
Evolution facts	14
Why vouchers were	20
Scientist against a	15
How to teach Ho	14
TDM vs. CDMA	10
Dot and O-2	7
Genes is wrong	7
The Abroad is Out	7
President Quot Ho	7
I got a question	6
A Faithful Deal	6
Freedom from Ho	5
Original Sin-Bad F	5
Evolution: Proof that C	5
SCIENTIFIC PRO	5
Christians teach a	4
Antimicrobials is	4
All post's remain	4
Car Theft	4
1000Ho vs. 1000	4
Vouchers K12	4

Email Mountain [Viegas]

Conversation by person over time (who x when).





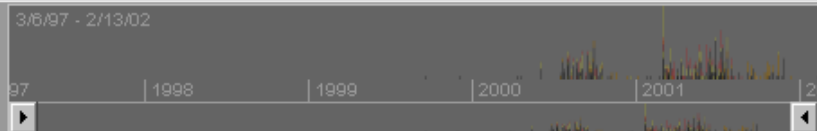
Enron E-Mail Corpus

[Heer]



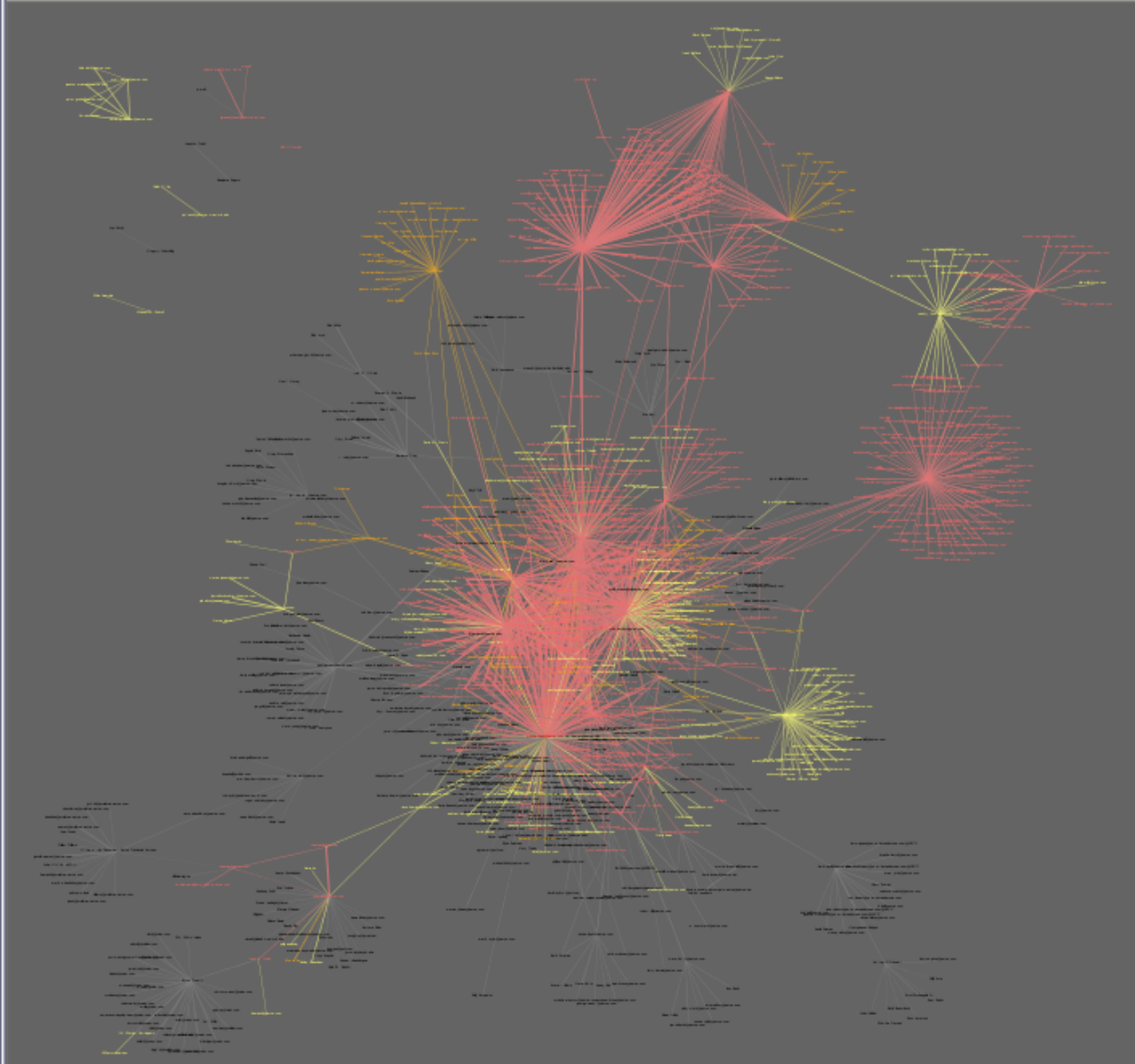
connectivity >>

community >>



search >>

search >>



steven.kean@enron.com

- 2000-09-01 04:25:00.0 Linda Jenkins on "Jerry's Show" Mond
- 2000-09-02 10:14:00.0 Re: The Governors' Natural Gas Summ
- 2000-09-08 10:03:00.0
- 2000-09-10 14:07:00.0 CPUC Hearing in SD on 9/8
- 2000-09-10 16:20:00.0 Re: Fletcher School/Enron
- 2000-09-13 00:57:00.0 Re: Contact

0 2 0 0 0 0 1 0 0 0 0 0 0 0

ID: 174285

Subject:

From: <steven.kean@enron.com>

Date: 2000-09-08 10:03:00.0

To: <kmagrude@enron.com>

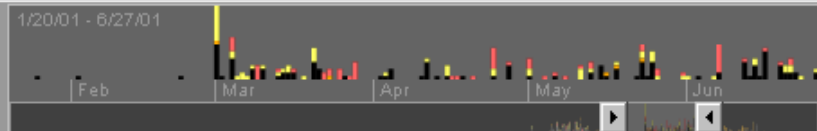
Cc: Richard Shapiro <rshapiro@enron.com>

Got your message. I'm testifying at the Congressional hearing and Dasovich is covering **FERC**. I think Jeff's comments were taken out of context. He said policymakers do need to take care of small customers whose bills are tripling. Frankly, we'd get slaughtered if we said anything else. But he also said there is a right way and a wrong way to do it. Enron and others had provided a market based answer by offering a fixed price deal to SDG&E (which would have enabled them to cap rates to those who had not switched. **California** elected instead to cap rates and deficit spend (ie create a deferral account). I don't think we can stand for anything that doesn't protect the small customers, but we can continue to emphasize the market based solutions. One of the messages in my testimony will be: customers should be encouraged to choose. Those who did are doing fine.

