

Text and Document Visualization 2



CS 4460 – Intro. to Information Visualization
November 15, 2017
John Stasko

Learning Objectives



- Explain what word concordance is & how WordTree representation works
- List different queries/tasks often needed on document collections
- List various analytic metrics often calculated on documents
- List different aspects of documents often visualized
- Explain vector space document analysis (similarity calculation, search) & TFIDF
- Describe visual representation used by and contributions of these systems
 - SentenTree, TextArc, Themail, Jigsaw, ThemeScape/IN-SPIRE, ThemeRiver

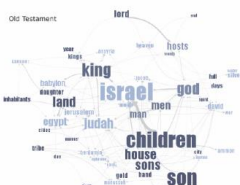
This Week's Agenda



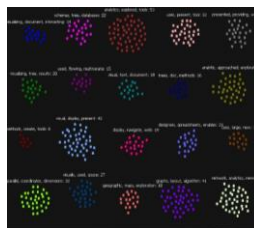
Visualizing text
Showing words,
combinations, and
context



Visualizing document sets
Words & sentences
Analysis metrics
Concepts & themes



Fall 2017



CS 4460

3

Beyond Individual Words



- The techniques (from last class) focus largely on words
 - Especially word clouds & wordles
- Can we show combinations of words, ie, actual phrases and sentences, in order to provide more context?

Fall 2017

CS 4460

4

Concordance



Definition

concordance - Definition from the Merriam-Webster Online Dictionary - Windows Internet Explorer

http://www.merriam-webster.com/dictionary/concordance

Merriam-Webster Online

concordance

One entry found.

Concordance

Find the Benefits of Concordance Software by LexisNexis. Buy Now!
law.lexisnexis.com

Main Entry: con-cord-ance

Pronunciation: 'kan-'kor-'dʌn(t)s, kən-'

Function: noun

Etymology: Middle English, from Anglo-French, from Medieval Latin *concordantia*, from Latin *concordant-, concordans*, present participle of *concordare* to agree, from *concord-, concors*

Date: 14th century

1 : an alphabetical index of the principal words in a book or the works of an author with their immediate contexts

2 : CONCORD, AGREEMENT

Fall 2017

CS 4460

5

Concordance in Text



Concordance - Larkin, Concordance

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 t
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 Gerts	Essential Bea
HEARTY	1	A harbour for the	heart	against distress.	Bridge for the
HEAT	6	These I would choose my	heart	to lead	After-Dinner F
HEAT-HAZE	1	Time in his little cinema of the	heart		Time and Spa
HEATH	1	This petrified	heart	has taken,	A Stone Churc
HEATS	1	How should they sweep the girl clean...	heart	,	I see a girl dra
HEAVE	1	Hands that the	heart	can govern	Heaviest of fit
HEAVE	1	For the	heart	to be loveless, and as col...	Dawn
HEAVEN	4	With the unguessed-at	heart	riding	One man walk
HEAVEN-HOLDING	1	If hands could free you,	heart	,	If hands could
HEAVIER-THAN...	1	That overflows the	heart		Pour away thi

