

CS448B :: 17 Nov 2011

# Text Visualization



Jason Chuang Stanford University

## Why visualize text?

## 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

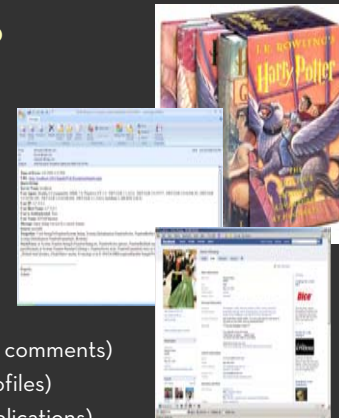
## What is text data?

### Documents

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

### Collection of documents

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



## 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.

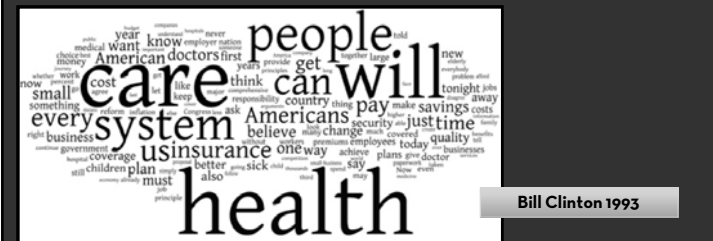
But we did not come here just to clean up crises. We came to build a future. So tonight, I return to speak to all of you

## Tag Clouds: Word Count

President Obama's Health Care Speech to Congress [New York Times]



[economix.blogs.nytimes.com/2009/09/09/obama-in-09-vs-clinton-in-93](http://economix.blogs.nytimes.com/2009/09/09/obama-in-09-vs-clinton-in-93)



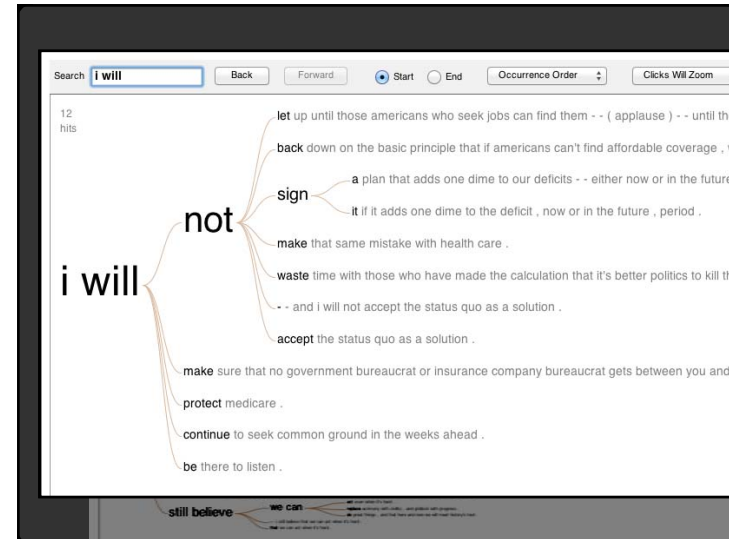
Bill Clinton 1993



Barack Obama 2009

[economix.blogs.nytimes.com/2009/09/09/obama-in-09-vs-clinton-in-93](http://economix.blogs.nytimes.com/2009/09/09/obama-in-09-vs-clinton-in-93)

## WordTree: Word Sequences



## A Double Gulf of Evaluation

Many (most?) 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?

## Challenges of Text Visualization

- 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 are (not) nominal?

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

## Tips: Tokenization and 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

## WordCount (Harris 2004)

<http://wordcount.org>



## Keyword Weighting

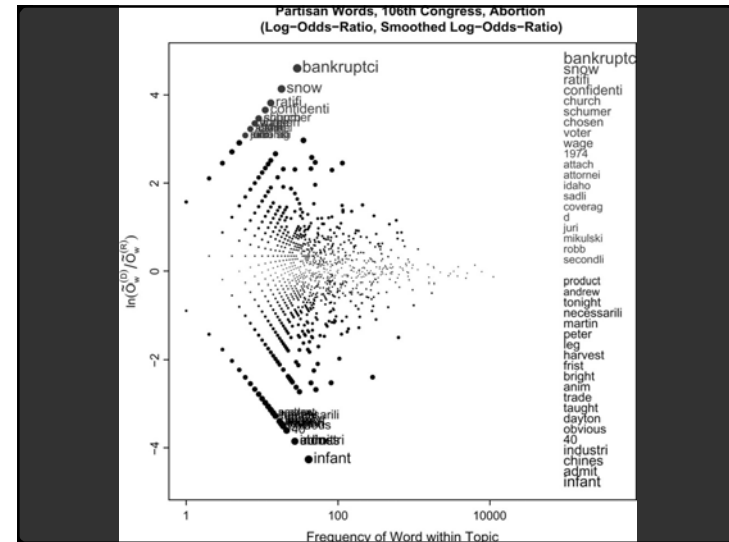
### Term Frequency

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

### 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

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

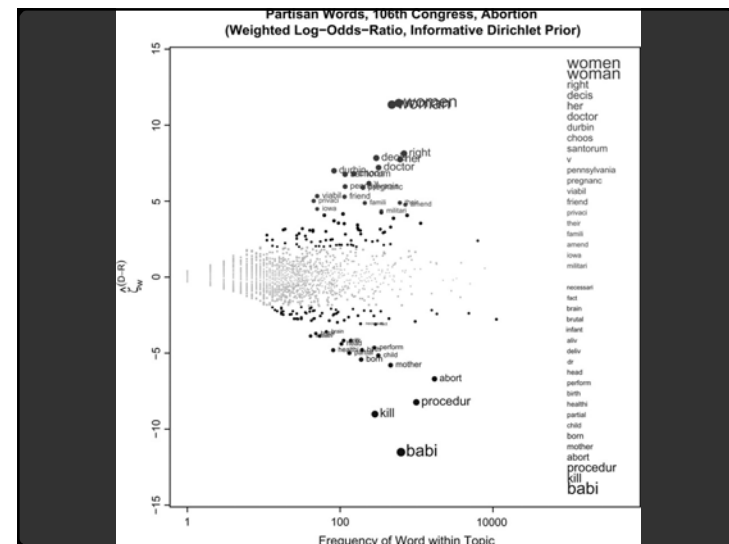
$$df_t = \# \text{ docs containing } t; N = \# \text{ of docs}$$

### G<sup>2</sup>: Probability of different word frequency

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

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

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



## 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

A “bag of words” ignores additional information

- Grammar / part-of-speech
- Position within document
- Recognizable entities

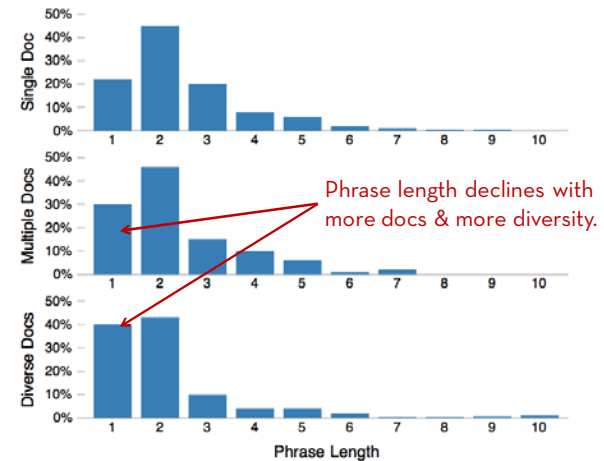
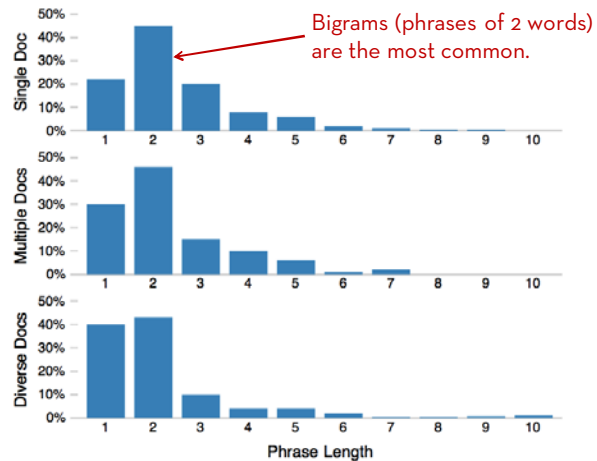
## How do people describe text?