Messages

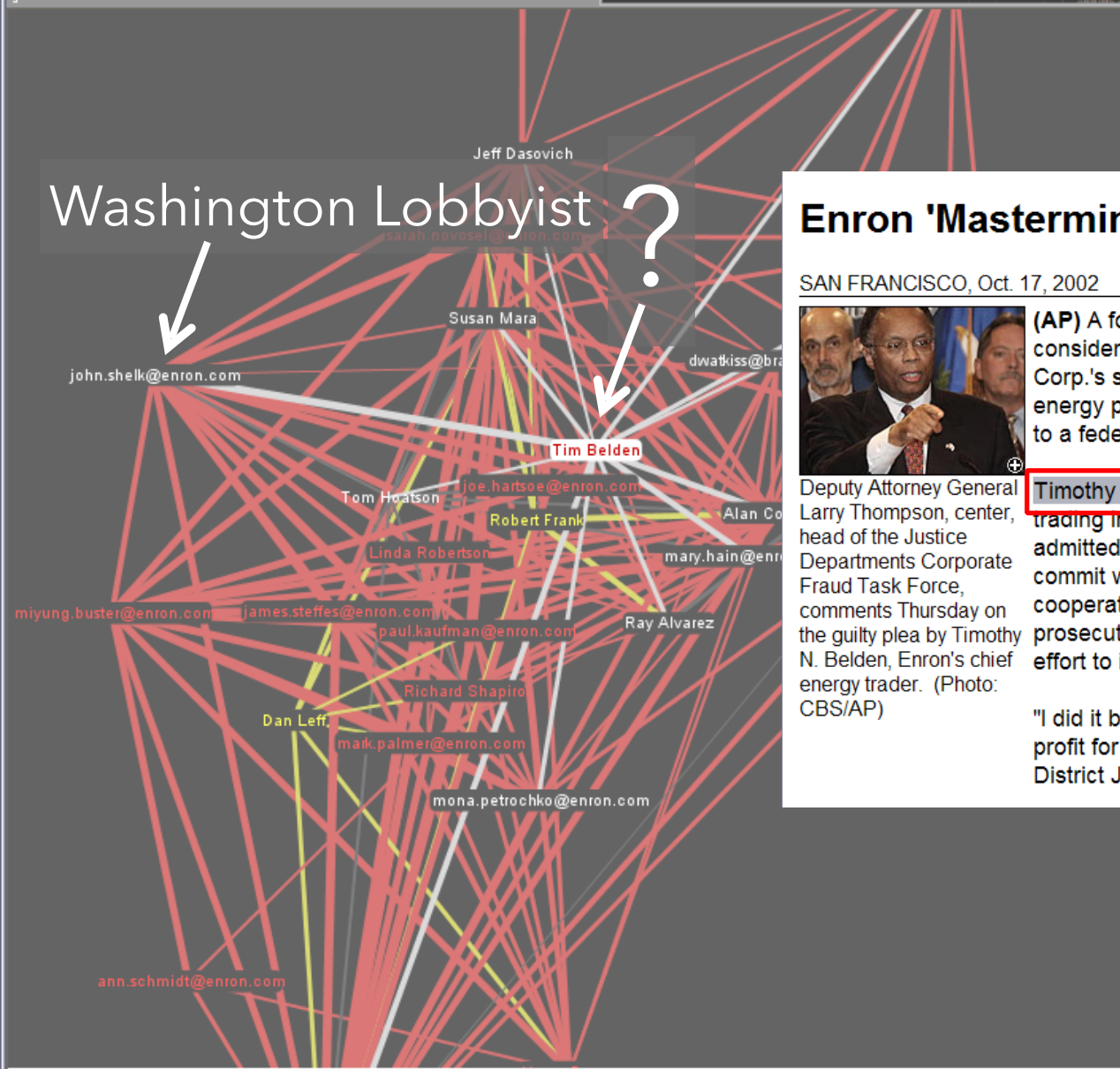
connectivity >>

community >>



search >> california

search >> ferc



Washington Lobbyist ?

- Tim Belden**
- 2001-06-06 16:48:00.0 ISO's Response to BPA Rebuttal of Sh
 - 2001-06-07 11:00:00.0 Legislative Update -- Two Track In The
 - 2001-06-18 00:15:00.0 White House To Support FERC Action
 - 2001-06-19 04:22:00.0 NEWS FLASH ON THIS MORNING'S
 - 2001-06-20 10:37:00.0 Today's Senate Hearing
 - 2001-06-21 02:15:00.0 More on FERC Refunds

Enron 'Mastermind' Pleads Guilty

SAN FRANCISCO, Oct. 17, 2002



(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.

Deputy Attorney General Larry Thompson, center, head of the Justice Departments Corporate Fraud Task Force, comments Thursday on the guilty plea by Timothy N. Belden, Enron's chief energy trader. (Photo: CBS/AP)

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.

from four western governors -- those from Arizona, North Dakota, Utah and Wyoming -- saying that since FERC has acted, there is no need for Congress to pursue price control legislation.

There were a series of questions and comments on details and technical aspects of the orders. I will do an e-mail on these items later today. Please advise if you have any questions or comments.

