

G53FIV: Fundamentals of Information Visualization

Lecture 11: Visualizing Text and Documents

Ke Zhou
School of Computer Science
Ke.Zhou@nottingham.ac.uk

<https://moodle.nottingham.ac.uk/course/view.php?id=68644>

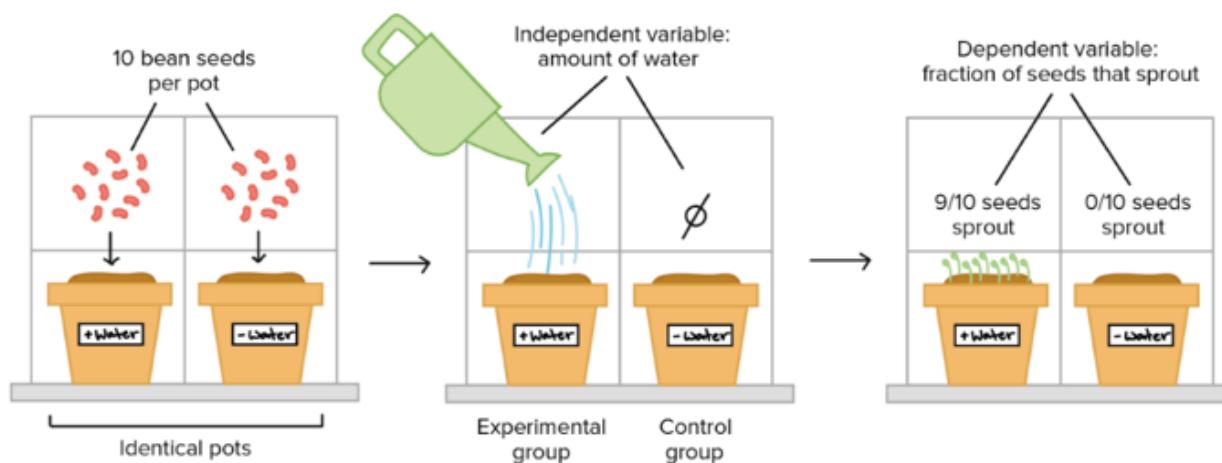
Last Lecture

Evaluation

Quantitative Methods: Controlled Experiments

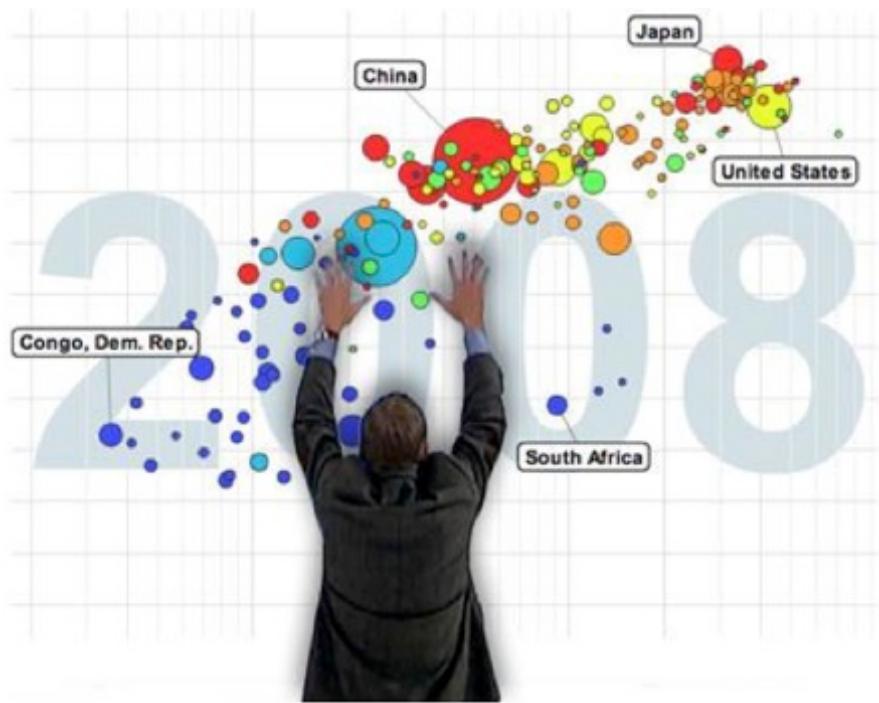
- Good for measuring performance or comparing multiple techniques
- What do we measure?
 - Performance, time, errors,

...



An Example: Controlled Experiment

- Run an experiment to evaluate three visualization strategies
 - Animation
 - Small multiples
 - Traces
- Especially interested in examining whether animated bubble charts are beneficial for analysis and presentation



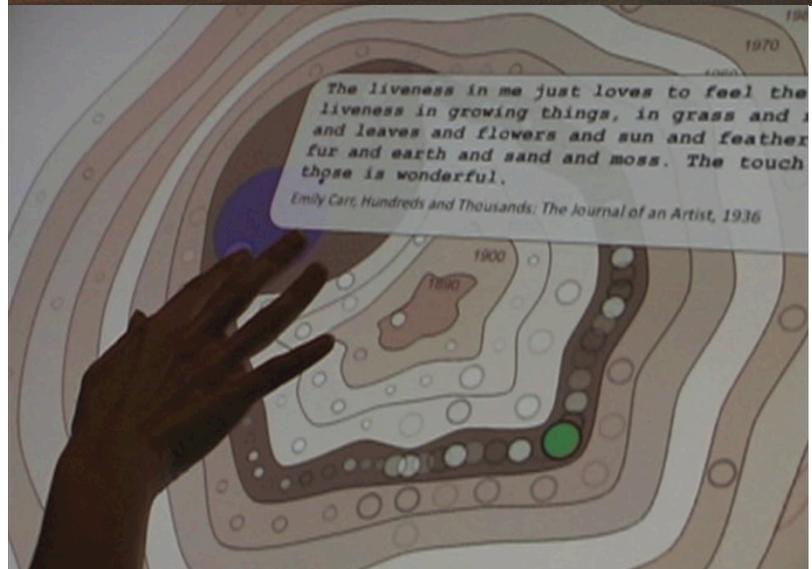
*Do you remember Hans Rosling's TED talk?
(Lecture 2)*

Qualitative Methods

- Types
 - Nested methods
 - Experimenter observation, think-aloud protocol, collecting participant opinions
 - Inspection evaluation methods
 - Heuristics to judge
- Observational context
 - In situ, laboratory, participatory
 - Contextual interviews is important

An Example: Subjective Assessments

- Evaluating a newly developed visualization system (EMDialog) at the museum (Emily Carr exhibit)
- Discourse visualization
 - Time
 - Context



Methodology vs. Desirable Features

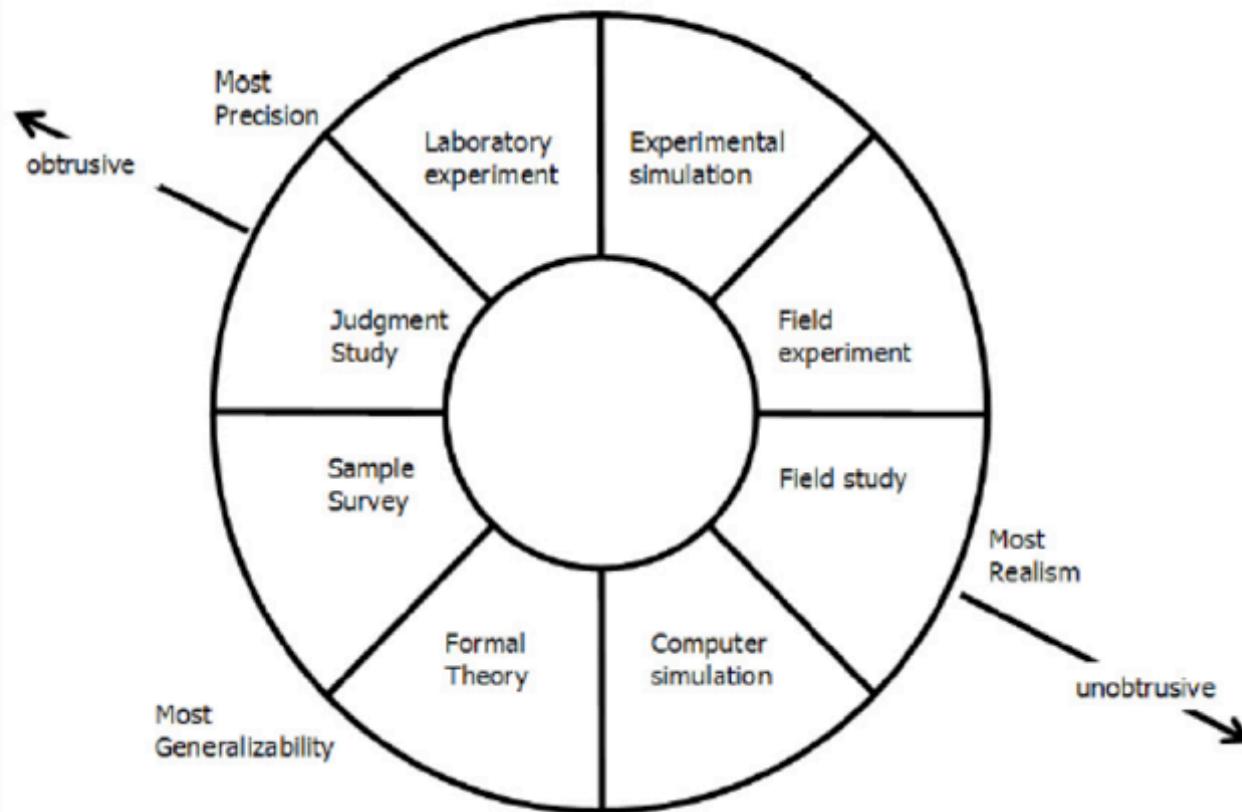


Fig. 1. Types of methodologies organized to show relationships to precision, generalizability and realism. (adapted, simplified from McGrath 1995)

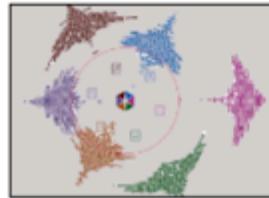
Overview



Visualizing text
Showing words,
phrases, and
sentences



- Content
- Context
- Relationship to others



Visualization for IR
Helping search



Why Visualize Text?

- What can information visualization provide to help users in understanding and gathering information from text and document collections?
- **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

Challenges

- High Dimensionality
 - Where possible use text to represent text...
 - ... which terms are the most descriptive?
- Context and 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.

Visualizing Text

How do we represent the words, phrases, and sentences in a document or set of documents?



Visualizing text

Showing words,
phrases, and
sentences



An Example

- Health care speech transcripts
 - Clinton in 1993
 - Obama in 2009
- What questions might you want to answer?
- What visualizations might help?