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, Heer & Manning, 2010]



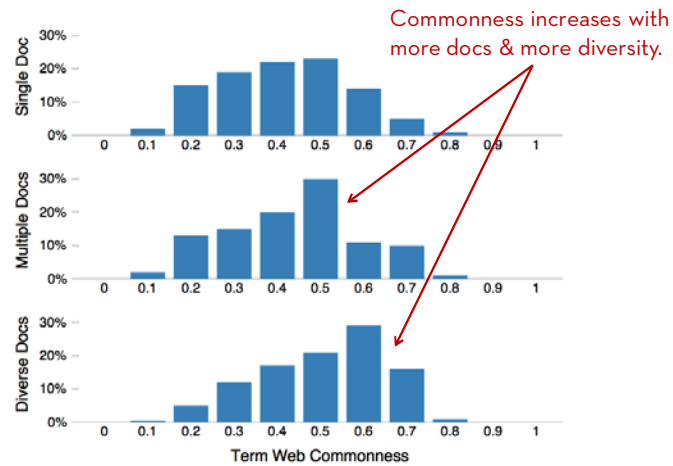
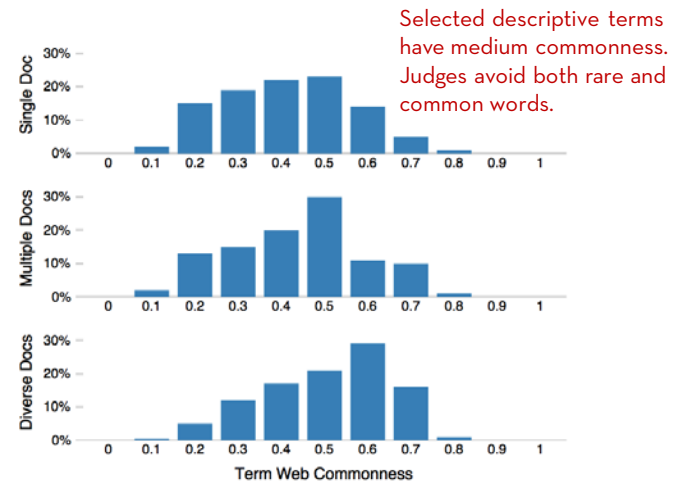


## Term Commonness

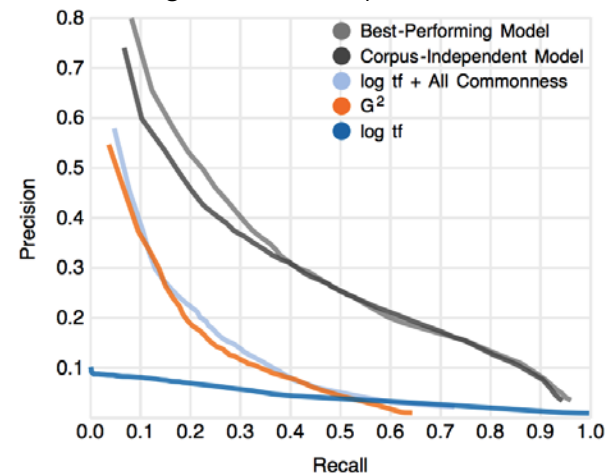
$$\log(\text{tf}_w) / \log(\text{tf}_{\text{the}})$$

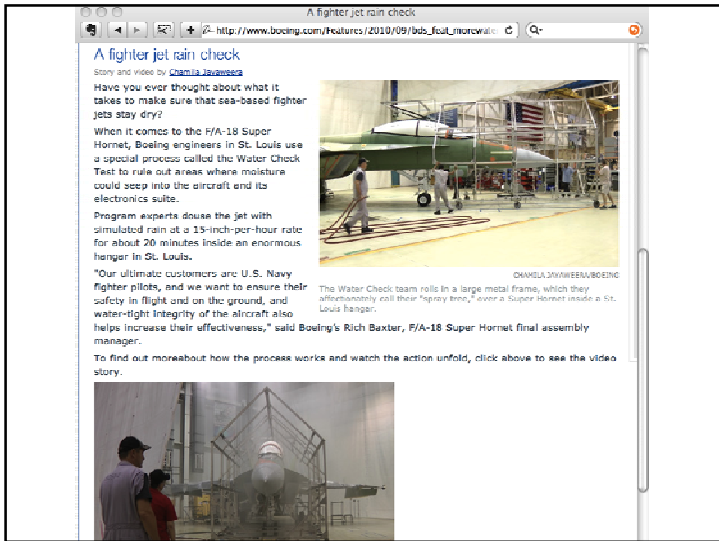
The normalized term frequency relative to the most frequent n-gram, e.g., the word “the”.

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



## Scoring Terms with Freq, Grammar & Position





$G^2$	Regression Model
fighter	Super Hornet
F/A	F/A -18
Hornet	fighter jet
Super	Boeing engineers
Boeing	special process
-18	rain check
rain	electronics suite
St.	Program experts
jet	simulated rain
Louis	ultimate customers
15-inch-per-hour	enormous hangar
douse	water-tight integrity
hangar	Rich Baxter
water-tight	15-inch-per-hour rate
Check	video story
Baxter	aircraft
sea-based	U.S. Navy fighter pilots
aircraft	Super Hornet final assembly manager
Rich	
seep	
click	
Navy	
sure	
Water	
moisture	
watch	
enormous	
dry	

## Yelp: Review Spotlight [Yatani 2011]

'09 amazing around baked bar bass best chef delicious eat elite everything favorite fish food fresh going hamachi hawaiian hour line love mango minutes mussels name night nighn order people ~~really~~ restaurant roll expensive or cheap? sake salmon sea seated service spicy stars size **sushi** table drink turn **wait** waitress worth

"long wait" or "no wait"? what type of sushi roll?

## Yelp: Review Spotlight [Yatani 2011]

'09 amazing around baked bar bass best chef delicious eat elite everything favorite fish food fresh going hamachi hawaiian hour line love mango minutes mussels name night nighn order people ~~really~~ restaurant roll expensive or cheap? sake salmon sea seated service spicy stars size **sushi** table drink turn **wait** waitress worth

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

baked mango long line hawaiian roll reasonable price

**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 .

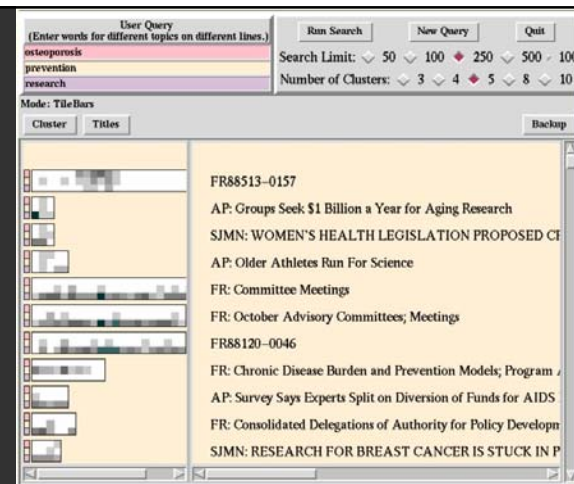
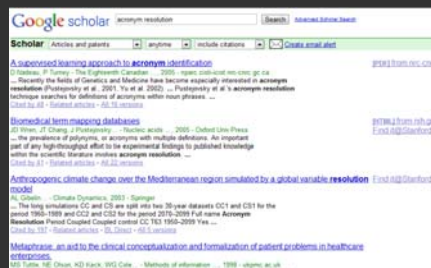
## Tips: Descriptive Keyphrases