Messages

Document Collections

<h1>Pakistan's hopes fade in international arena</h1>		<h1>Clinton criticises Israeli demolitions in E. Jerusalem</h1>		<h1>Cricket ref: Pakistan police fled during ambush</h1>		<h1>Was Fallon funny?</h1> <p>How Could Rihanna Take Back Chris Brown?</p> <p>Economy to Dominate Annual Chinese Gathering</p> <p>Michael Jackson seeks comeback in London</p>		<h1>CONCERT REVIEW: Britney Yet to Hit Stride</h1> <p>Jason & Molly: 'Fighting Like Hell to Make It Work'</p> <p>'High School Musical' gets a new cast</p> <p>Miley Movie News</p>		<h1>The Rush To Bash Rush</h1> <p>Marines: Multiple errors caused San Diego crash</p> <p>Obama: Contracting overhaul to save \$40 billion a year</p>		<h1>The Strengths and Weaknesses of the Chandra Levy Case</h1> <p>Chester Stiles found guilty on all counts</p> <p>Fed spending bill contains billions in earmarks for LI</p> <p>Supreme Court Divided Over Judicial Bias Case</p> <p>Campaign Aide Tapped to Head FCC</p>									
<h1>British Prime Minister to Address US Congress</h1> <p>Commander says Iran missiles can reach Israel atom sites</p> <p>Interim Leader Takes Over in Guinea-Bissau</p>		<h1>World court issues arrest warrant for Sudan's Bashir</h1> <p>Has Pakistan become the central front?</p> <p>Zimbabwe's Tsvangirai calls for end to sanctions</p>		<h1>Afghan Election Commission Says Early Vote Not Possible</h1> <p>South Africa: United States May Boycott UN Racism Conference</p> <p>Mexico troops enter drug war city</p> <p>China sees smaller defense boost amid economy woes</p> <p>Germany: 2 still missing after building collapse</p>		<h1>Bomb Kills Three Canadian Servicemen in Southern Afghanistan</h1> <p>Gandhi kin appeals for return of memorabilia</p> <p>Man hoping for 'world's best job'</p> <p>Cameroon: Gov't, Church to Share Cost of Pope's Visit</p> <p>Bad altimeter a factor in Netherlands plane crash</p>		<h1>Coast Guard Calls Off Search for Missing Players, Fellow Boater</h1> <p>Manly, LA close</p> <p>Girardi says hip bothered A-Rod last year</p> <p>Suns-Heat Preview</p> <p>Pacers-Trail Blazers Preview</p> <p>Red Wings-Avalanche Preview</p> <p>LA Clippers</p>		<h1>Beckham's crafty new loan deal exposes the truth</h1> <p>Free-agent center Matt Birk 'absolutely' considers Ravens</p> <p>Draw is given back injection</p>		<h1>Wednesday eye-opener: Do you sign Kurt Warner?</h1> <p>Sharks' jammed lineup no match for Stars</p> <p>Free-agent center Matt Birk 'absolutely' considers Ravens</p> <p>Draw is given back injection</p>		<h1>Oil Gains a Second Day on Speculation China Will Boost Stimulus</h1> <p>Bernanke's AIG Blast May Mean More Curbs on Risk, Concentration</p> <p>Ukraine May Miss Deadline for Gazprom Gas Payment</p>		<h1>Chrysler spends \$5566 per car on US incentives, takes top Canada ...</h1> <p>Hollywood mulls little without Blockbuster</p> <p>UBS Gets New Ammo With Villiger</p> <p>US vendor sector contracts again in February Q&M</p>		<h1>Apple introduces new Mac desktop computers</h1> <p>In-depth review: Kindle 2, the Apple TV of books</p> <p>CEBIT-IT industry will be back, Schwarzenegger says</p> <p>Phew! Asteroid's passing was a cosmic near-miss</p> <p>Tiny moon discovered orbiting Saturn</p> <p>Microsoft's Starburst looks like a new search engine - analyst says</p> <p>FDA Panel to Consider Whether Use of Side-Effect-Prone Drug OK</p>		<h1>Intel and TSMC: What are they thinking?</h1> <p>Senators push to boost FDA food safety system</p> <p>AIDS affecting older adults over 50s, WHO</p> <p>Former HP leader diagnosed with cancer</p> <p>Concerted Effort Needed to Fight Drug-Resistant Flu Strain</p> <p>NC doctor helps colleagues track online reviews</p> <p>My county's health care workers to have 9/11 baby</p> <p>Modular of claims linked to one tooth</p>	

ARCHIVED	THU	FRI	SAT	SUN	MON	YEST	TODAY	LIVE	HOW
00:00									
06:00									
12:00									
18:00									

NewsMap: Google News Treemap [Weskamp]

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?

person [dropdown] [Add all] [Clear]

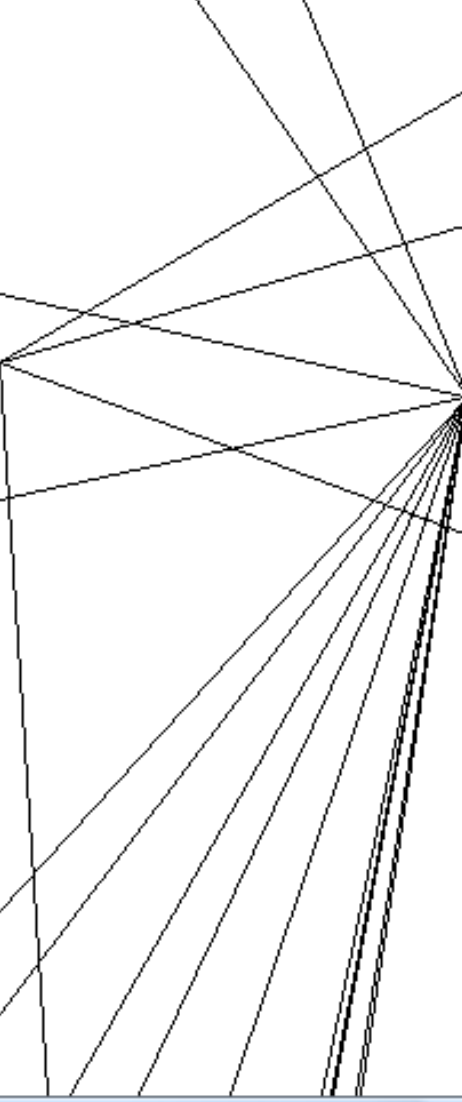
ABC [list icon] [list icon] [list icon] [list icon] [list icon] [list icon]

Show all connections

place [dropdown] [Add all] [Clear]

ABC [list icon] [list icon] [list icon] [list icon] [list icon] [list icon]

- | 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.
- | Jamat Sved



- 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

Similarity & Clustering

Compute vector distance among docs

For TF.IDF, typically cosine distance

Similarity measure can be used to cluster

Topic modeling

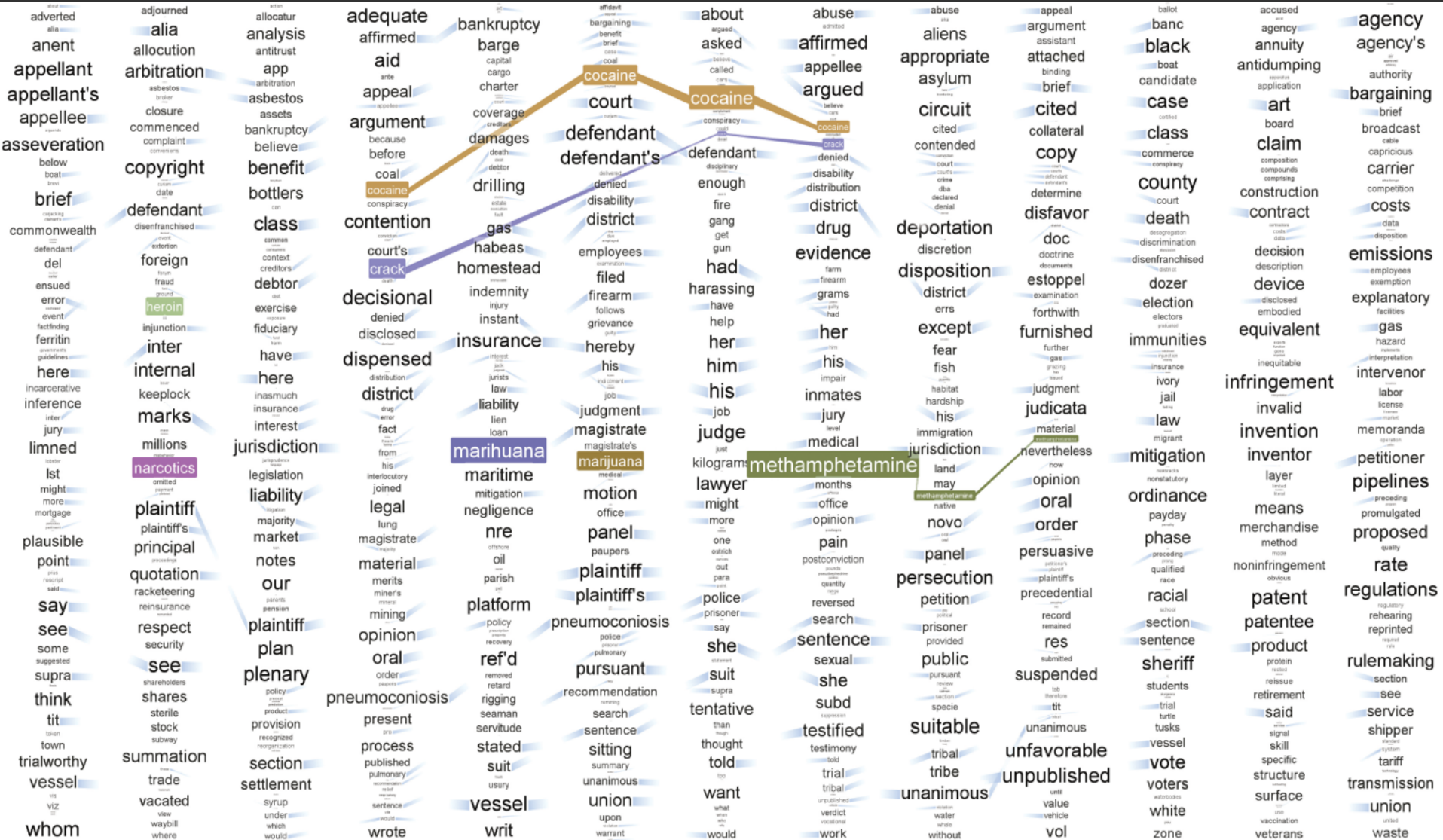
Assume documents are a mixture of topics