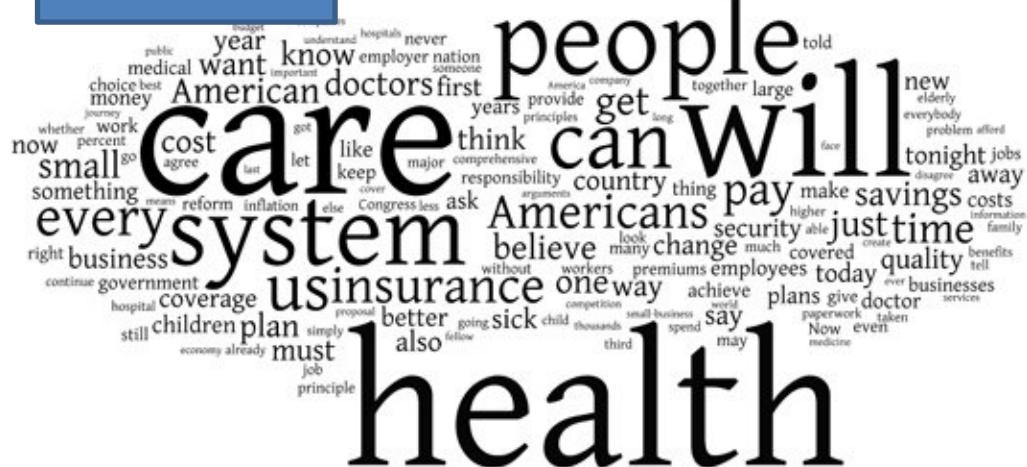
Bill Clinton on Health Care, 1993

Now we are in a time of profound change and opportunity – the end of the cold war, the information age, and the global economy have brought us both opportunity and hope and strife and uncertainty. Our purpose in this dynamic age must be to change – to make change our friend and not our enemy.

To achieve that goal, we must face all our challenges with confidence, with faith and with discipline, whether we're reducing the deficit, creating tomorrow's jobs and training our people to fill them, converting from a high-tech defense to a high-tech domestic economy, expanding trade, reinventing government, making our streets safer, or rewarding work over idleness, all of these challenges require us to change.

Tag Clouds: Word Count

Clinton 1993



Obama 2009

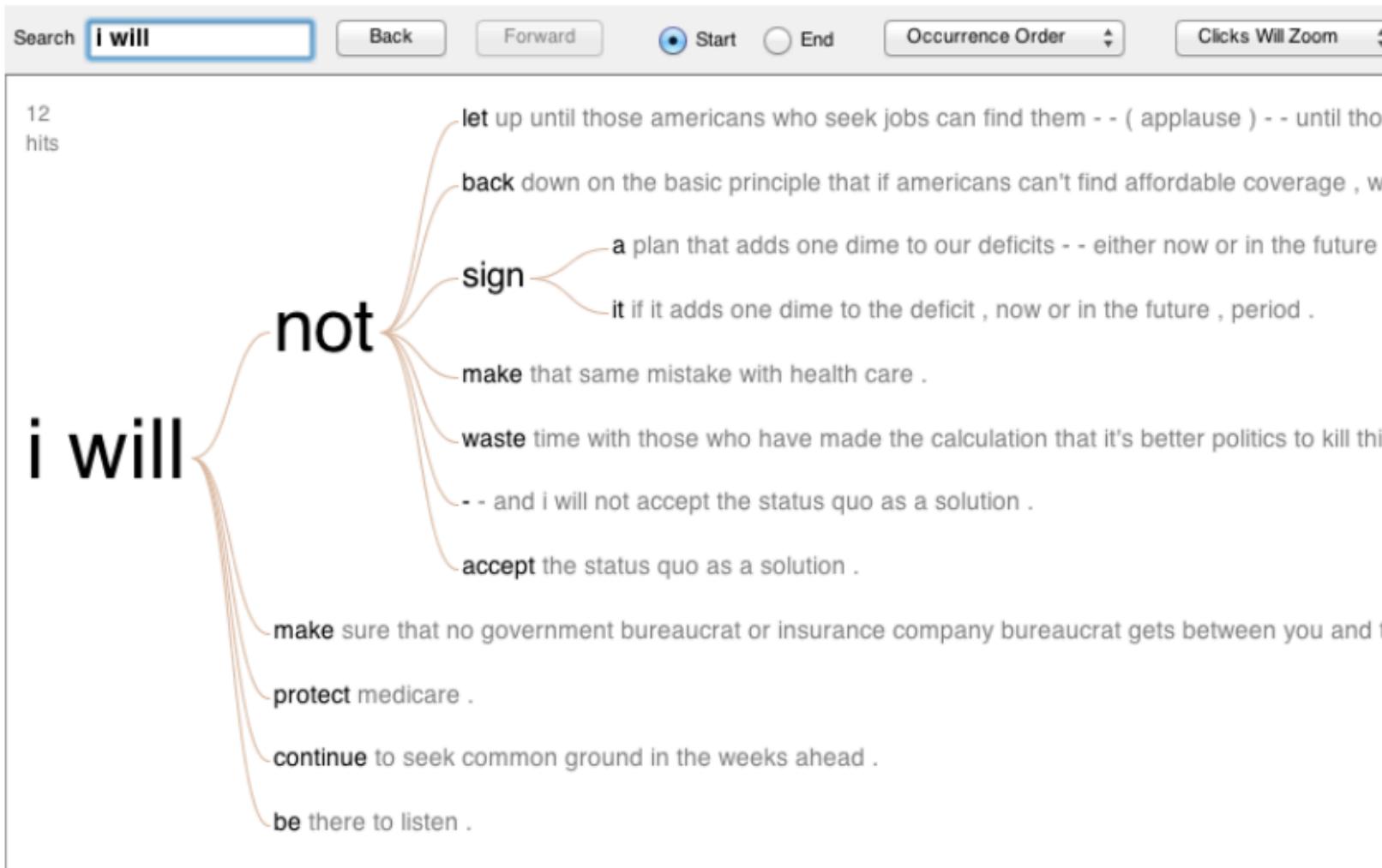


<https://economix.blogs.nytimes.com/2009/09/09/obama-in-09-vs-clinton-in-93>

Dr. Ke Zhou (<http://www.cs.nott.ac.uk/~pszkz/>)

Word Tree: Word Sequences

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



Language Model

- 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?

Challenges

- **High Dimensionality**
 - Where possible use text to represent text...
 - ... which terms are the most descriptive?
- **Context and 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.

Words as nominal data?

- High dimensional (10,000+)
- 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

- Tokenization
 - Segment text into terms.
 - Remove stop words? [a](#), [an](#), [the](#), [of](#), [to](#), [be](#)
 - Numbers and symbols? [#gocard](#), [@nottinghamforestfbball](#)
 - Entities? [Nottingham](#), [Trump](#).
- 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

Content

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

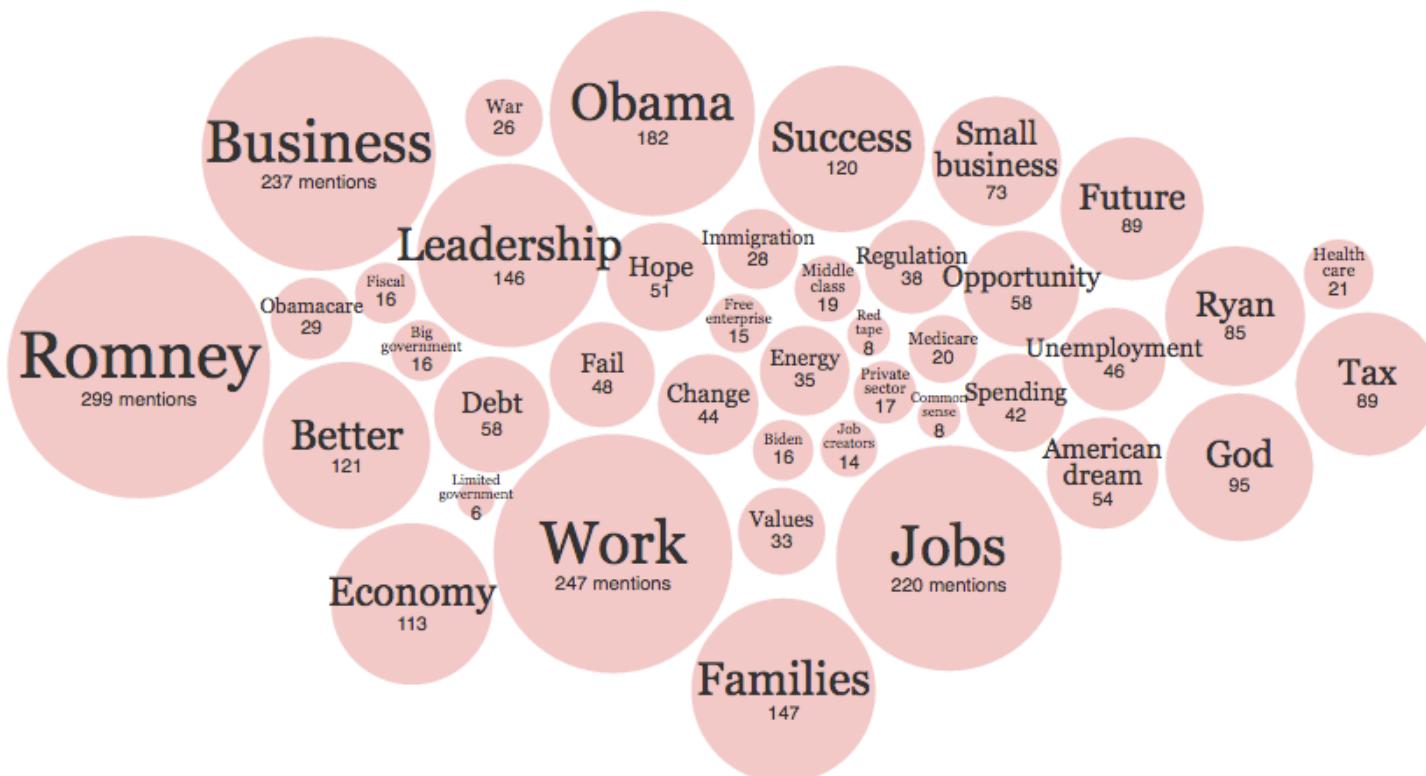
	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	157	73	0	0	0	0
Brutus	4	157	0	1	0	0
Caesar	232	227	0	2	1	1
Calpurnia	0	10	0	0	0	0
Cleopatra	57	0	0	0	0	0
mercy	2	0	3	5	5	1
worser	2	0	1	1	1	0

Word Counts

At the Republican Convention, the Words Being Used

A look at how often speakers at the Republican National Convention have used certain words and phrases so far, based on an analysis of transcripts from the Federal News Service.