- 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

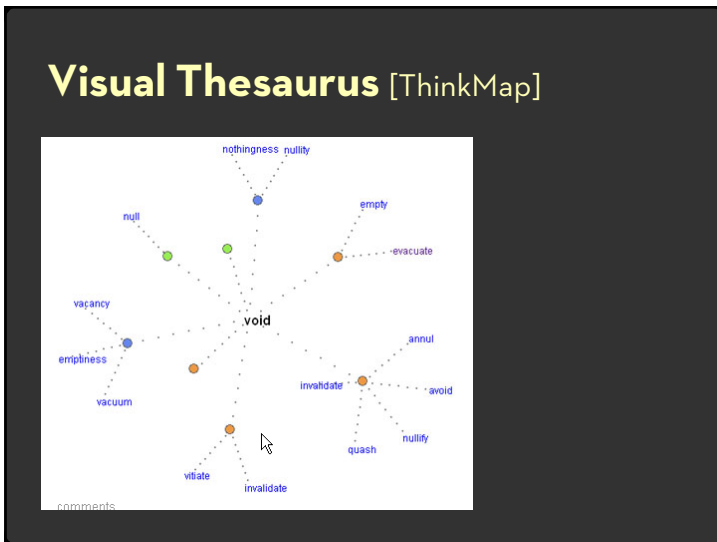
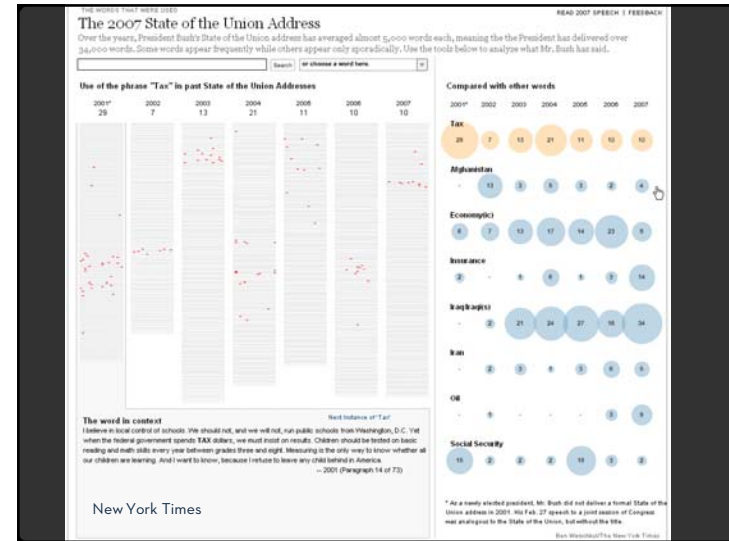
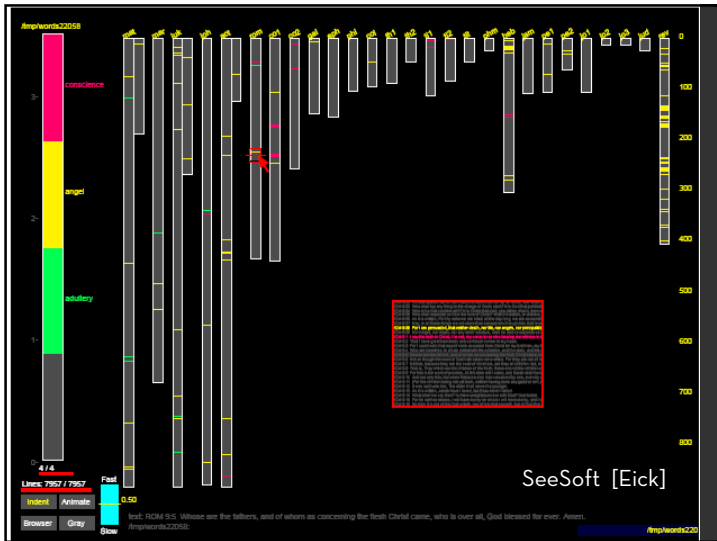
## Visualizing Document Content

## Information Retrieval

- Search for documents
  - Match query string with documents
- Contextualized search



TileBars [Hearst]



Google danse libre

Search About 16,400,000 results (0.21 seconds)

Everything The Academy of Danse Libre  
dansefire.org

Images Performing group for 19th and early 20th century social dance, based in California. Includes information about upcoming performances, and photos from past ...

Videos About Us M is for Mazurka  
For more information on performance scheduling and fees, or ... The Academy of Danse Libre presents "M is for Mazurka ..."

News Performances Danseurs  
First United Methodist Church, Palo Alto, Wednesday Night ... Previous to joining Danse Libre and performing in Stanford's ...

Shopping Auditions Repertoire  
Stanford, CA Danse Libre brings history to life, by choreographing and ... Danse Libre brings the vivacious atmosphere of the ...

Change location More results from dansefire.org >

All results  
Sites with images  
More search tools

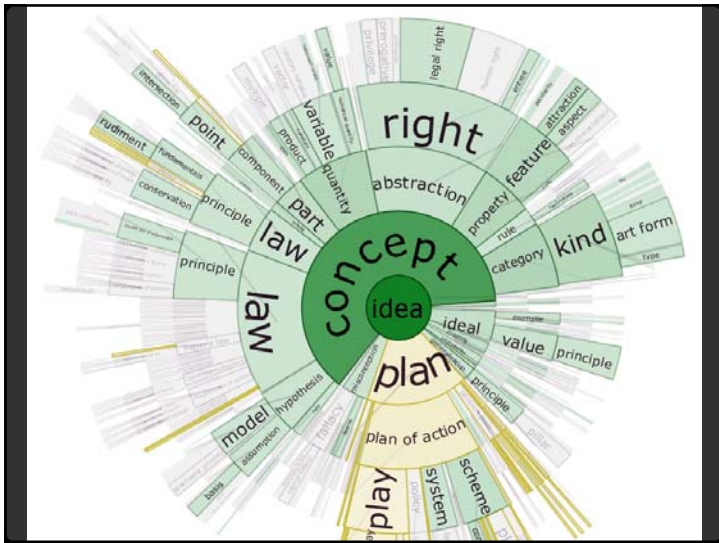
Friday Night Waltz  
www.fridaynightwaltz.com  
Danse Libre is performing tonight. \$12 for both the classes and the dance. \$8 students. No partner necessary. All Ages. Casual comfortable attire. Under 18 must ...  
Mike Brzozowski shared this on Blogger May 17, 2005

Elena Waltz (Academy of Danse Libre, 2006-12-08) - YouTube  
www.youtube.com/watch?v=MM2G0Rd...  
Mar 31, 2007 - 5 min - Uploaded by PanoramaSplash  
The Academy of Danse Libre performs a vintage dance called Elena Waltz at the weekly social dance party ...

Danse Libre - Brahms Mazurka Sextille at Stanford Viennese Ball ...  
www.youtube.com/watch?v=MM2G0Rd...  
Apr 3, 2011 - 2 min - Uploaded by academyofdansefire  
The Academy of Danse Libre performs the Brahms Mazurka Sextille at the Stanford Viennese Ball on ...

More videos for danse libre >

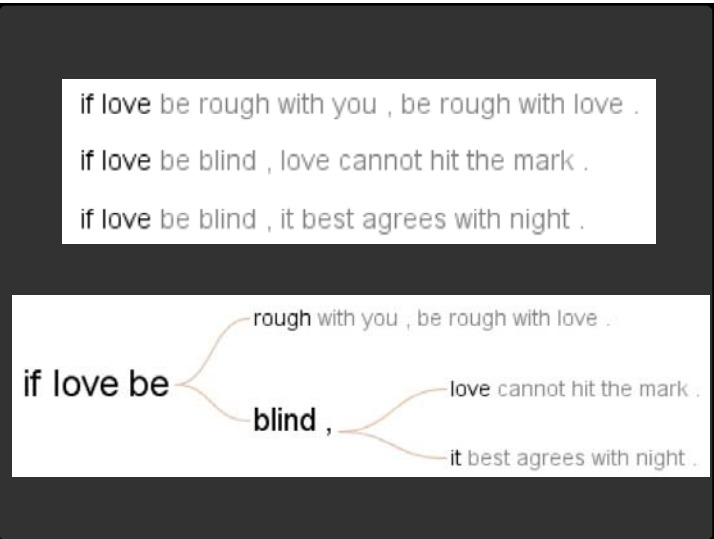
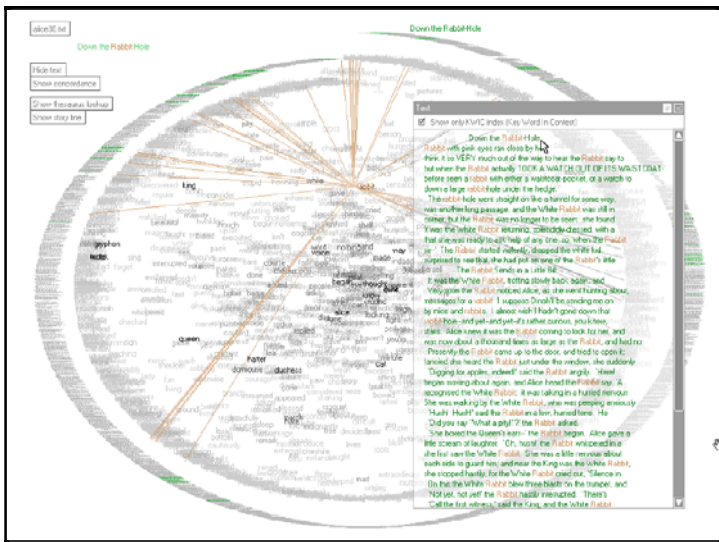
Danse Libre at Bin Dance



# Concordance

What is the common local context of a term?