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

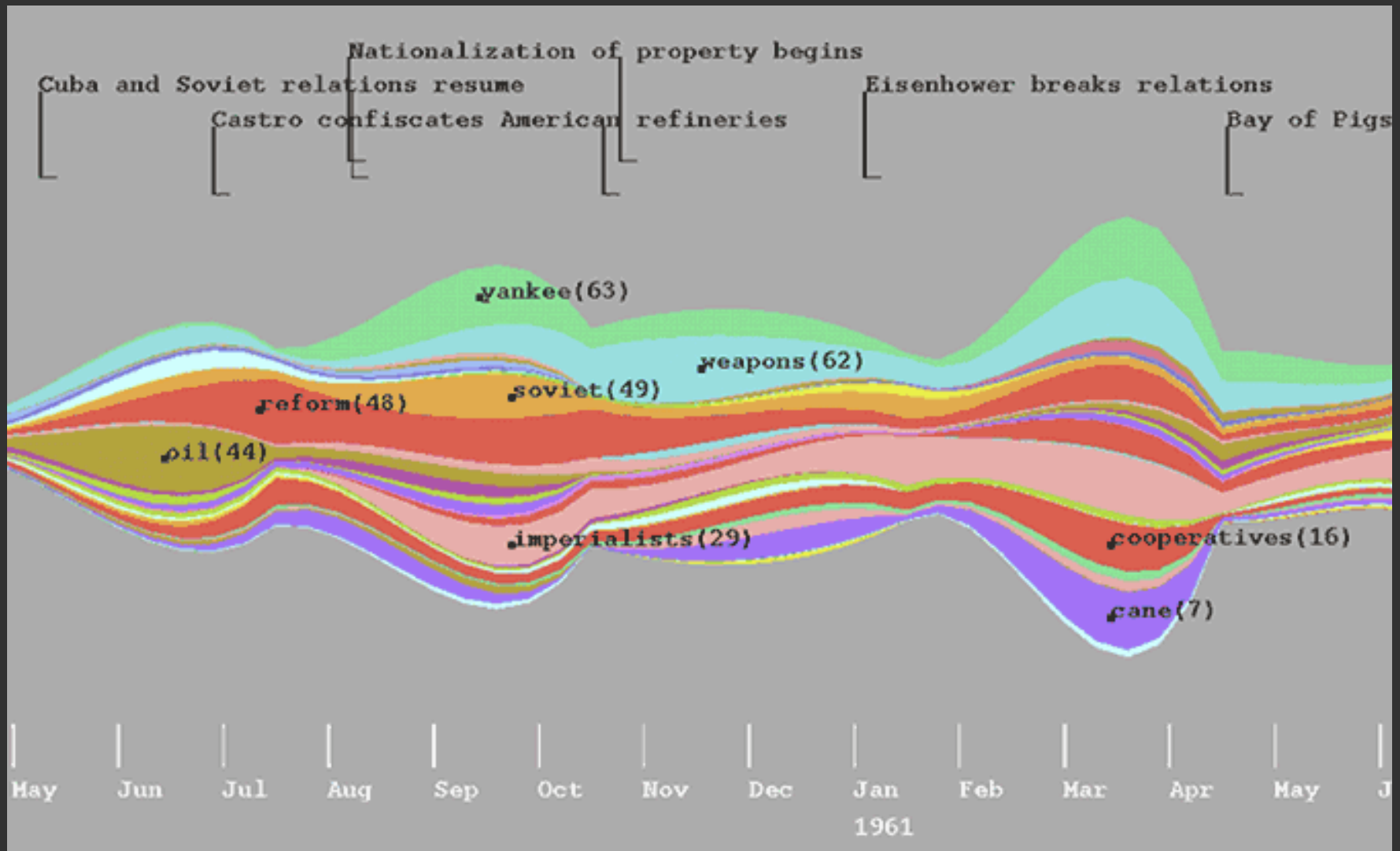
Latent Semantic Analysis (LSA): reduce term matrix