Add word or phrase



Word Counts



<http://www.wordcount.org>

Tag/Word Clouds



Wordle Tag Clouds

- Tightly packed words, sometimes vertical or diagonal
- Word size is linearly correlated with frequency (typically square root in cloud)
- Multiple color palettes
- User gets some control

Viegas, Wattenberg, &
Feinberg TVCG (InfoVis) '09

Tag Clouds: Pros and Cons

- 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

Descriptive Words

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

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 $(tf_{td} / \sum_t tf_{td})$
- 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

An Example

Document 1

Term	Term Count
this	1
is	1
a	2
sample	1

Document 2

Term	Term Count
this	1
is	1
another	2
example	3

$$\text{tf}("this", d_1) = \frac{1}{5} = 0.2$$

$$\text{tf}("this", d_2) = \frac{1}{7} \approx 0.14$$

$$\text{idf}("this", D) = \log\left(\frac{2}{2}\right) = 0$$

$$\text{tfidf}("this", d_1) = 0.2 \times 0 = 0$$

$$\text{tfidf}("this", d_2) = 0.14 \times 0 = 0$$

$$\text{tf}("example", d_1) = \frac{0}{5} = 0$$

$$\text{tf}("example", d_2) = \frac{3}{7} \approx 0.429$$

$$\text{idf}("example", D) = \log\left(\frac{2}{1}\right) = 0.301$$

$$\text{tfidf}("example", d_1) = 0 \times 0.301 = 0$$

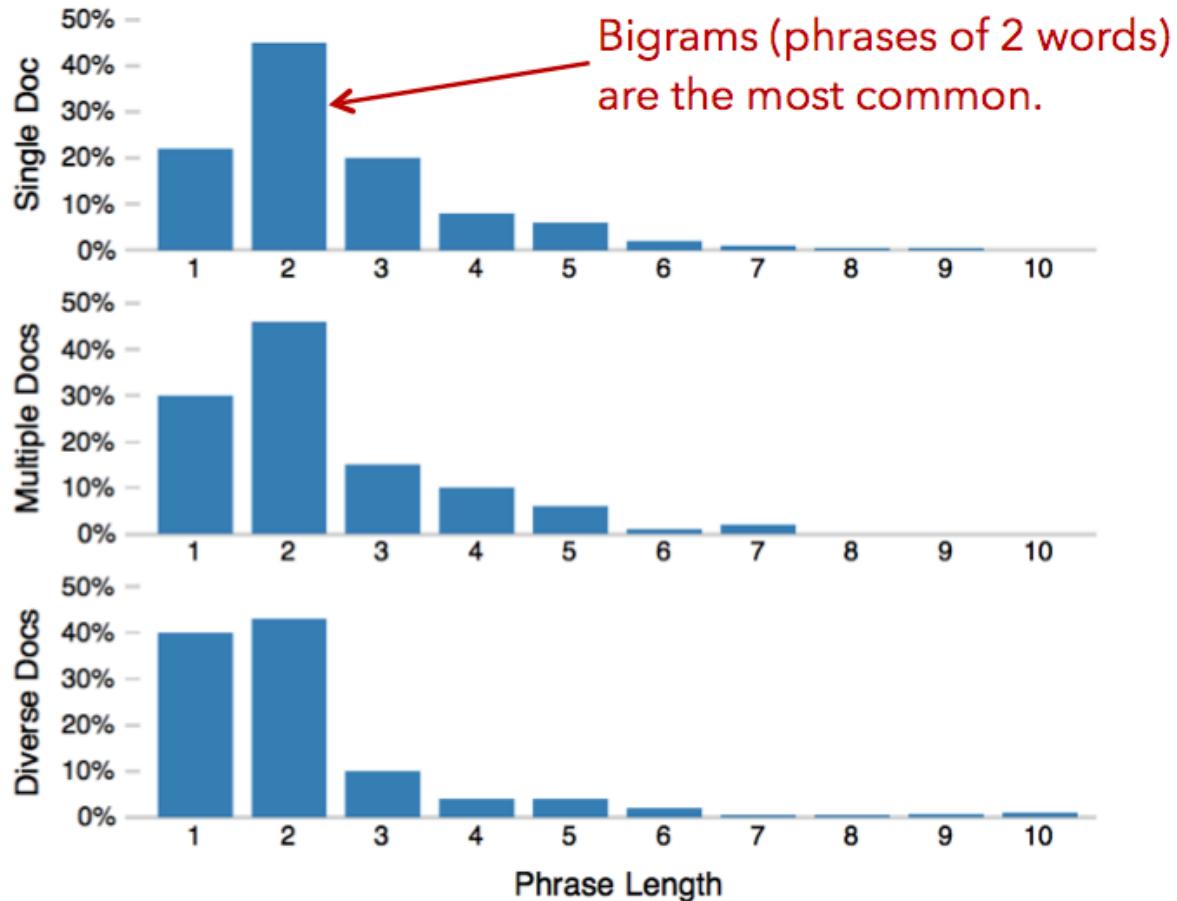
$$\text{tfidf}("example", d_2) = 0.429 \times 0.301 \approx 0.13$$

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?

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



Word vs. Phrases

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 moreabout 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.



Word vs. Phrases

Word (e.g. TFIDF)

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

Key Phrase Extraction

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



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

(Optional Reading) Automatic Keyphrase Extraction: A Survey of the State of the Art: <http://www.aclweb.org/anthology/P/P14/P14-1119.xhtml>

<http://bdewilde.github.io/blog/2014/09/23/intro-to-automatic-keyphrase-extraction/>

Dr. Ke Zhou (<http://www.cs.nott.ac.uk/~pszkz/>)

Context

Challenges

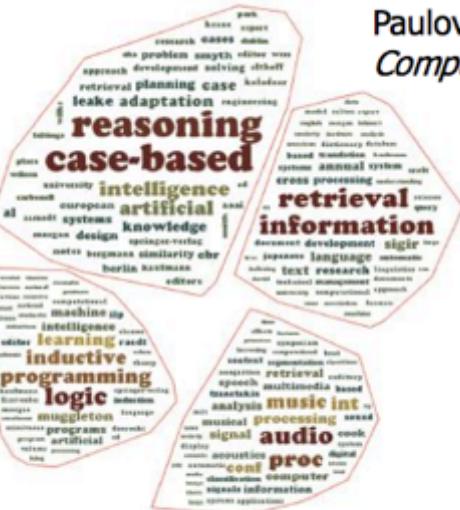
- High Dimensionality
 - Where possible use text to represent text...
 - ... which terms are the most descriptive?
- Context and Semantics
 - Provide relevant context to aid understanding.
 - Show (or provide access to) the source text.
- Modeling Abstraction
 - Determine your analysis task.
 - Understand abstraction of your language models.
 - Match analysis task with appropriate tools and models.

Semantic/Context Word Clouds

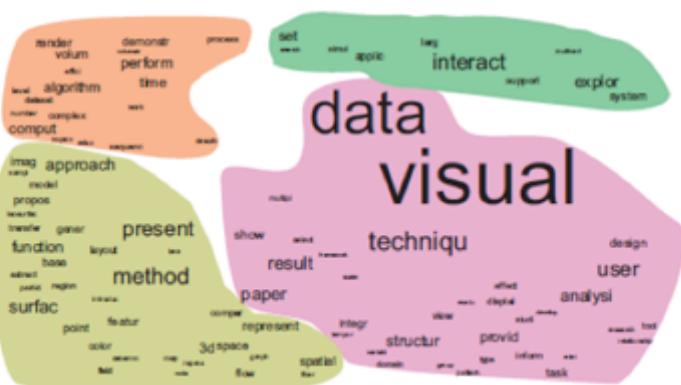
Wang et al
Graphics Interface '14



Paulovich et al
Computer Graphics Forum '12

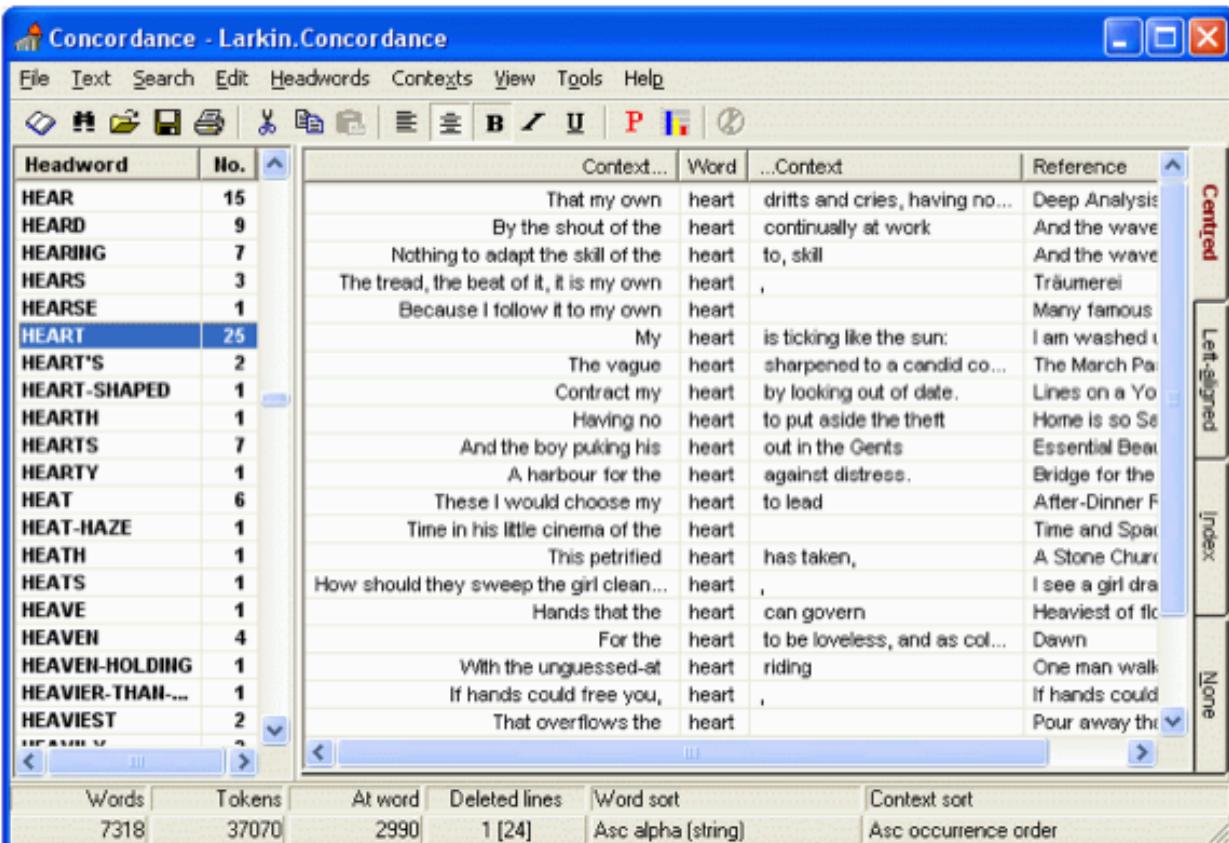


Wu et al
Computer Graphics Forum '11



Word / Phrase Context

- Concordance
 - What is the common local context of a term?