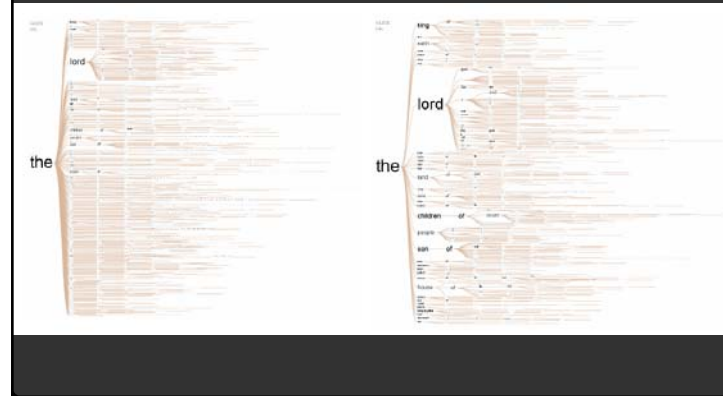
Headword	No.	Context	Word	Reference
HEAR	15	That my own heart	drifts and cries, having no...	Deep Analysis
HEARD	9	By the shout of the heart	continually at work	And the waves
HEARING	7	Nothing to adapt the skill of the heart	to, still	And the waves
HEARS	3	The bread, the loaf of it, it is my own heart		Travellers
HEARSE	1	Because I follow it to my own heart		Many famous
HEART	21	My heart	is Soking like the sun,	I am washed / sharpened to a candle co...
HEART'S	2	The vague heart	by looking out of date	Lines on a Yo
HEART-SHAPED	1	Contract my heart	to put aside the theft	Home is so So
HEARTH	1	Having no heart	out in the Gents	Essential leav...
HEARTS	7	And the boy putting his 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 life schema of the heart	has taken,	Time and Spo...
HEATH	1	They gambled heart		A Stone Chan...
HEATS	1	How should they sweep the girl clean...	heart	I see a girl da
HEAVE	1	Hands that the heart	can govern	Heaviest of fo...
HEAVER	4	For the heart	to be loveless, and as col...	Down
HEAVER-HOLDING	1	With the unguessed heart	riding	One man wak...
HEAVER-THAIR...	1	If hands could free you, heart		If hands could
HEAVEN	2	That overflows the heart		Four away th...



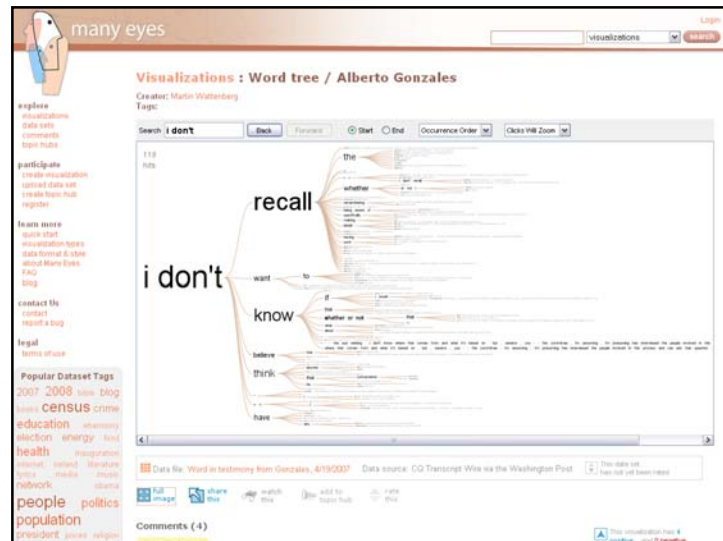
## WordTree (Wattenberg et al)



## Filter infrequent runs



## Recurrent themes in speech



## 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>

## Phrase Nets [van Ham et al]

Look for specific linking patterns in the text:

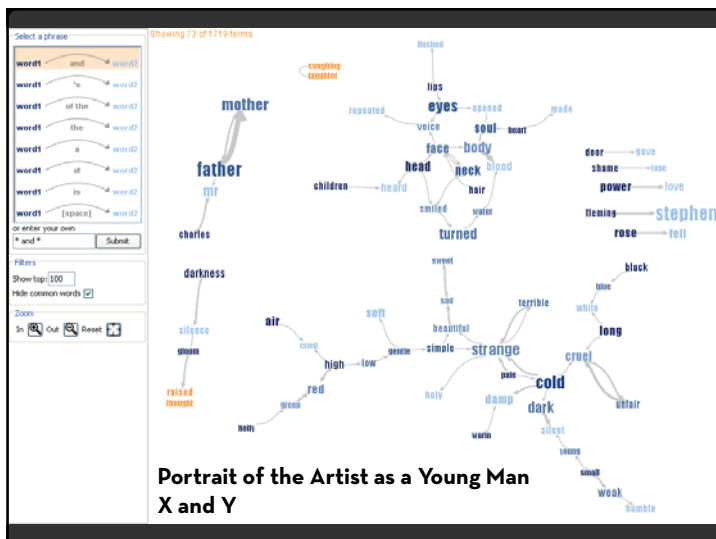
'A and B', 'A at B', 'A of B', etc

Could be output of regexp or parser.

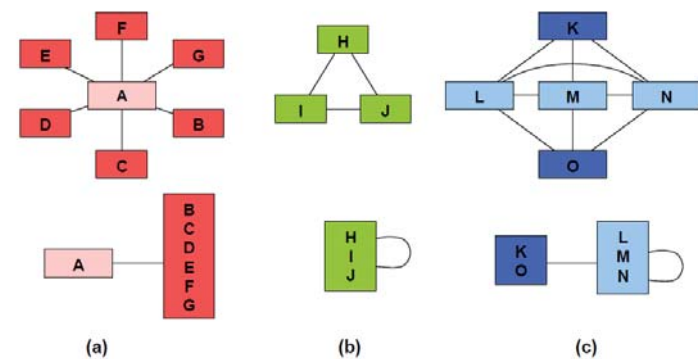
Visualize extracted patterns in a node-link view