Latent Dirichlet Allocation (LDA): statistical model

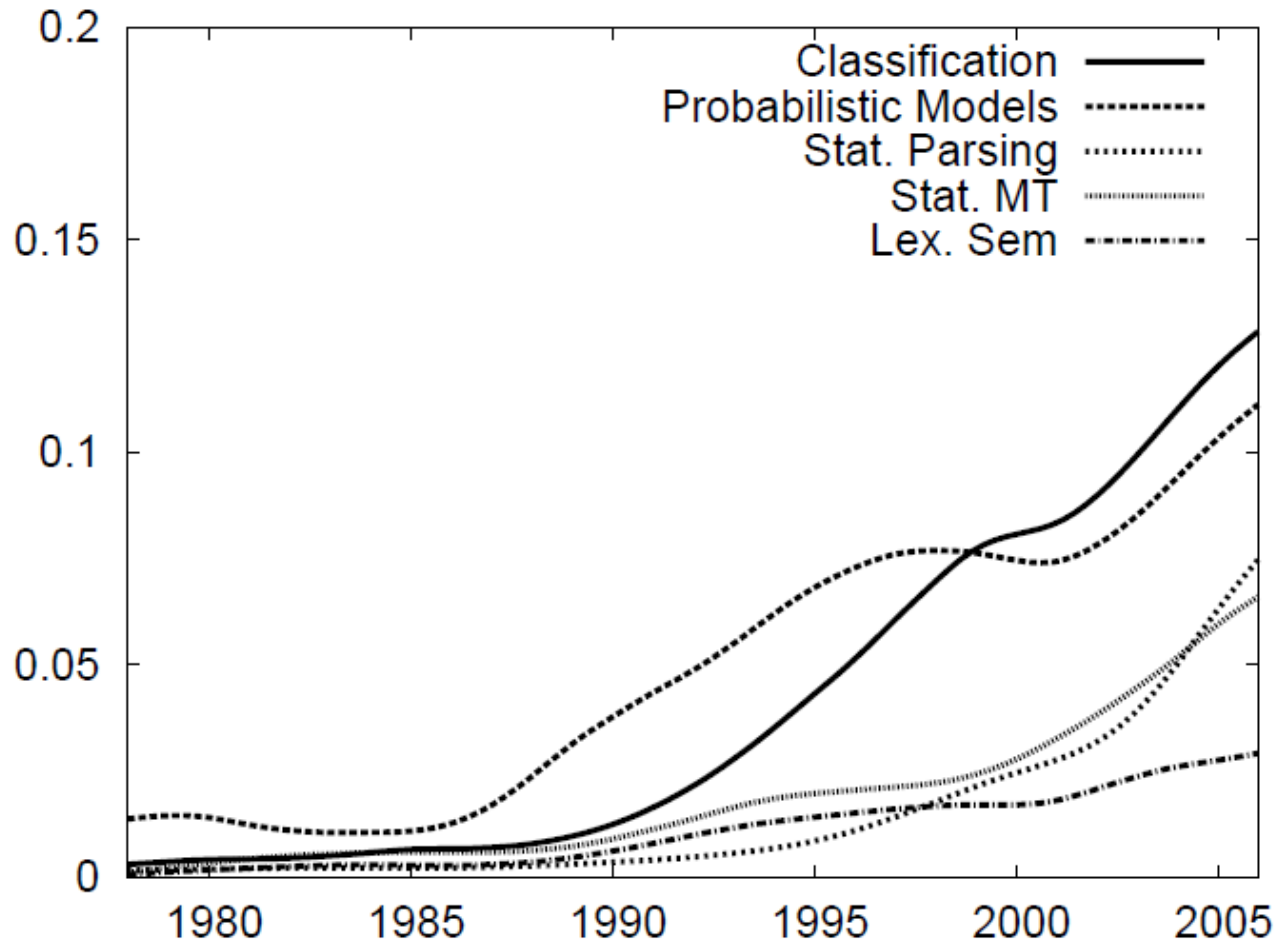
Parallel Tag Clouds [Collins et al.]



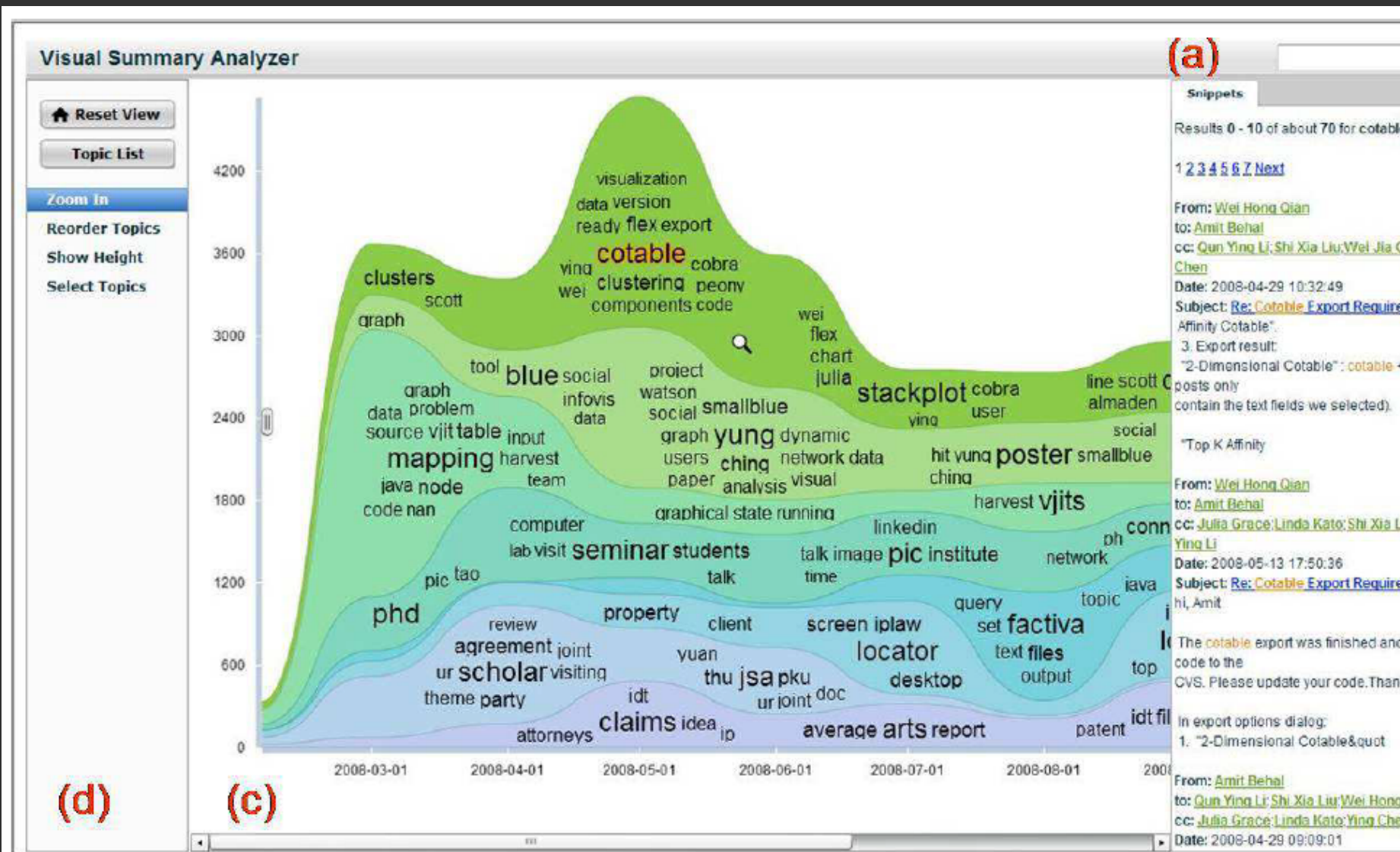
Theme River [Havre et al.]