The screenshot shows the Larkin.Concordance software interface. On the left is a list of headwords with their frequencies (e.g., HEAR 15, HEART 25). The main window displays a grid of contexts for the word 'HEART'. Each row shows a context snippet, the word 'heart', and a reference line. The right side of the interface has dropdown menus for alignment ('Centred', 'Left-aligned', 'Index', 'None') and sorting ('Context sort', 'Word sort', '...Context', 'Reference').

Headword	No.	Context...	Word	...Context	Reference
HEAR	15	That my own	heart	dritts 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 flo...
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 th...

Word Tree

love the



Word Tree

- Shows context of a word or words
 - Follow word with all the phrases that follow it
- Font size shows frequency of appearance
- Continue branch until hitting unique phrase
- Clicking on phrase makes it the focus
- Ordered alphabetically, by frequency, or by first appearance

(Optional) Wattenberg & Viégas TVCG (InfoVis) '08

Word Tree Interaction



Phrase Nets

- Concordances show local, repeated structure, but what about other types of patterns?
 - Lexical: <A> at
 - Syntactic: <Noun> <Verb> <Object>
- 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

Phrase Nets

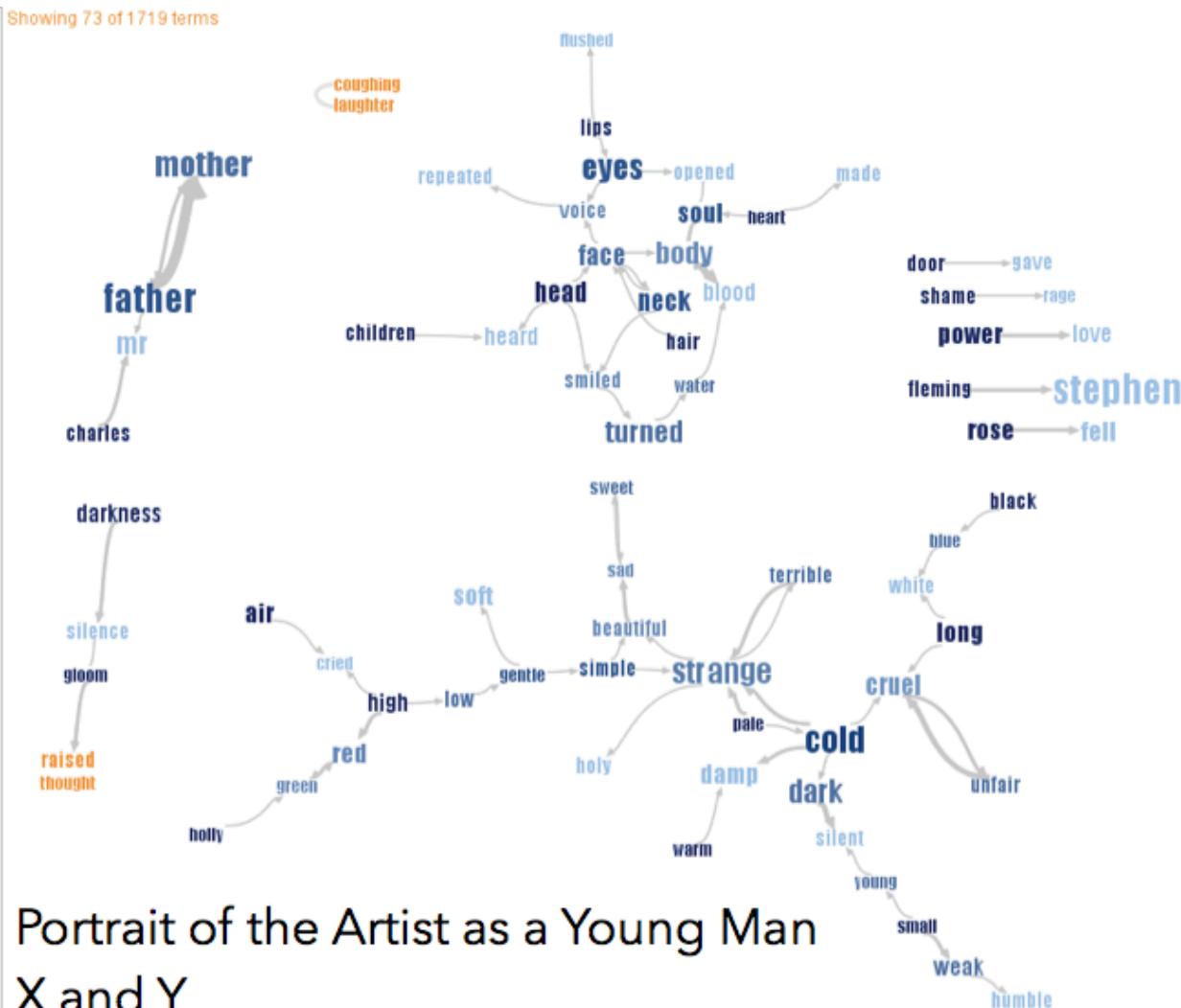
Select a phrase

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

or enter your own
 * and *

Filters
 Show top: 100
 Hide common words

Zoom

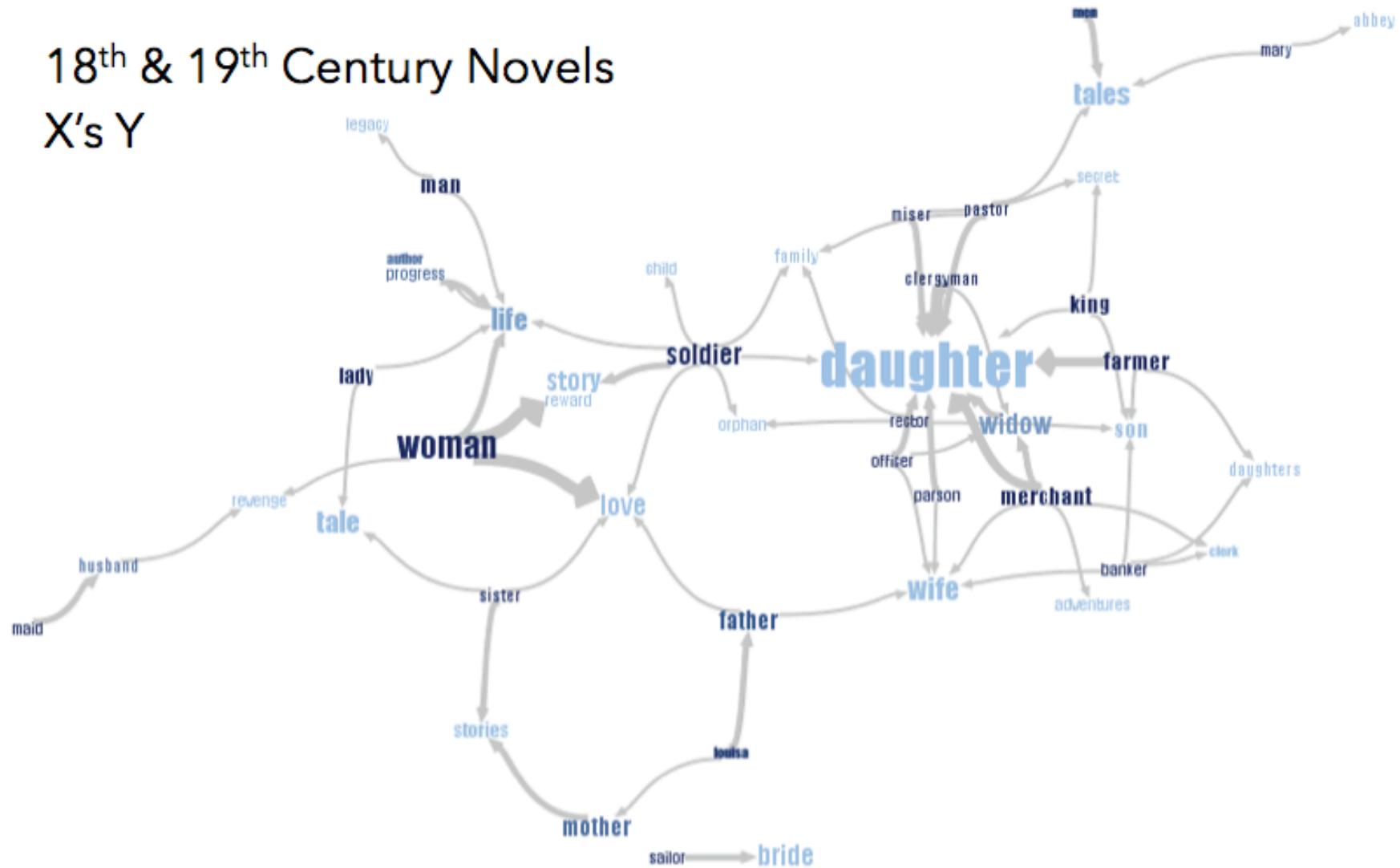


Portrait of the Artist as a Young Man
 X and Y

Phrase Nets

18th & 19th Century Novels

X's Y



SentenTree

- Elements of word clouds and word trees
 - Highlight keywords using size
 - Show sentence fragments
 - Provide a summary of the dataset
 - Enable drill-down into details



Hu, et al. TVCG '17 (InfoVis '16)

Summary of 189,450 tweets (108,702 unique) posted in a 15 minute time window around the first goal of the opening game of the 2014 Soccer World Cup