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

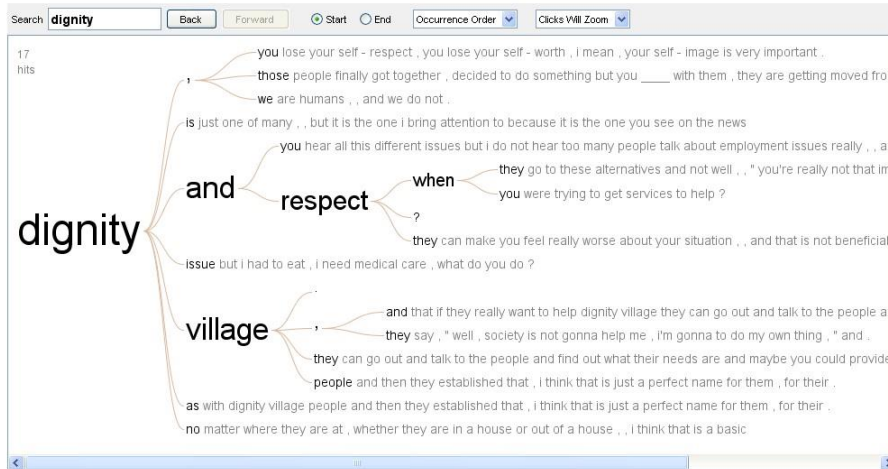
http://www.concordancesoftware.co.uk

Fall 2017

CS 4460

6

Word Tree



Fall 2017

CS 4460

7

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

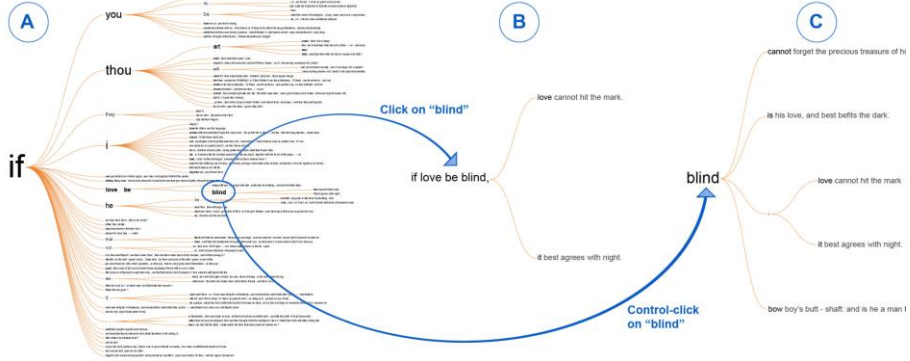
Wattenberg & Viégas
TVCG (InfoVis) '08

Fall 2017

CS 4460

8

Interaction

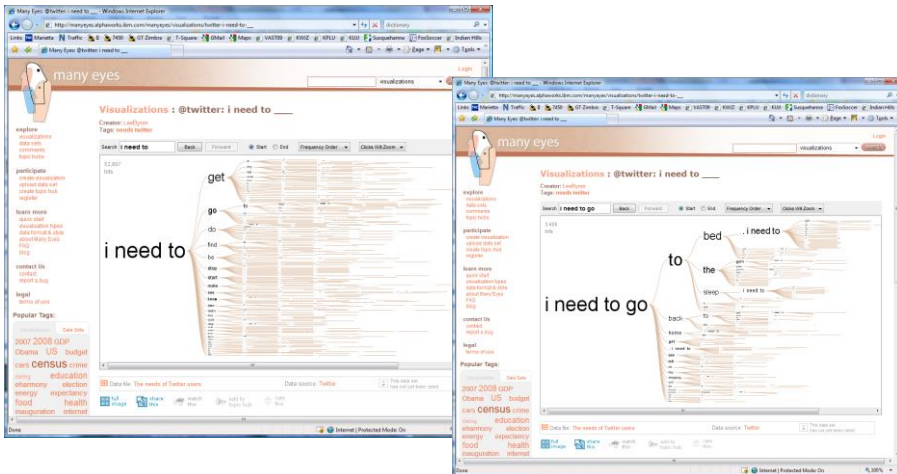


Fall 2017

CS 4460

9

Many Eyes' WordTree



Fall 2017

CS 4460

10

Words and Context



- Can we show most frequent words like a word cloud but also provide context?
 - Should each word appear one time?
 - But then how to show context?
 - If appears multiple times, how to make that work?

Fall 2017

CS 4460

11

Hu, Wongsuphasawat, and Stasko
TVCG '17 (InfoVis '16)

SentenTree



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



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

Fall 2017

CS 4460

12

Example



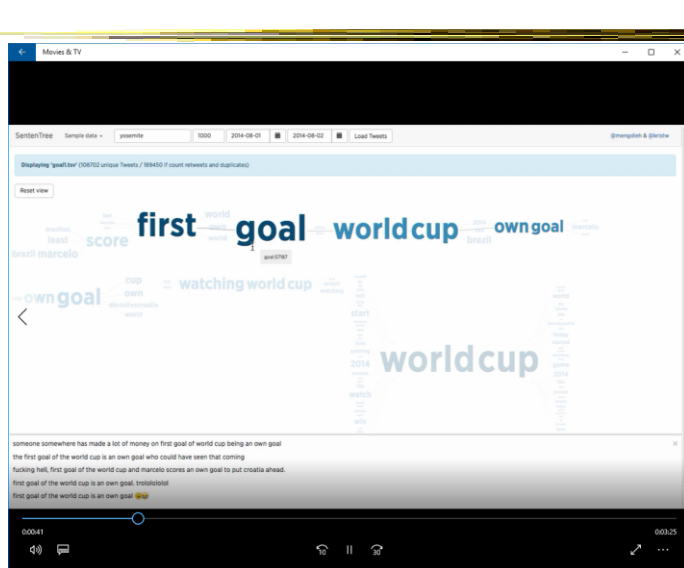
Tweets mentioning word "Yosemite" from Aug 1, 2014

Fall 2017

CS 4460

13

Video



Fall 2017

CS 4460

14

Another Challenge



- Visualize an entire book
- What does that mean?
 - Word appearances
 - Sentences
 - ...