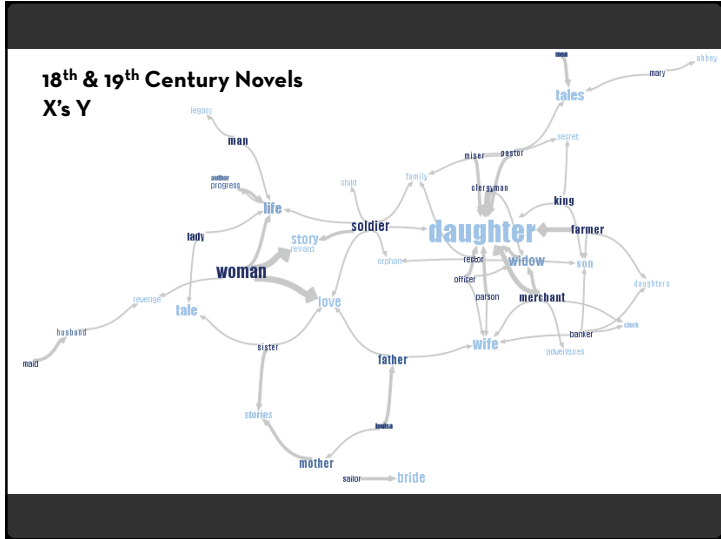
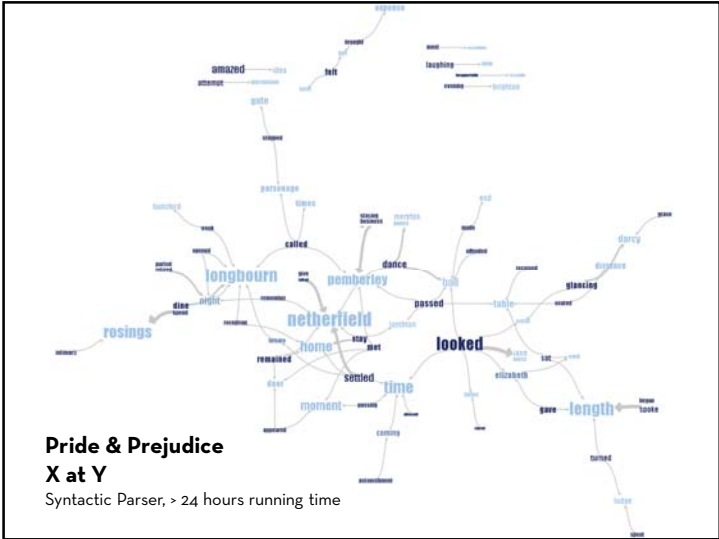
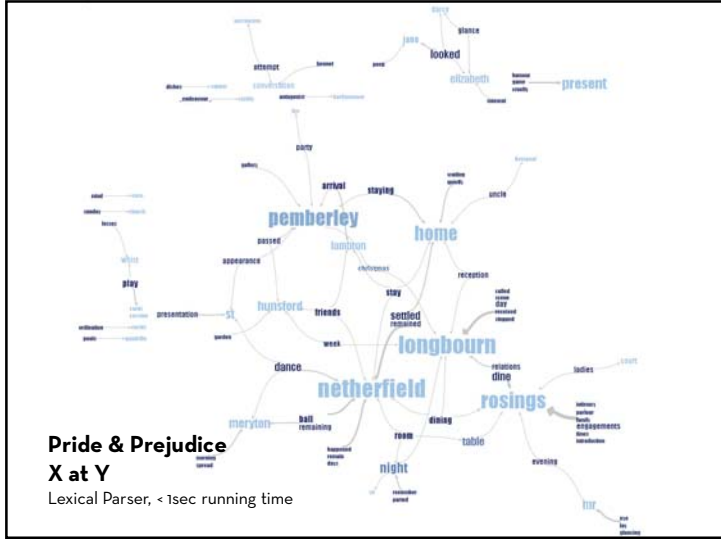
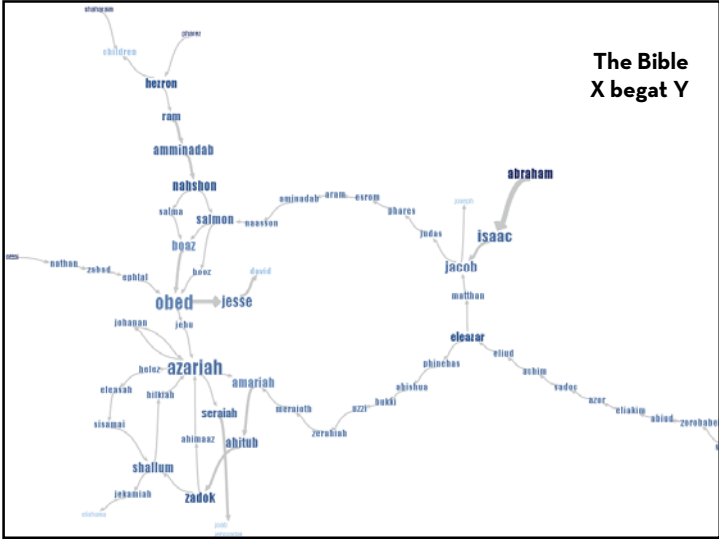
Occurrences → Node size

Pattern position → Edge direction



## Node Grouping









## Final Project

### Design a new visualization technique or system

Many options: new system, interaction technique, design study  
6-8 page paper in conference paper format  
2 Project Presentations

### Schedule

Project Proposal: **Tuesday, Nov 15** (end of day)

Initial Presentation: **Tuesday, Nov 29**

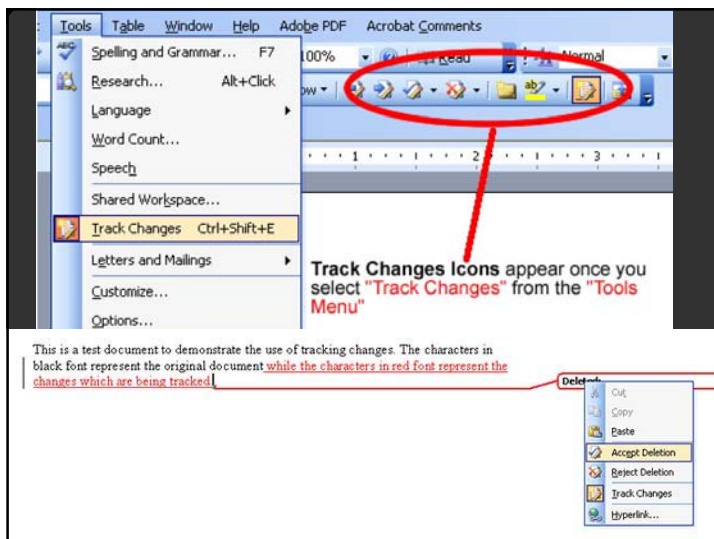
Poster Presentation: **Tuesday, Dec 13** (5-7pm)

Final Papers: **Thursday, Dec 15** (end of day)

### Logistics

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

## Evolving Documents



Track Changes Icons appear once you select "Track Changes" from the "Tools Menu"

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

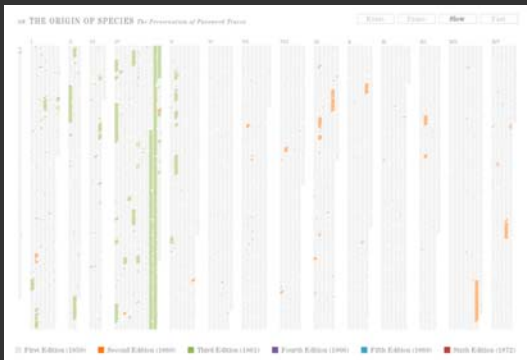
### Chocolate

Revision history

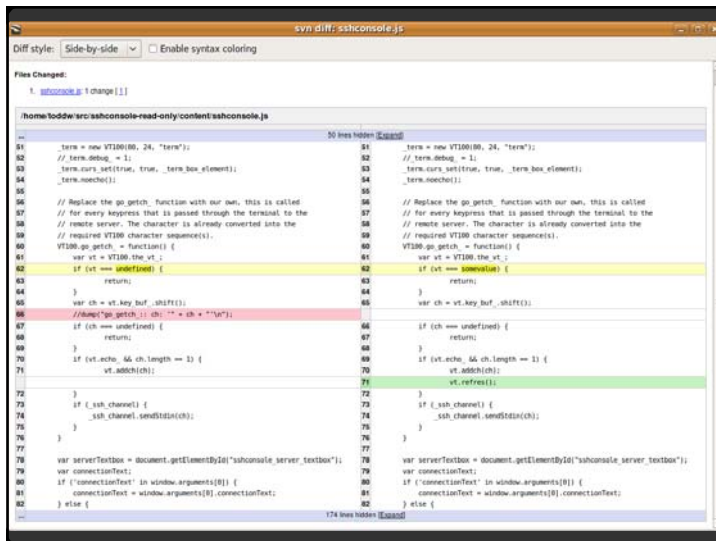
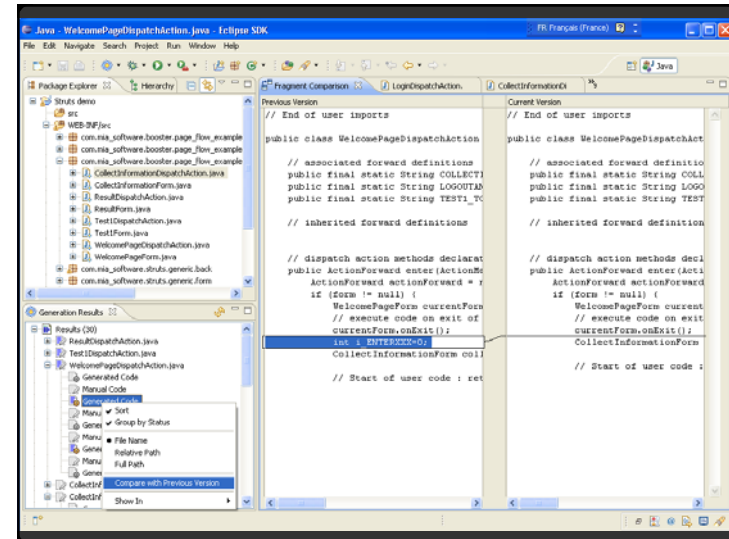
Legend: (cur) = difference with current version, (last) = difference with preceding version, M = minor edit