Relationship to Other Texts

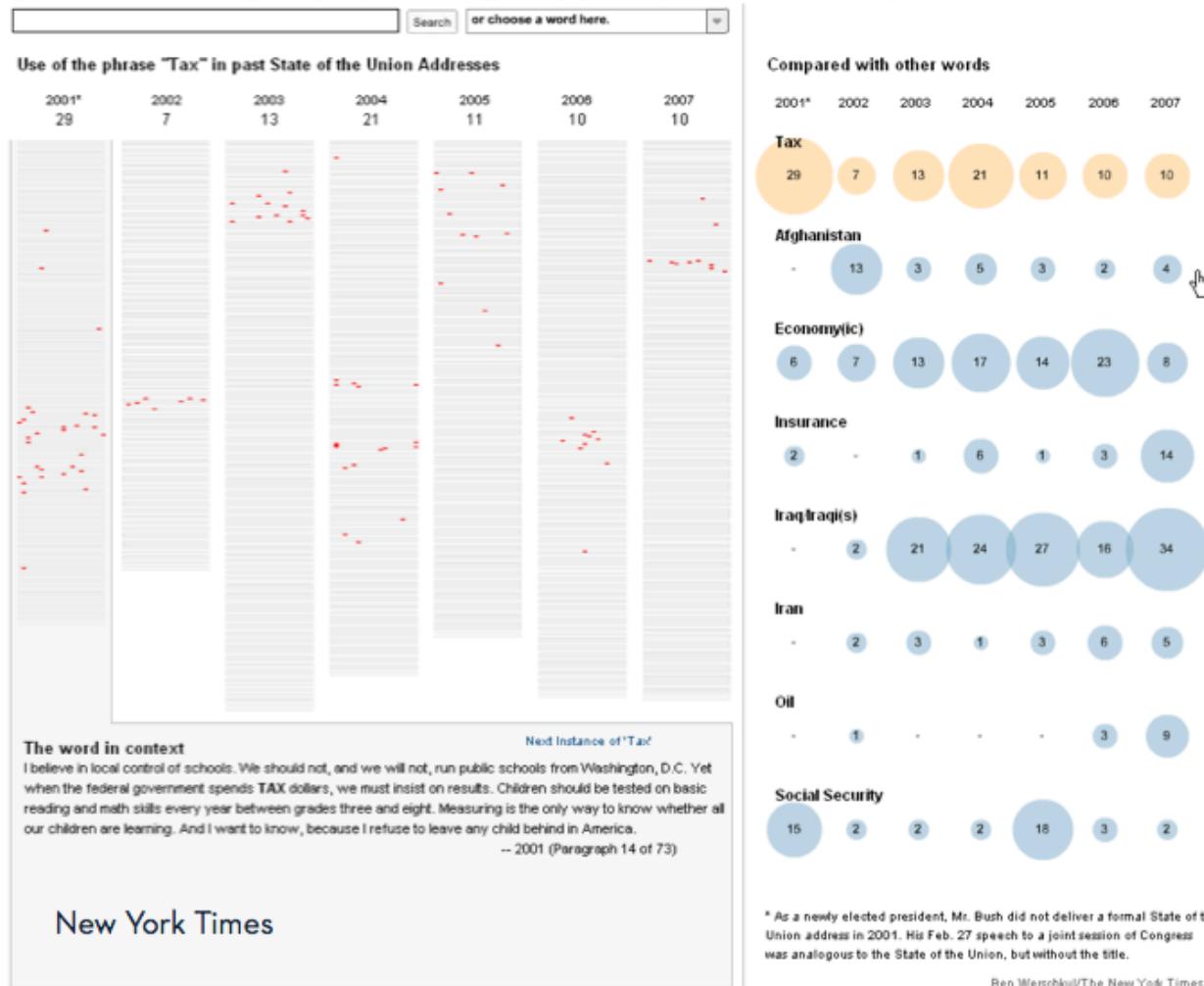
Compare to Others

THE WORDS THAT WERE USED

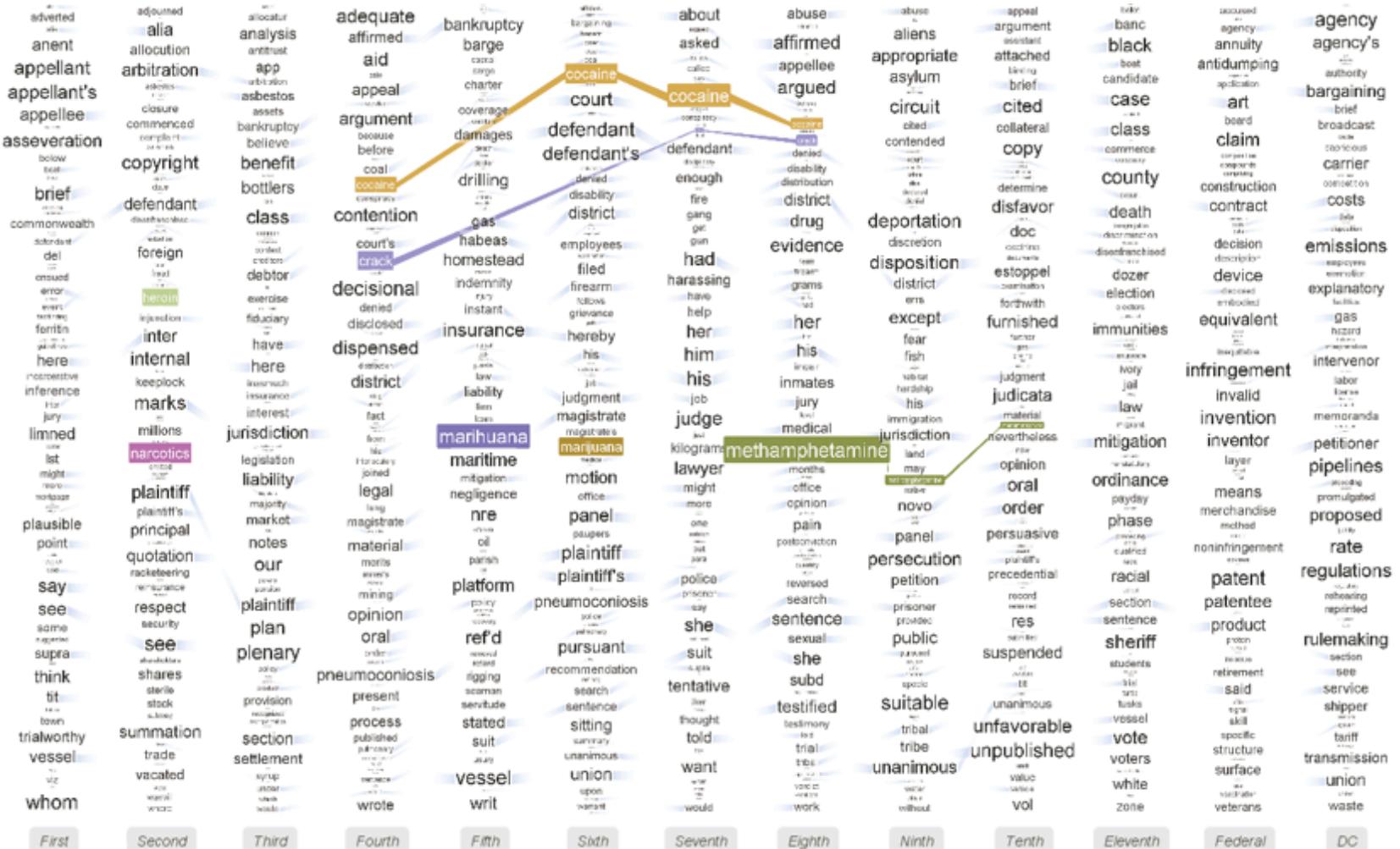
The 2007 State of the Union Address

READ 2007 SPEECH | FEEDBACK

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.



Parallel Tag Clouds



Collins et al VAST '09

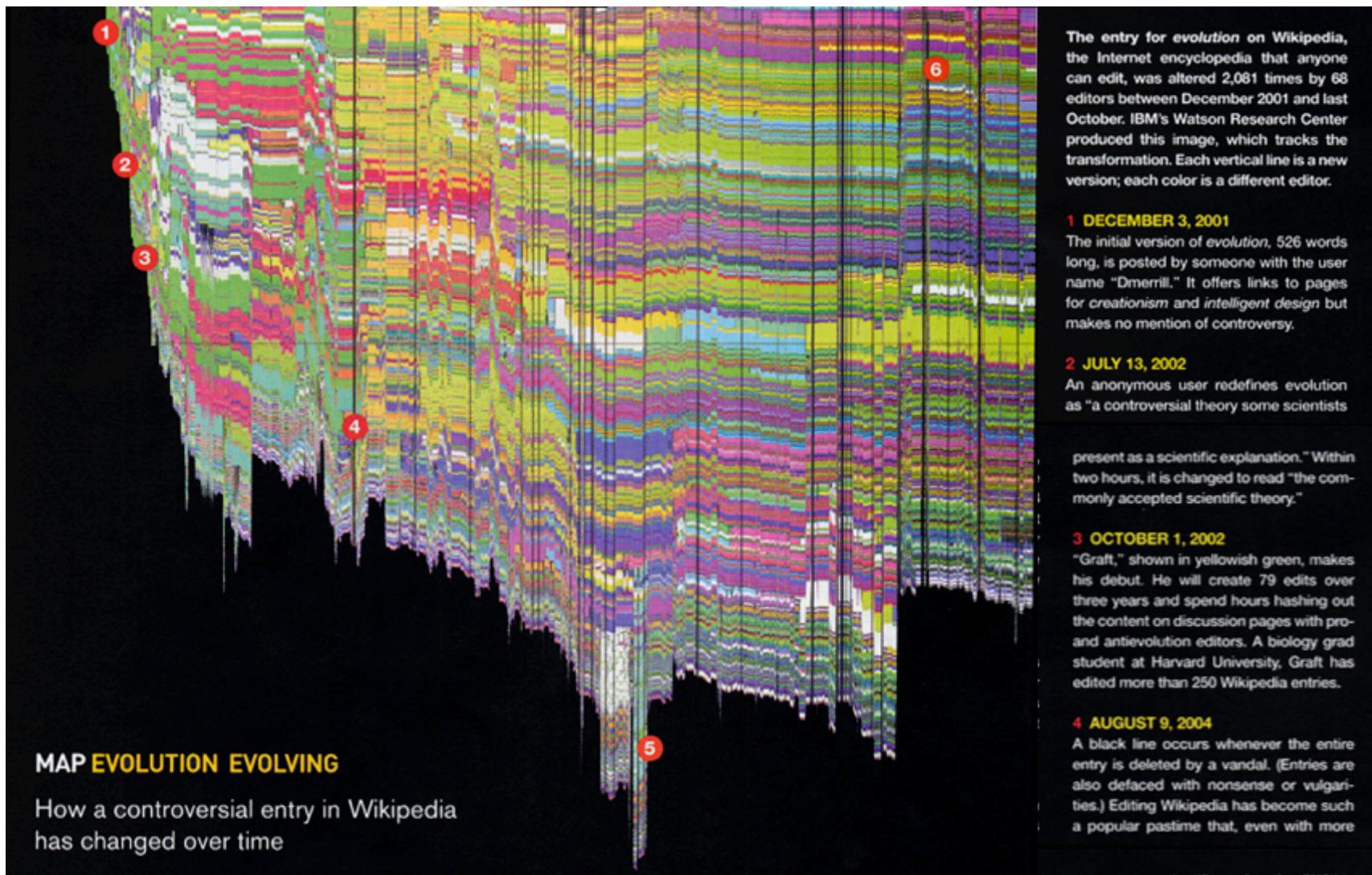
Video: <https://www.youtube.com/watch?v=rL3Ga6xBgLw>

Dr. Ke Zhou (<http://www.cs.nott.ac.uk/~pszkz/>)