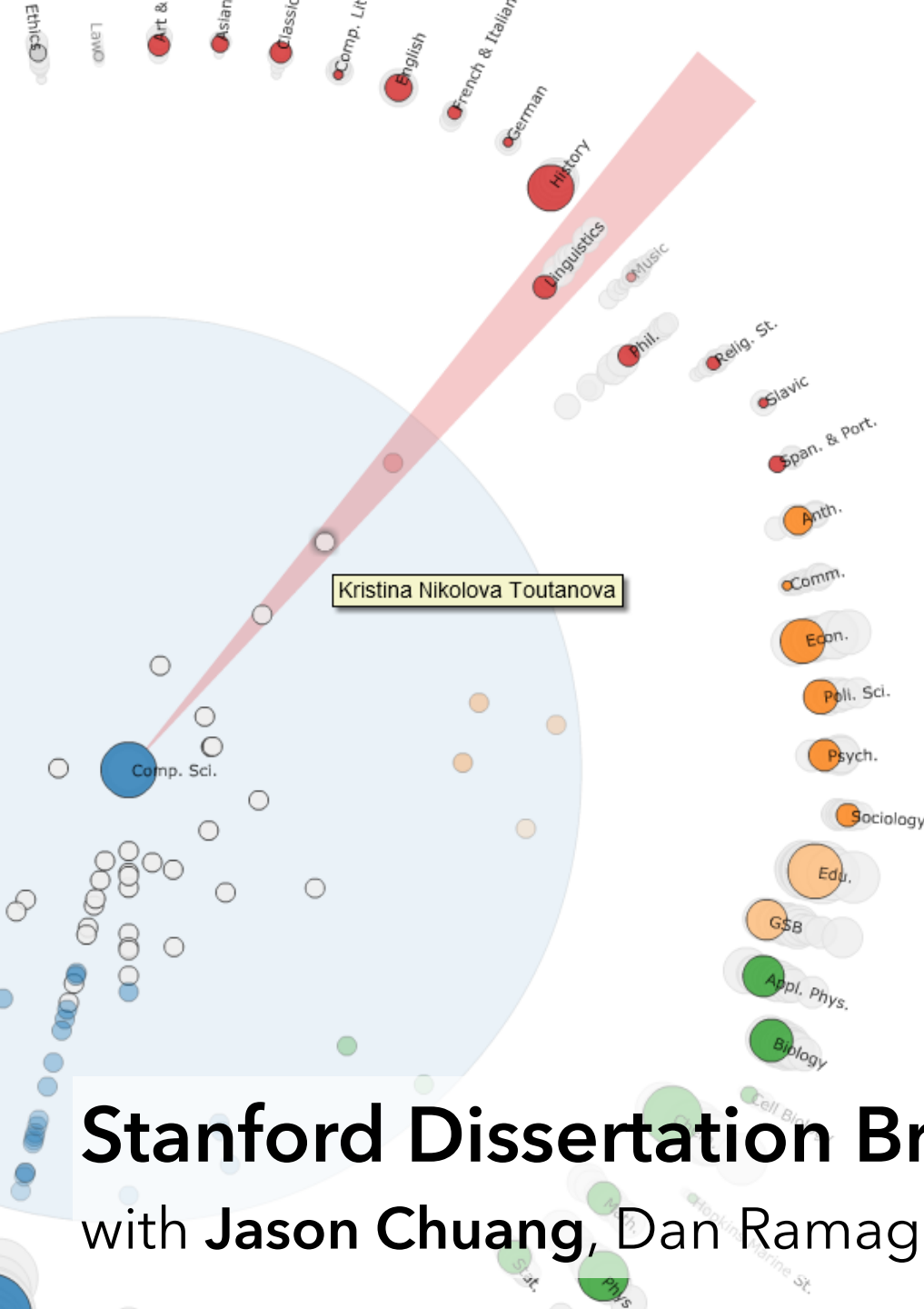


History of Comp. Ling. [Hall et al.]



Tiara [Wei et al.]





Effective statistical models for syntactic and semantic disambiguation

Student: Kristina Nikolova Toutanova
Advisor: Christopher D. Manning

Computer Science (2005)