- (cur) (last) . 12.01, 20 Aug 2003 . [Dysprosia](#) (*heaten 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*)

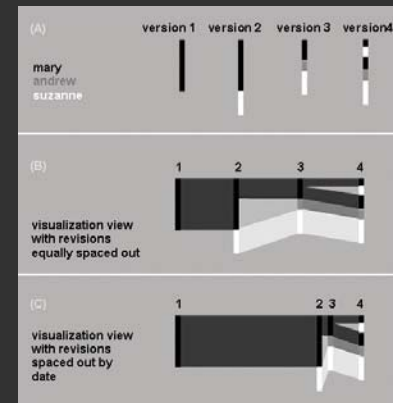
# Animated Traces [Ben Fry]

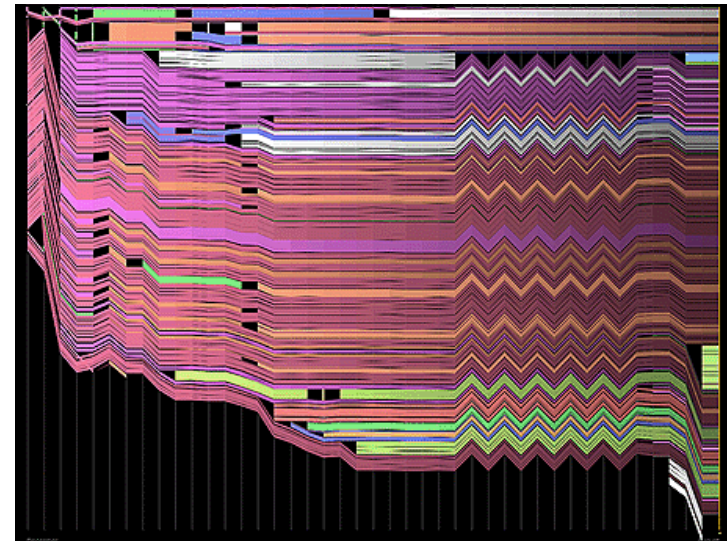
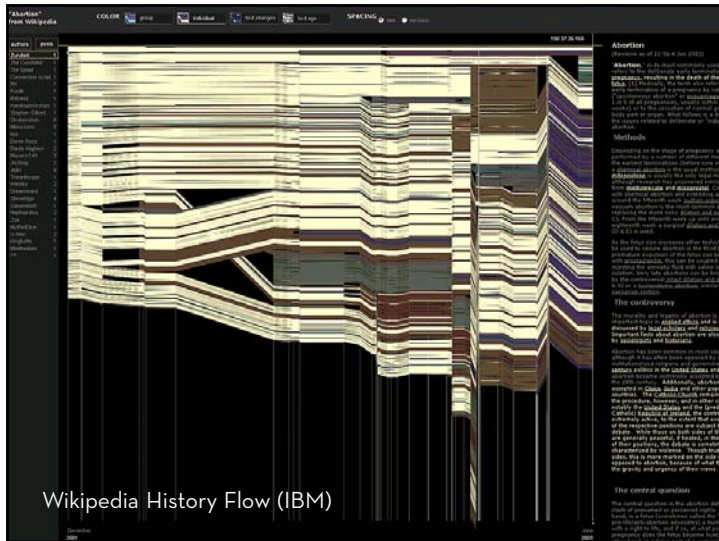


<http://benfry.com/traces/>



# History Flow (Viégas et al)





## Tips: Evolving documents

- High-level understanding
- Provide context
  - Show text within source document
  - Cross reference with other dimensions

## Visualizing Conversation

# Visualizing Conversation

Many dimensions to consider:

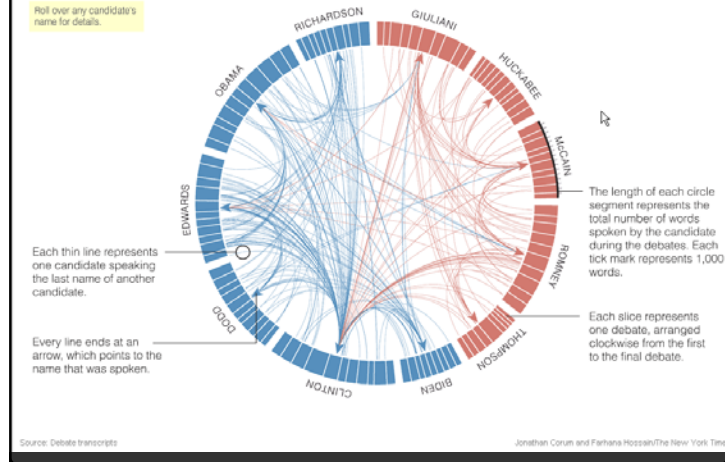
- Who (senders, receivers)
- What (the content of communication)
- When (temporal patterns)

Interesting cross-products:

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

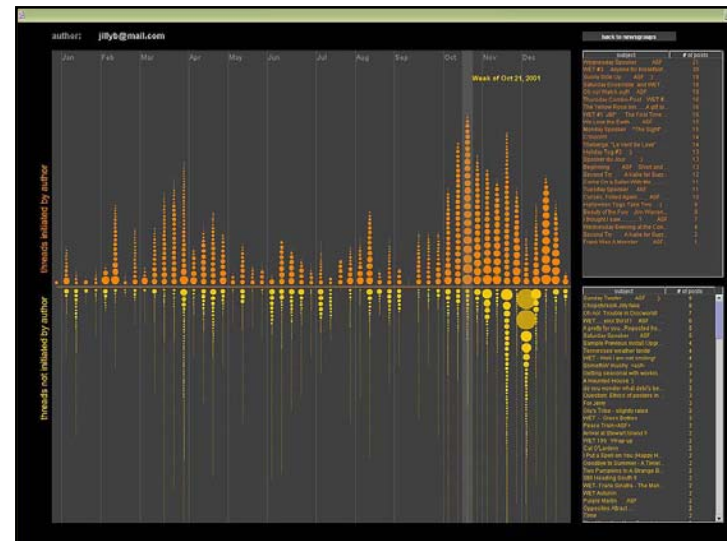
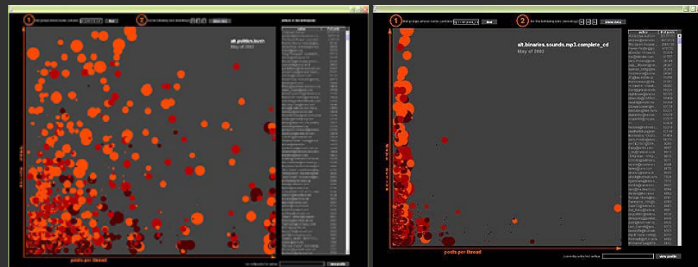
## Naming Names