Animated Traces: evolving documents



Wikipedia Edit Evolution

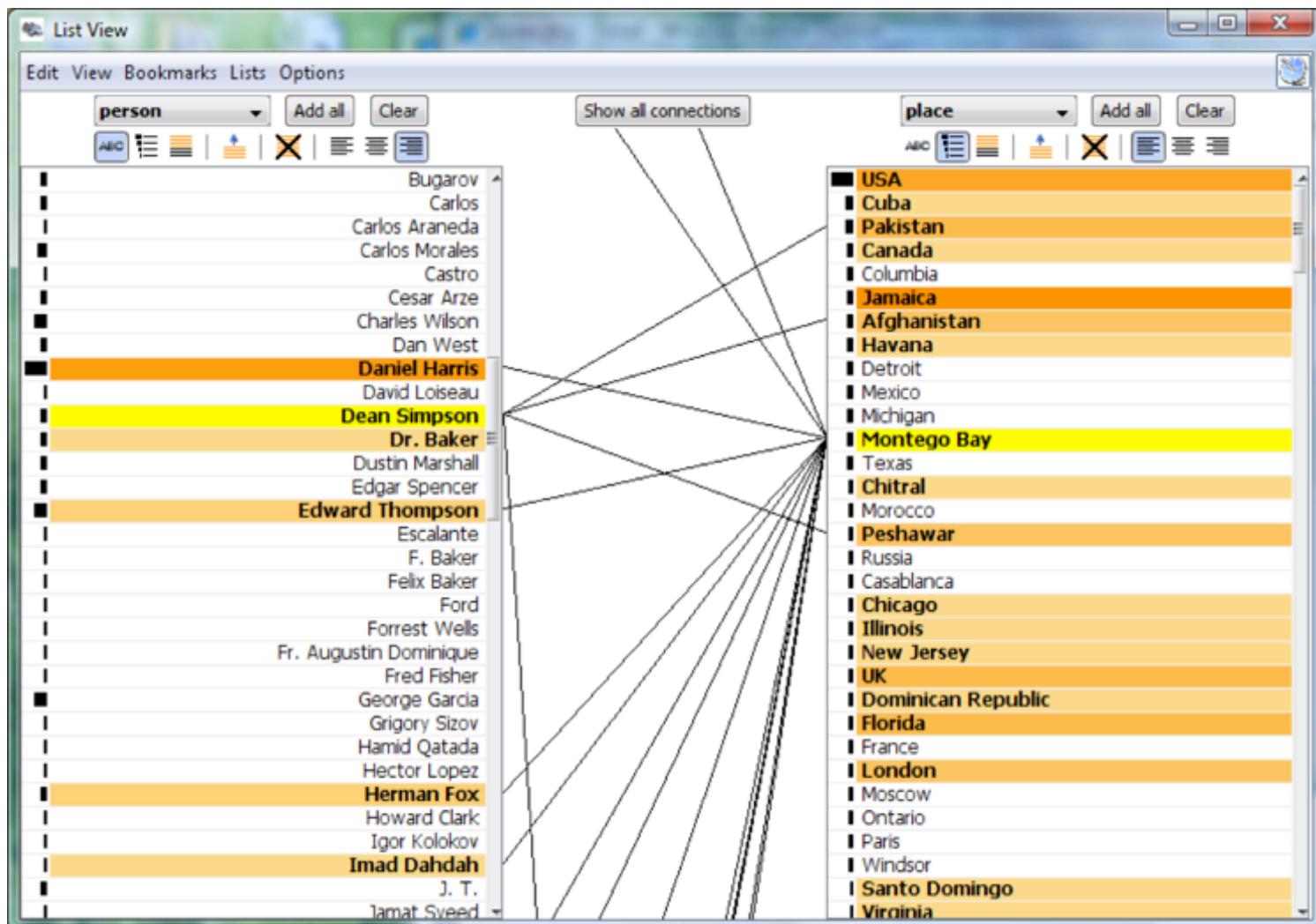


Visualizing Document Collections

Named Entity Recognition

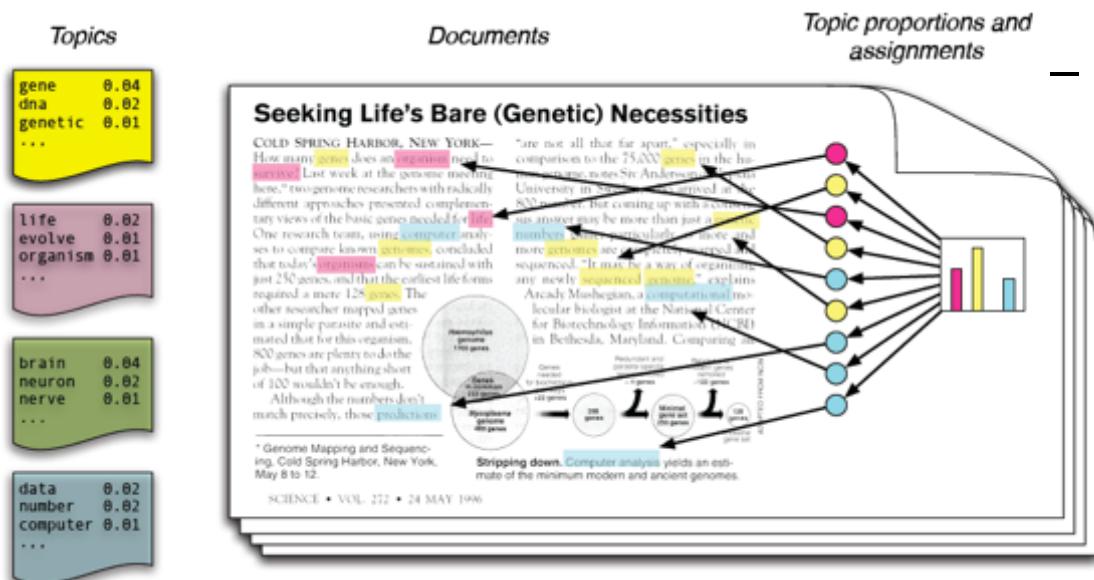
- Label named entities in text:
 - John Smith -> PERSON
 - Soviet Union -> COUNTRY
 - 353 Serra St -> ADDRESS
 - (555) 721-4312 -> PHONE NUMBER
- Entity relations: how do the entities relate?
- Simple approach: do they co-occur in a small window of text?

Entity Linkage



Similarity & Clustering

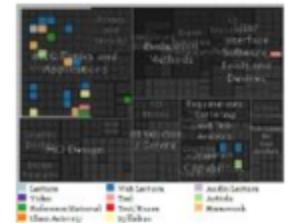
- Compute vector distance among docs
 - For TF.IDF, typically cosine distance Similarity measure can be used to cluster
- Topic modeling
 - Assume documents are a mixture of topics
 - Topics are (roughly) a set of co-occurring terms
 - Latent Semantic Analysis (LSA): reduce term matrix
 - Latent Dirichlet Allocation (LDA): statistical model



Visualization for Information Retrieval

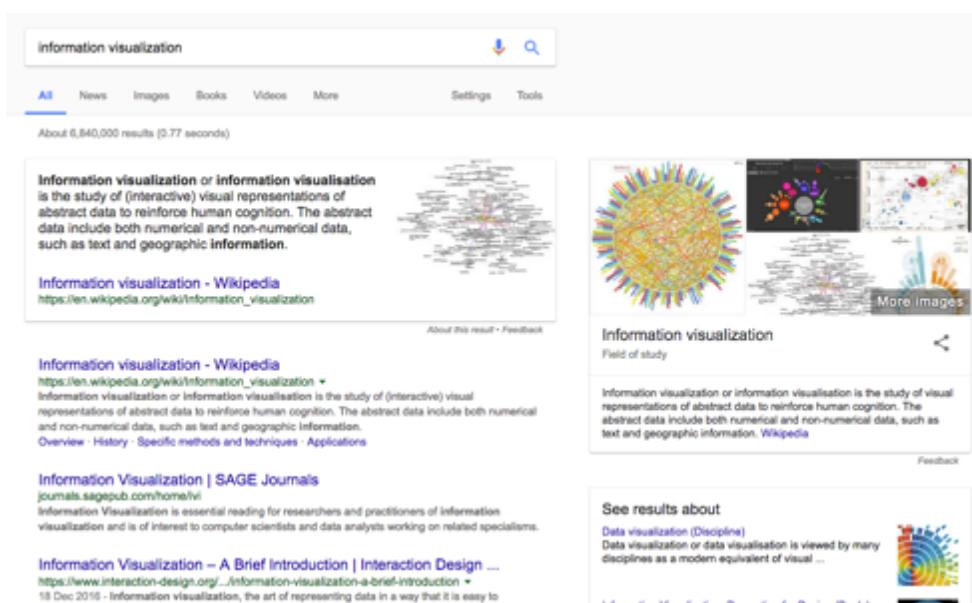


Visualization for IR
Helping search



Information Retrieval (IR)

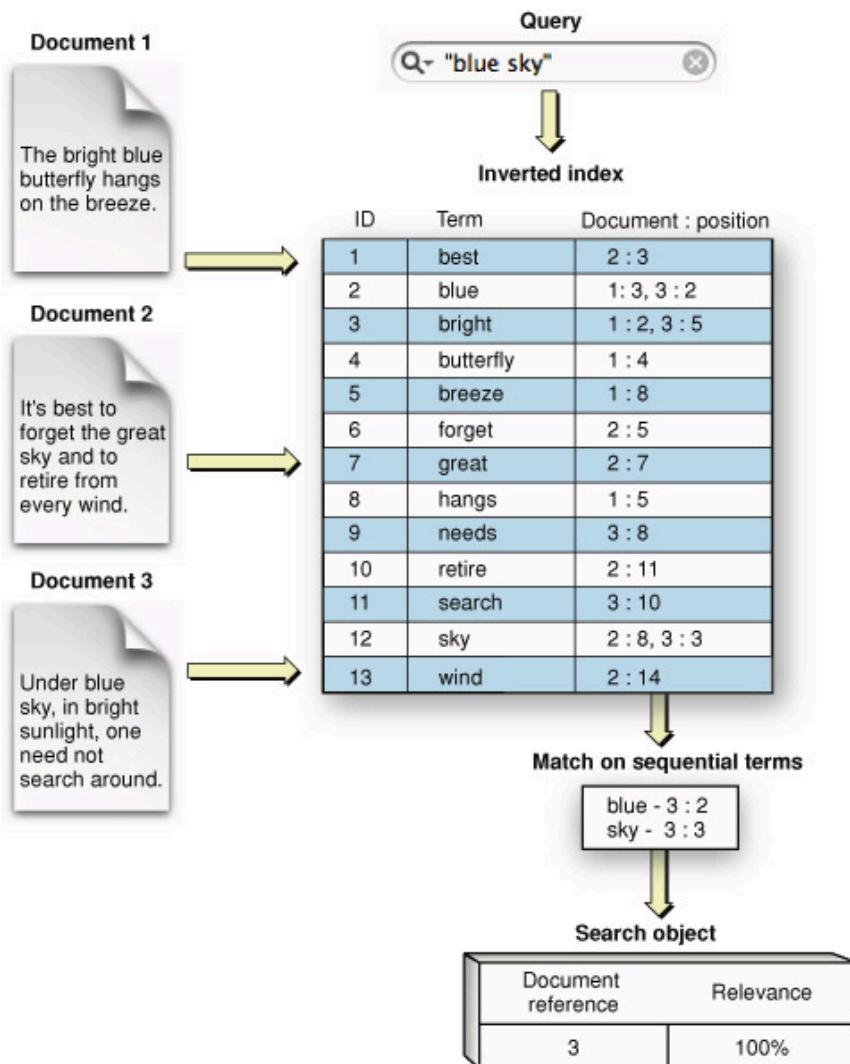
- Search for documents
- Visualization to contextualize query-doc matching results
- Can Information visualization help IR?
 - A (large) set of documents



A screenshot of a Google search results page for the query "information visualization". The search bar at the top contains the query. Below it, a navigation bar offers options: All, News, Images, Books, Videos, More, Settings, and Tools. The main content area shows approximately 6,840,000 results found in 0.77 seconds. The first result is a summary of "Information visualization or information visualisation" from Wikipedia, featuring a complex network graph visualization. Below this are several other links, including another Wikipedia entry, a journal article from SAGE Journals, and a brief introduction to interaction design. To the right of the search results, there's a sidebar titled "Information visualization" under "Field of study" with a link to the Wikipedia page. At the bottom right, there's a "Feedback" link.