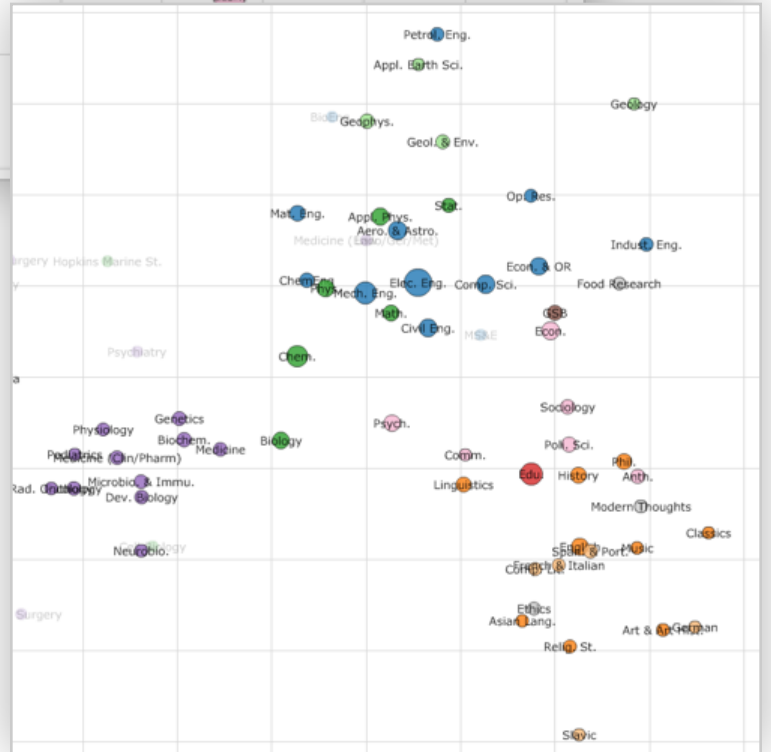
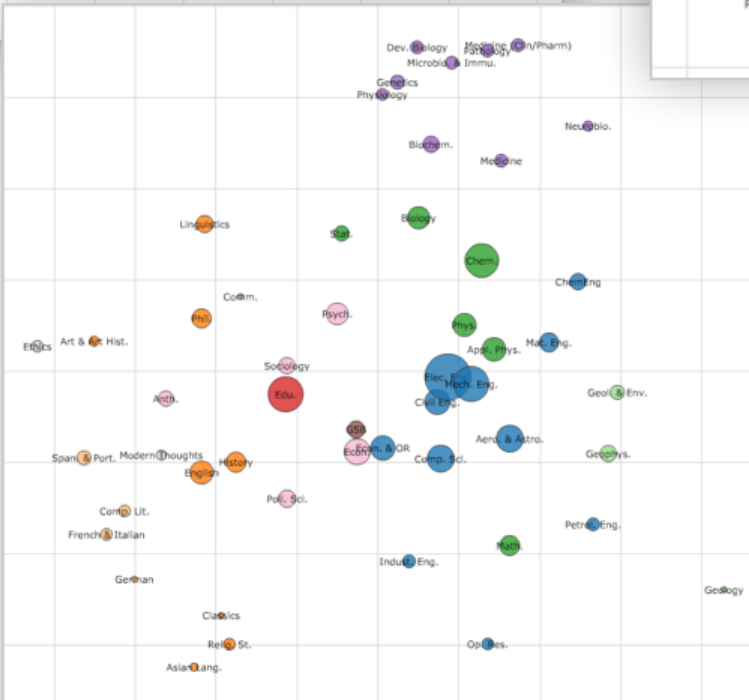
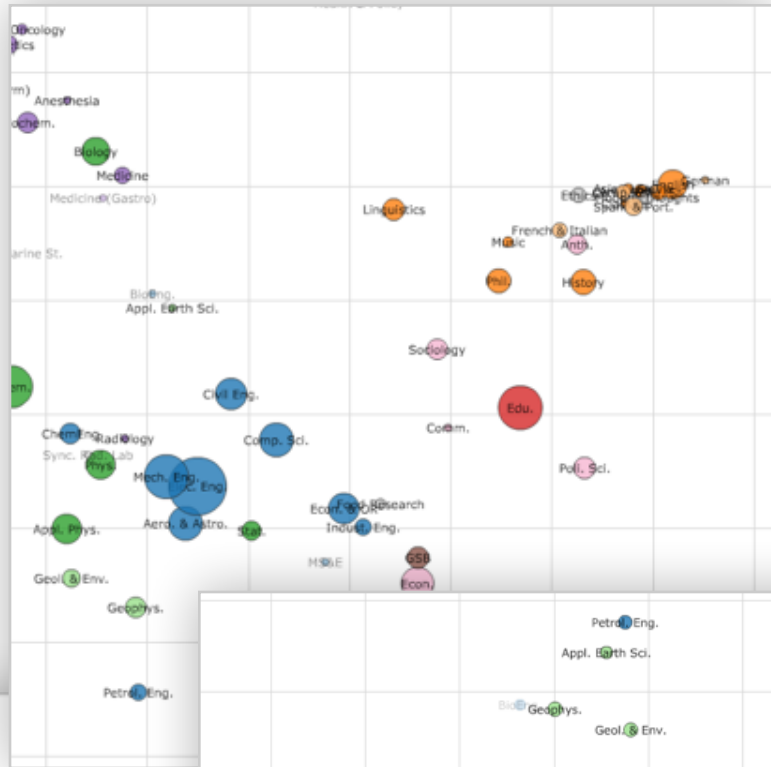
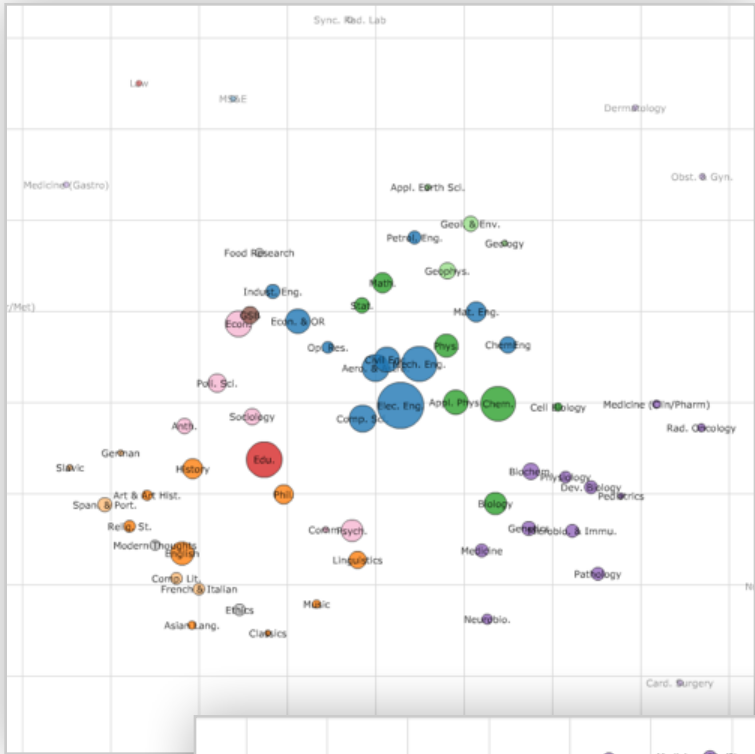
Keywords: Syntactic, Semantic, Tree kernels, Parsing

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.

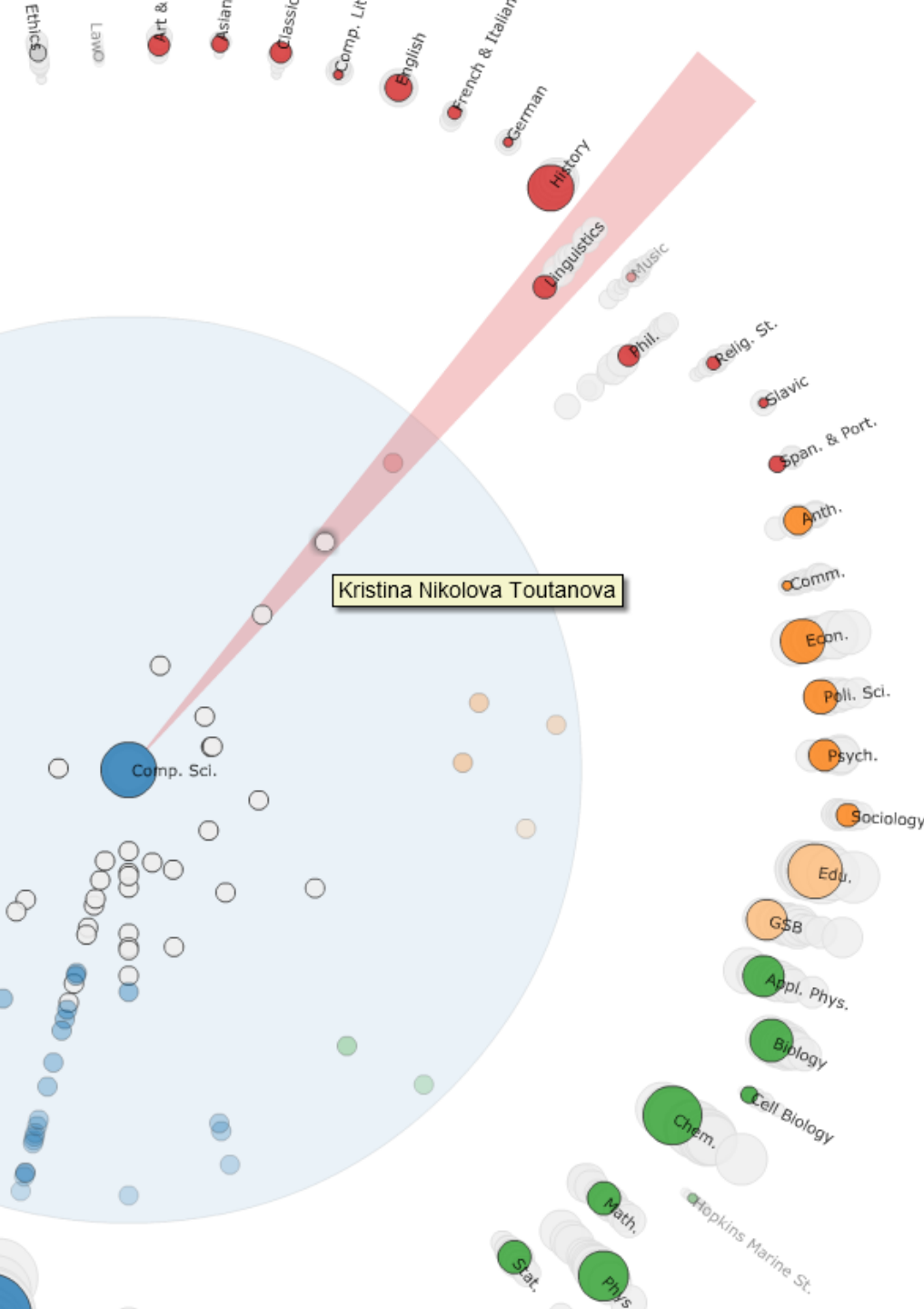
Stanford Dissertation Browser

with Jason Chuang, Dan Ramage & Christopher Manning





Oh, the humanities!