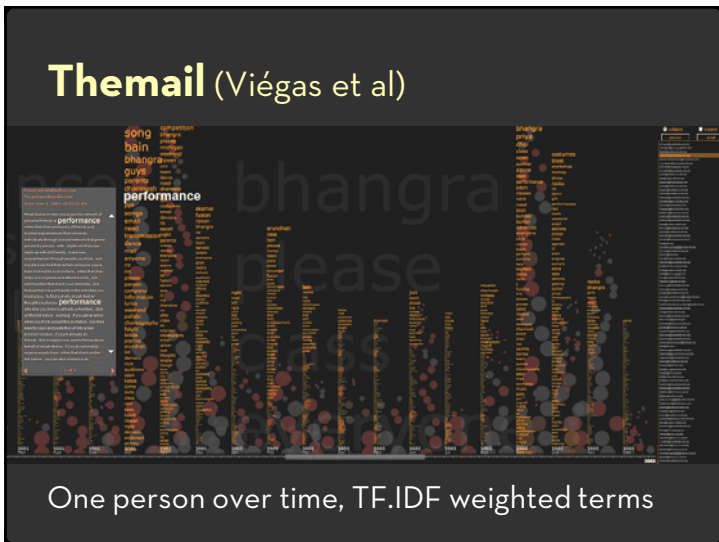
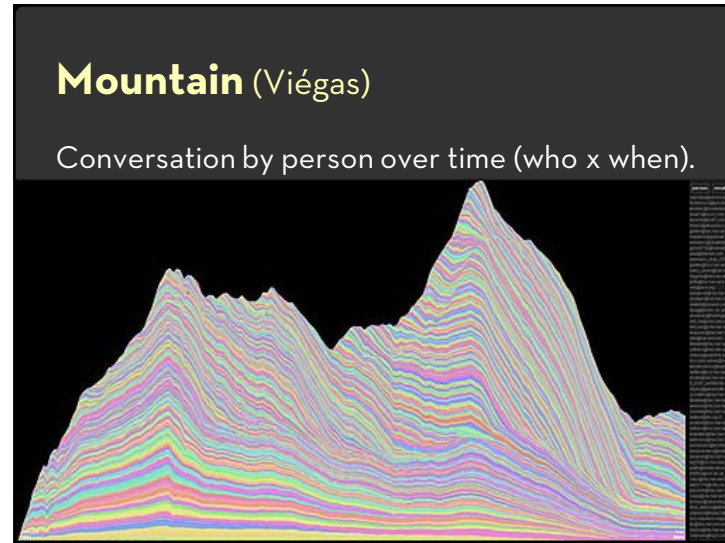
Names used by major presidential candidates in the series of Democratic and Republican debates leading up to the Iowa caucuses.

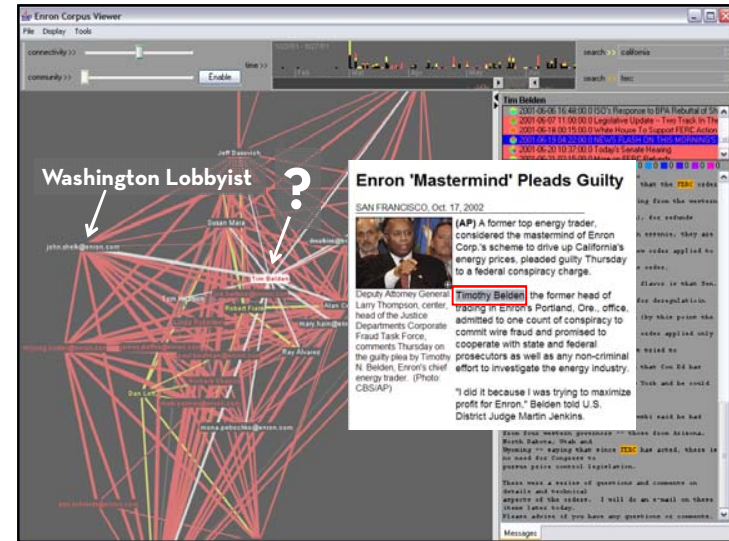
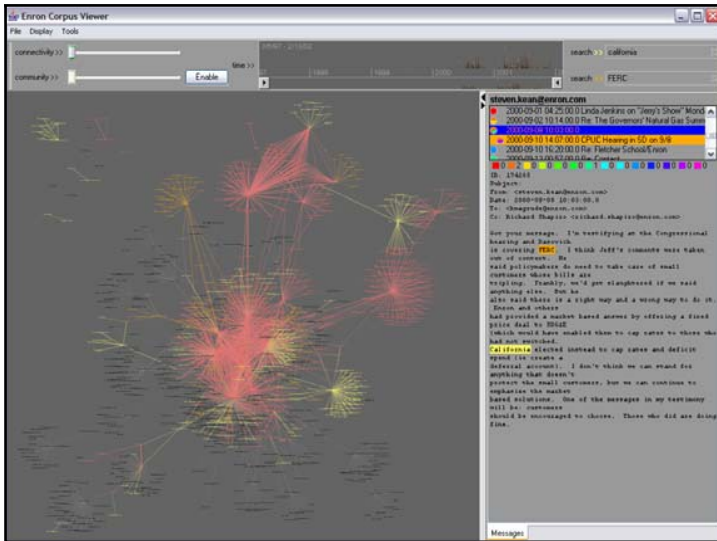


# Usenet Visualization (Viégas & Smith)

Show correspondence patterns in text forums  
Initiate vs. reply; size and duration of discussion



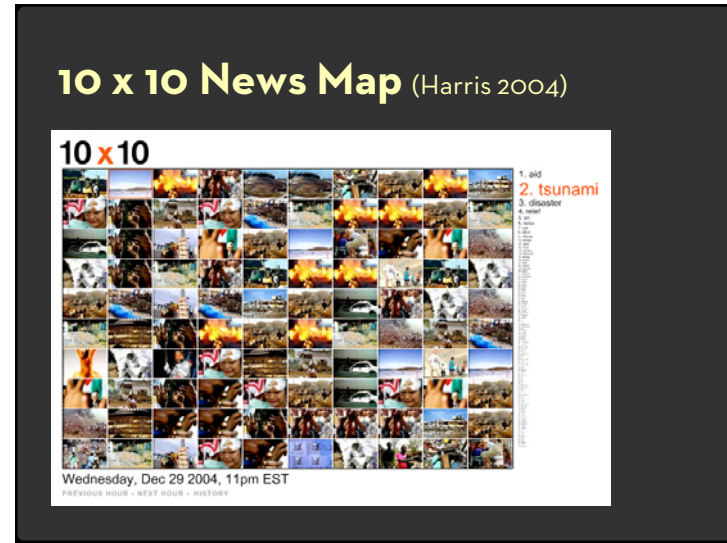
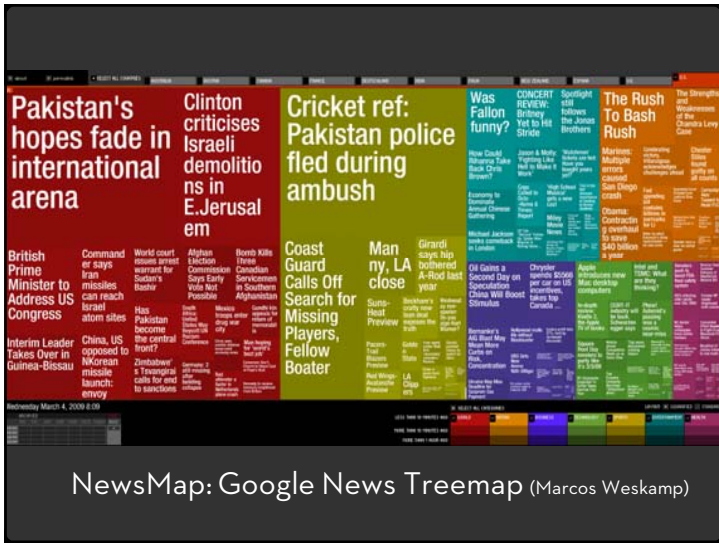




## Tips: Conversations

- Understand your units of analysis
  - Extract entities and relationships relevant to analysis task.
  - Cross-reference with other data dimensions.