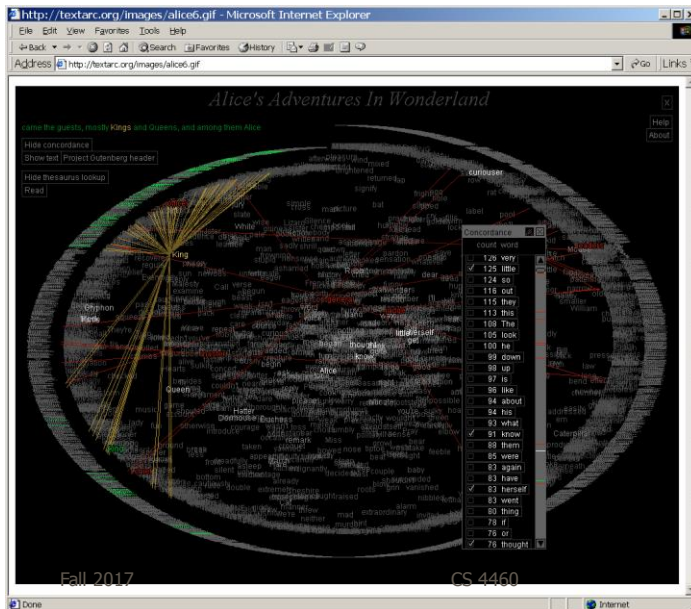
Fall 2017

CS 4460

15

TextArc

<http://textarc.org>



Sentences laid out
in order of appearance

Words near to where
they appear

Significant interaction

Brad Paley

16

Related Topic - Sensemaking



- Sensemaking
 - Gaining a better understanding of the facts at hand in order to take some next steps
 - (Better definitions in VA lecture)
- InfoVis can help make a large document collection more understandable more rapidly

Exercise



- Suppose you have 300 reviews of a TV/LCD monitor on Amazon?
 - Want to build a visualization of these...but
- What tasks/goals/questions would someone want to accomplish with the vis?
- Let's generate a list...

Questions/Tasks



What is the tone of a review? (sentiment analysis)
What words show up in the reviews, especially in pertinent subsets?
What was avg score?
Facilitate comparisons.
Find features, see what words go with each.
Distribution of scores.
How have scores changed over time?
Surface pictures of product.
Show attributes/metrics of diff reviews
Represent writer's background/history.
If available, compare/contrast public vs expert views
Check out certain important words "recommend", "suggest"
Look for frequencies of certain key words.
Make actual text available.

Fall 2017

CS 4460

19

Overlaps & Similarities



- Are some items in our list in the same "category"?
 - Can we generalize a little and narrow the list down to some core questions/tasks?

Fall 2017

CS 4460

20

Evaluate a Vis



- Use the generated list to evaluate a visualization for this problem
- Think about this for your P5...

Example Tasks & Goals



- Which documents contain text on topic XYZ?
- Which documents are of interest to me?
- Are there other documents that are similar to this one (so they are worthwhile)?
- How are different words used in a document or a document collection?
- What are the main themes and ideas in a document or a collection?
- Which documents have an angry tone?
- How are certain words or themes distributed through a document?
- Identify "hidden" messages or stories in this document collection.
- How does one set of documents differ from another set?
- Quickly gain an understanding of a document or collection in order to subsequently do XYZ.
- Understand the history of changes in a document.
- Find connections between documents.

Various Document Metrics



- Different variables for literary analysis
 - Average word length
 - Syllables per word
 - Average sentence length
 - Percentage of nouns, verbs, adjectives
 - Frequencies of specific words
 - Hapax Legomena – number of words that occur once

How would you visualize this?