Effective statistical models for syntactic and semantic disambiguation

Student: Kristina Nikolova Toutanova

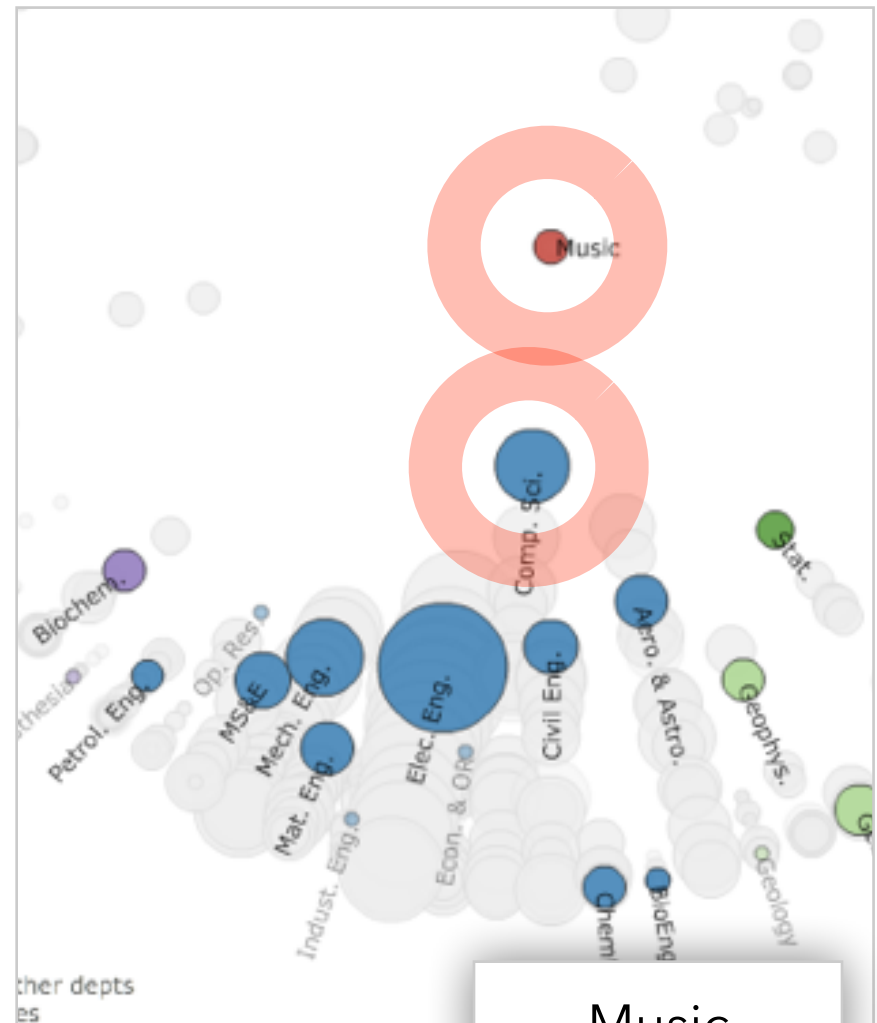
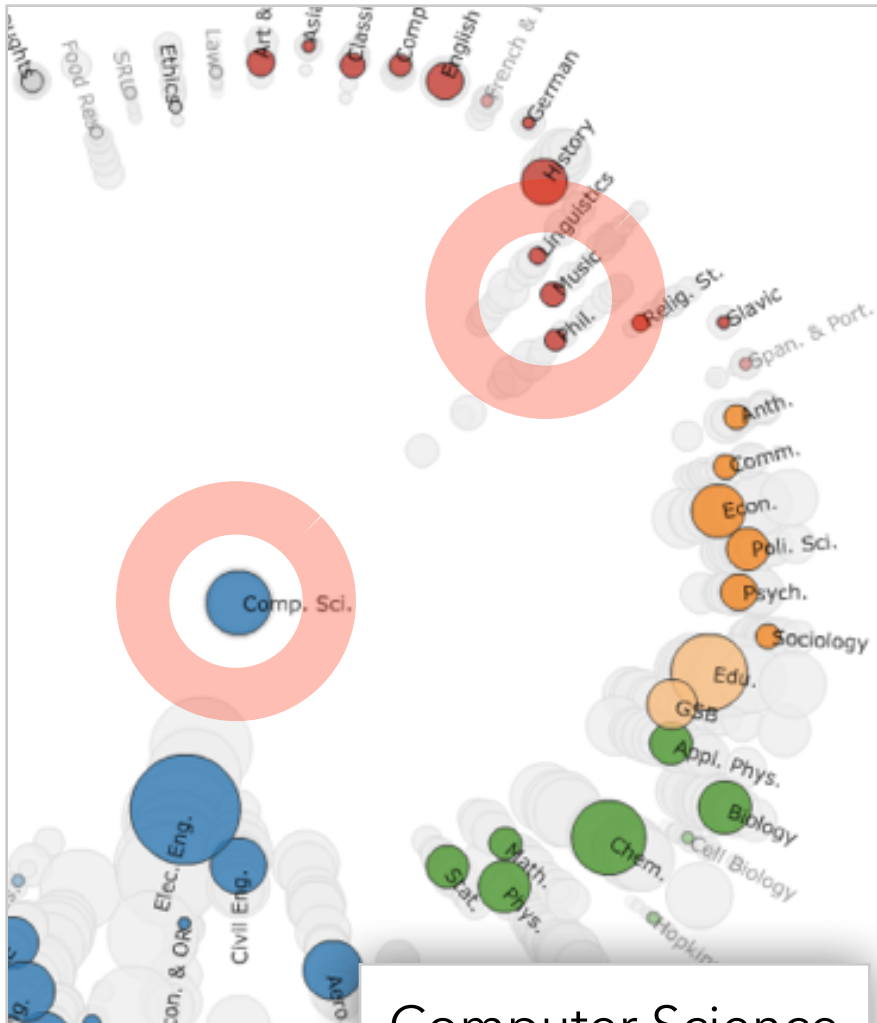
Advisor: Christopher D. Manning

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.



“Word Borrowing” via Labeled LDA

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

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