## Visualizing Document Collections



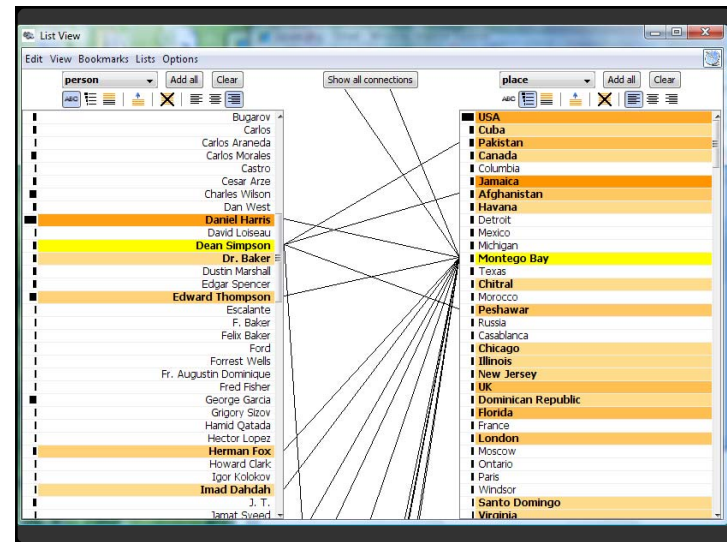
## Named Entity Recognition

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 a small window of text?





## Doc. Similarity & Clustering

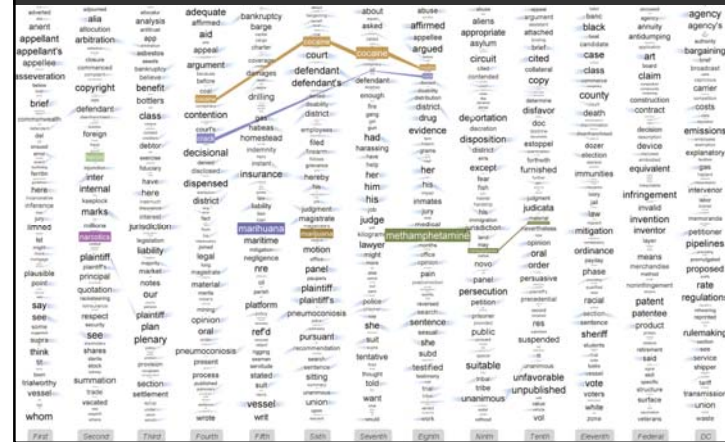
In vector model, compute distance among docs

- For TF.IDF, typically cosine distance
- Similarity measure can be used to cluster

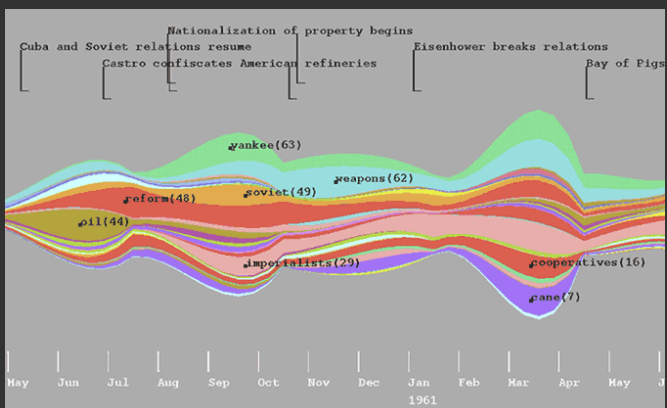
Topic modeling approaches

- 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

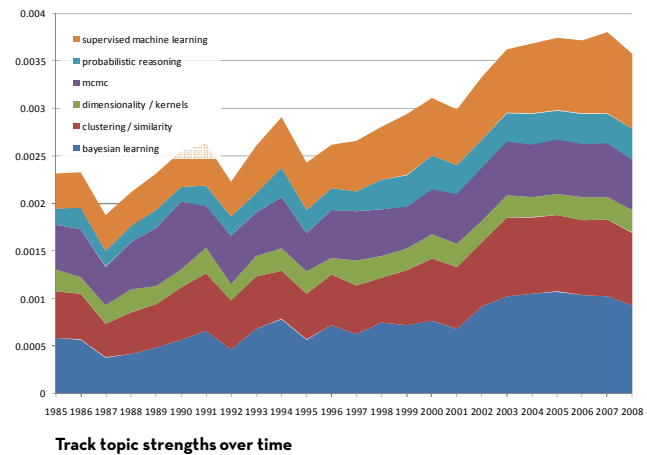
## Parallel Tag Clouds [Collins et al 09]



## ThemeRiver [Havre et al 99]



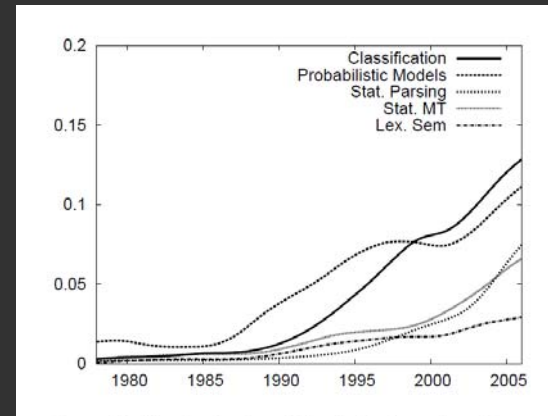
## Statistical Machine Learning in Pubmed



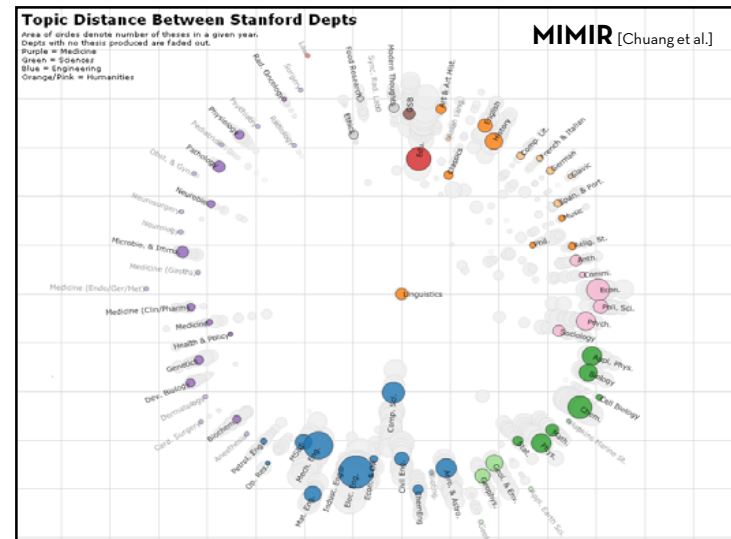
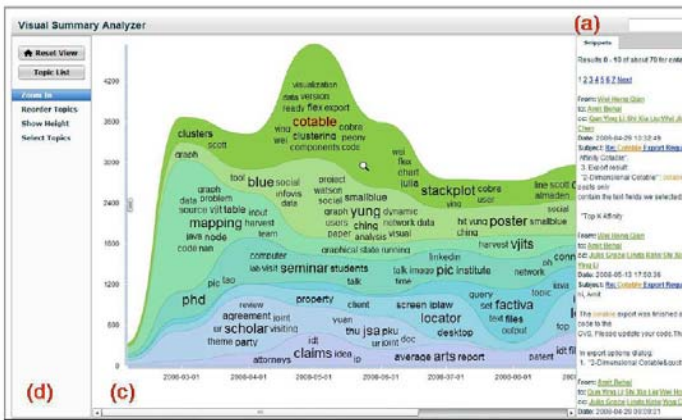
# Interpretation and Trust?

- Interpretable topics?
- Trust the topics?

# History of Comp Linguistics [Hall et al 06]



# TIARA [Wei et al. 09]



## Challenges of Text Visualization

- 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.

## Lessons for Text Visualization

- Align analysis task with appropriate model.
- Handle high dimensionality...
  - Semantically
    - Interpretation: Longer phrases
    - Restaurant reviews: Adjective-noun word pairs
    - Relationships: Word sequences, hierarchy, clustering, ...
    - Topic models: **with care**
  - Visually
    - Word position within document
    - High-level structures in document collection
    - Visual representation matching semantic relationships

## Lessons for Text Visualization

- Align analysis task with appropriate model.
- Provide context and semantics...
  - Apply appropriate text processing: stemming, named entities, etc.
  - Reverse stem for presentation
  - Show text within source document
  - Interaction to enable analysis cycle
  - Allow users to express contextual or domain knowledge
  - Cross-reference with other data dimensions