Paper Handout: <http://searchuserinterfaces.com/book/>, Chapter 10

Information Retrieval



Ten Blue Links to More Complex Data and Visualization

Google search results for "homebrew podcast":

- Basic Brewing™ : Home Brewing Beer Podcast and DVD ...**
www.basicbrewing.com/radio/ ▾
We visit Modern Times Brewing in San Diego to talk to Jacob McKean and Michael Tonsmeire about starting a brewery based on homebrew, roasting and ...
Basic Brewing Radio™ 2014 - Basic Brewing Radio™ 2006
- The Brewing Network | Beer Radio for Brewers and Beer ...**
www.thebrewingnetwork.com/ ▾
We went (like true homebrewing patriots) and partied at Club Night as sponsored by ...
Podcast (jami-show): Play in new window | Download (Duration: 1:02:45 ...
Brew Strong - Shows - Beer Forum - Dr. Homebrew
- BeerSmith Home Brewing Podcast**
www.beersmith.com ▾ BeerSmith Home Brewing Forum ▾ Our Web Sites ▾
BeerSmith Home Brewing Podcast. ... by BeerSmith - All Podcast Episodes available on iTunes now! ... Episode #96 - Mastering Homebrew with Randy Mosher.
- Homebrew Podcast | Homebrew Podcasts**
www.hogtownbrewers.org/podcasts.cfm ▾
Hosted by homebrew jockeys Ron & John, Homebrew Talk is a monthly local Gainesville podcast about all things homebrewing. Get your brew on by exploring ...
- Brewing Podcasts - Brew Your Own**
<https://byo.com/hops/item/302-brewing-podcasts> ▾ Brew Your Own ▾
Homebrewing podcasts run the gamut from offering tips on the basics of homebrewing, to walking you through award-winning recipes, to attempting innovative ...
- Become a Better Brewer with the 5 Best Homebrewing ...**
blog.kegoutlet.com/become-a-better-brewer-with-the-5-best-homebrew... ▾
Here are the 5 best homebrewing podcasts that I have found and listen to regularly. These 5 podcasts will make you a better brewer. Subscribe to any or all of ...



Google search results for "nottingham":

- The University of Nottingham - a world top 1% university**
<https://www.nottingham.ac.uk/> ▾
The University of Nottingham ... Celebrating MRI in Nottingham. Marking 25 years of the Sir Peter Mansfield Magnetic Resonance Imaging Centre ...
- Nottingham - Wikipedia**
<https://en.wikipedia.org/wiki/Nottingham> ▾
Nottingham is a city and unitary authority area in Nottinghamshire, England, located 30 miles (48 km) south of Sheffield and 30 miles (48 km) north of Leicester. Nottinghamshire: University of Nottingham - List of people from Nottingham - Arnold
- The Top 10 Things to Do in Nottingham 2017 - TripAdvisor**
<https://www.tripadvisor.co.uk/.../England-Nottingshire-Nottingham> ▾
Top Things to Do in Nottingham, Nottinghamshire - Nottingham Attractions, ... Nottingham weather essentials ... Historic Sites, Science Museums, Points of Interest & Landmarks.
- Nottingham Post: Nottingham News, Sports & Events**
www.nottinghampost.com ▾
Get the latest news from the Nottingham Post online. Plus breaking news updates, sports, events and local businesses in Nottinghamshire.

Top stories



Six things Bros told us before cancelling Nottingham concert

Nottingham Post · 17 hours ago



This is what it's like to work in Nottingham's oldest sex shop

Nottingham Post · 19 hours ago



New cave is discovered under Nottingham city centre

Nottingham Post · 52 mins ago

Points of interest



Galleries of Justice Museum



Wollaton Hall



City of Caves



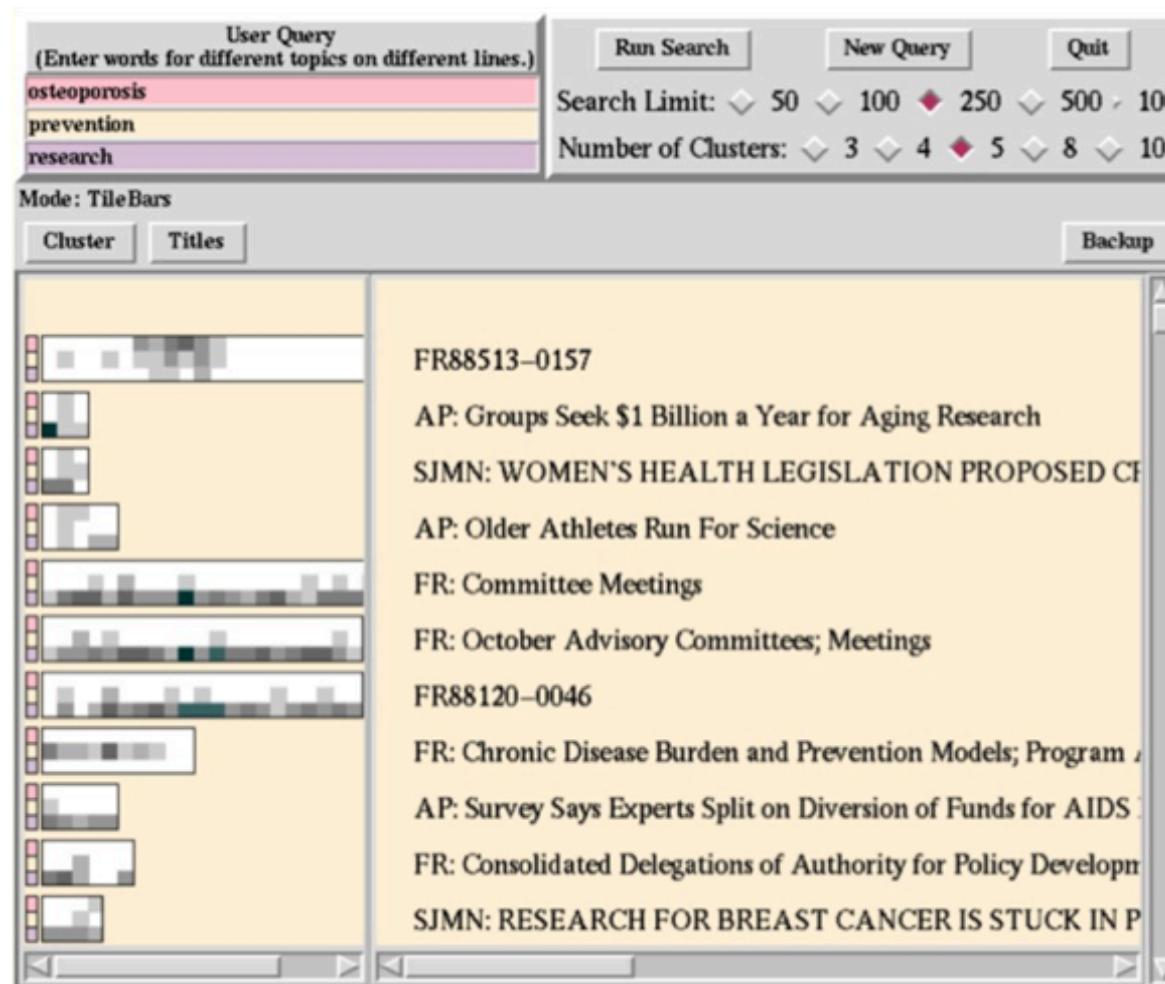
Nottingham Contemporary



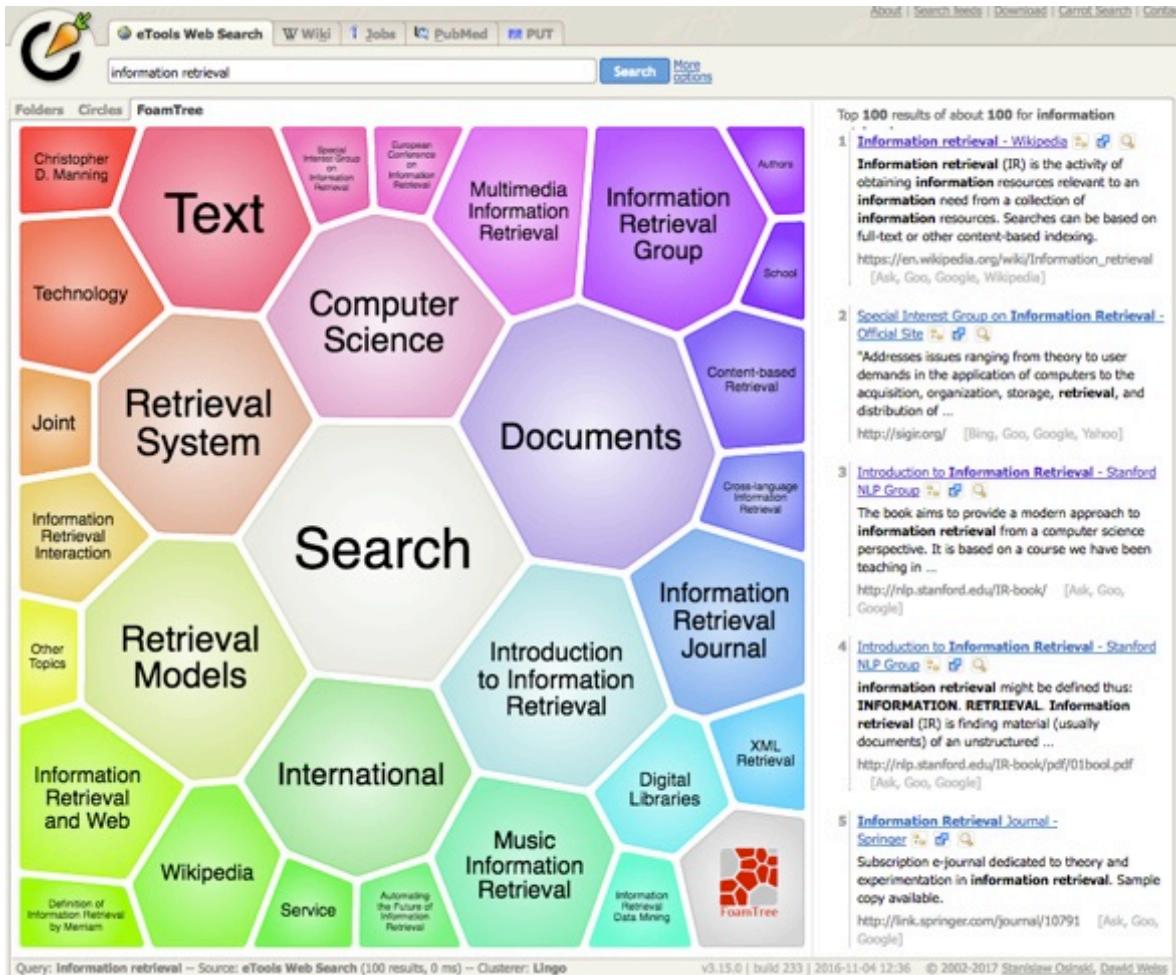
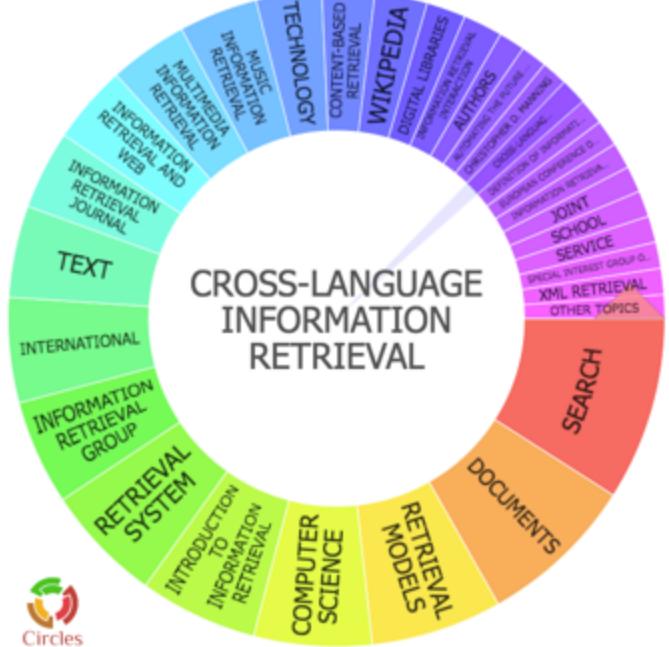
Lace Market

[View 10+ more](#)

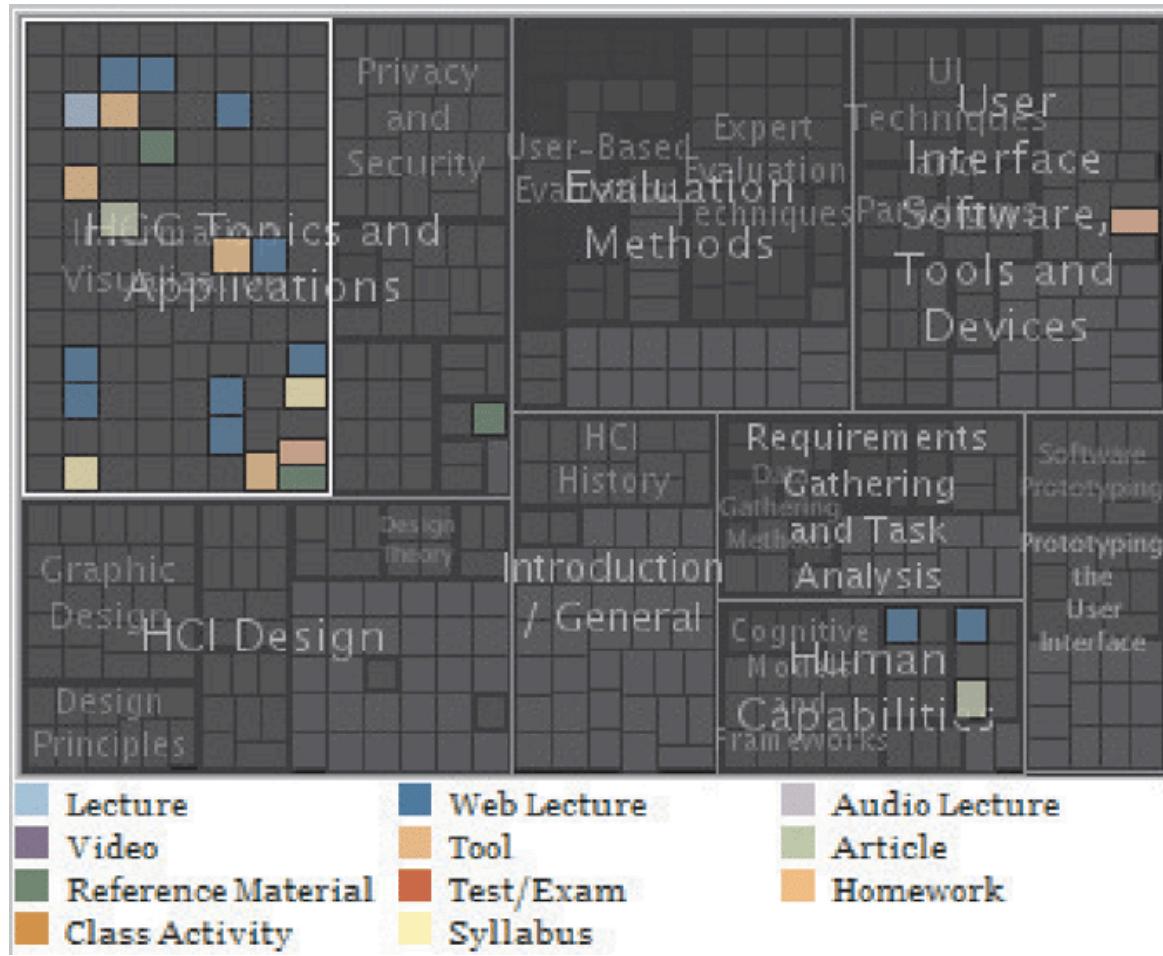
TileBars



Cluster Based Search Visualization

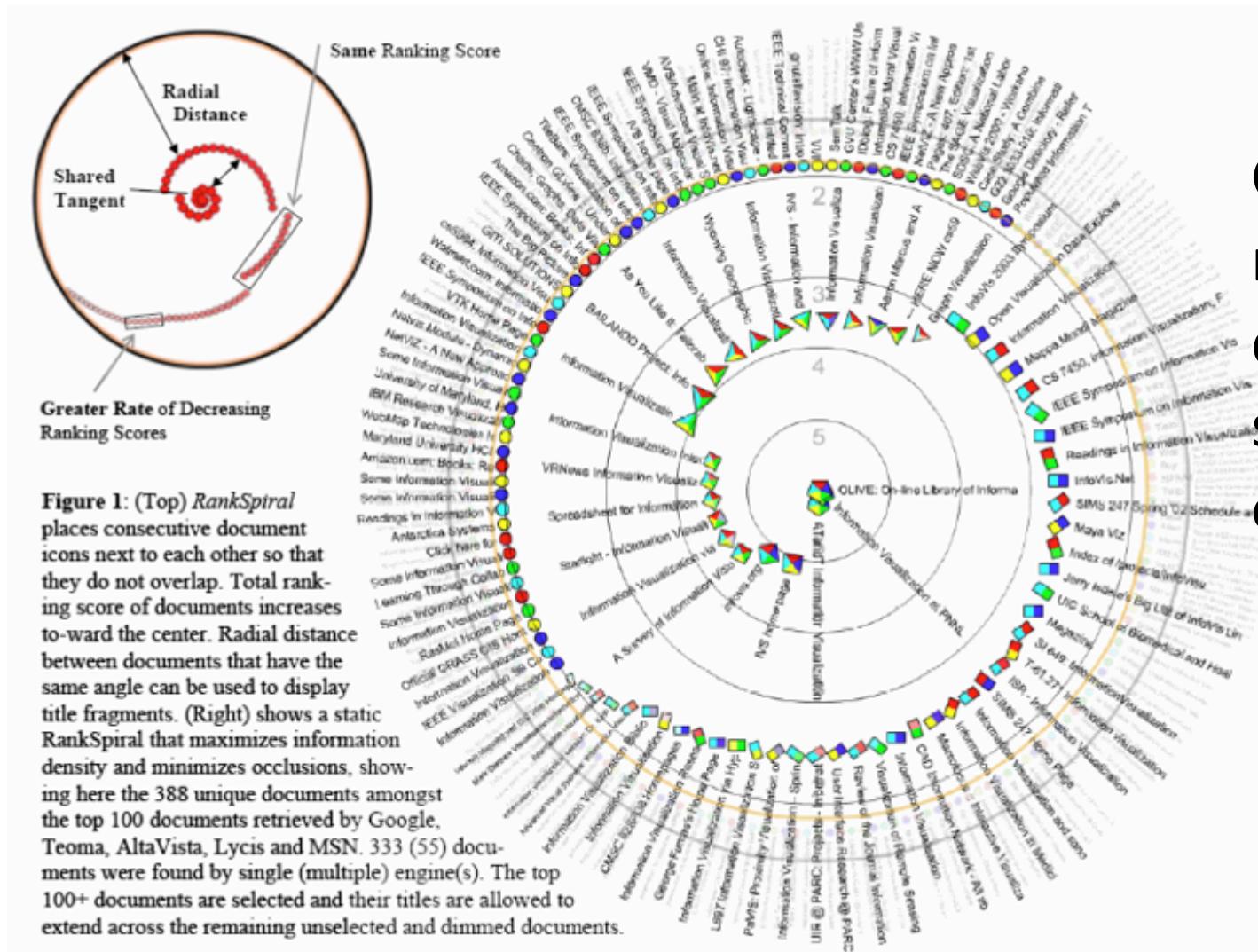


Result Maps



Treemap-style visualization for showing query results in a digital library

RankSpiral



Color
represents
different
search
engines

Summary

- High Dimensionality
 - Where possible use text to represent text...
 - ... which terms are the most descriptive?
- Context & Semantics
 - Provide relevant context to aid understanding.
 - Show (or provide access to) the source text.
- Modeling Abstraction
 - Understand abstraction of your language models.
 - Match analysis task with appropriate tools and models.
 - From bag-of-words to vector space embeddings

Next Lecture

- Topic:
 - Visualizing Time Series,
Trees and Graphs