Keim & Oelke
VAST '07

Fall 2017

CS 4460

23

Vis

Each block represents a contiguous set of words, eg, 10,000 words

Do partial overlap in blocks for a smoother appearance

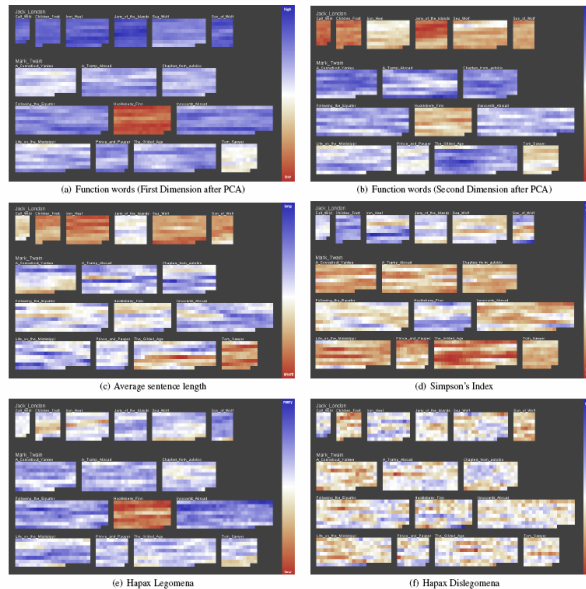


Figure 2: Fingerprints of books of Mark Twain and Jack London. Different measures for authorship attribution are tested. If a measure is able to discriminate between the two authors, the visualizations of the books that are written by the same author will equal each other more than the visualizations of books written by different authors. It can easily be seen that this is not true for every measure (e.g. Hapax Dislegomena). Furthermore, it is interesting to observe that the book *Huckleberry Finn* sticks out in a number of measures as if it is not written by Mark Twain.

Fall 2017

CS 4460

24

The Bible

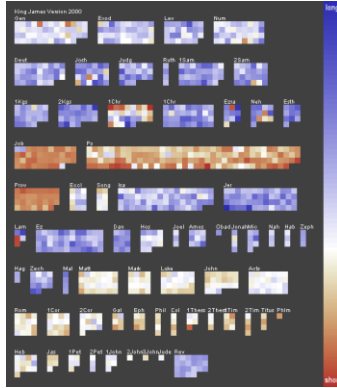


Figure 4: Visual Fingerprint of the Bible. Each pixel represents one chapter of the bible and color is mapped to the average verse length. Interesting characteristics such as the generally shorter verses of the poetry books, the inhomogeneity of the 1. Book of Chronicles or the difference between the Old Testament and the New Testament can be perceived.

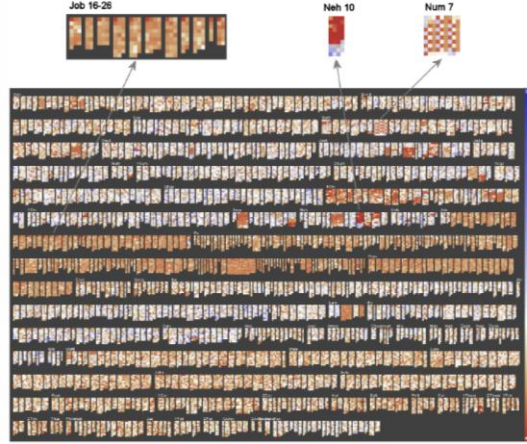


Figure 5: Visual Fingerprint of the Bible. More detailed view on the bible in which each pixel represents a single verse and verses are grouped to chapters. Color is again mapped to verse length. The detailed view reveals some interesting patterns that are camouflaged in the averaged version of fig. 4.

Fall 2017

CS 4460

25

Bohemian Bookshelf

Video



Serendipitous browsing



Fall 2017

CS 4460

Thudt et al
CHI '12

26

Themail



- Visualize one's email history
 - With whom and when has a person corresponded
 - What words were used
- Answer questions like:
 - What sorts of things do I (the owner of the archive) talk about with each of my email contacts?
 - How do my email conversations with one person differ from those with other people?

Viégas, Golder & Donath
CHI '06

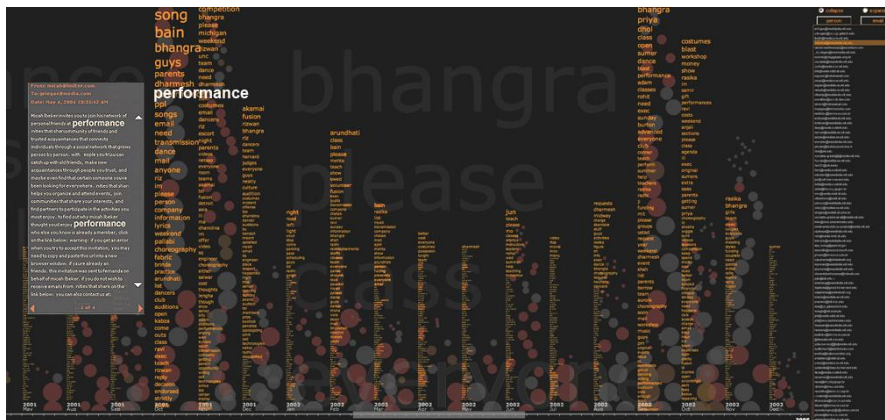
Fall 2017

CS 4460

27

Interface

Text analysis to seed visualization
Monthly & yearly words



Fall 2017

CS 4460

28

More Document Info



- Highlight entities within documents
 - People, places, organizations
- Document summaries
- Document similarity and clustering
- Document sentiment

Fall 2017

CS 4460

29

Jigsaw

Saw earlier in term



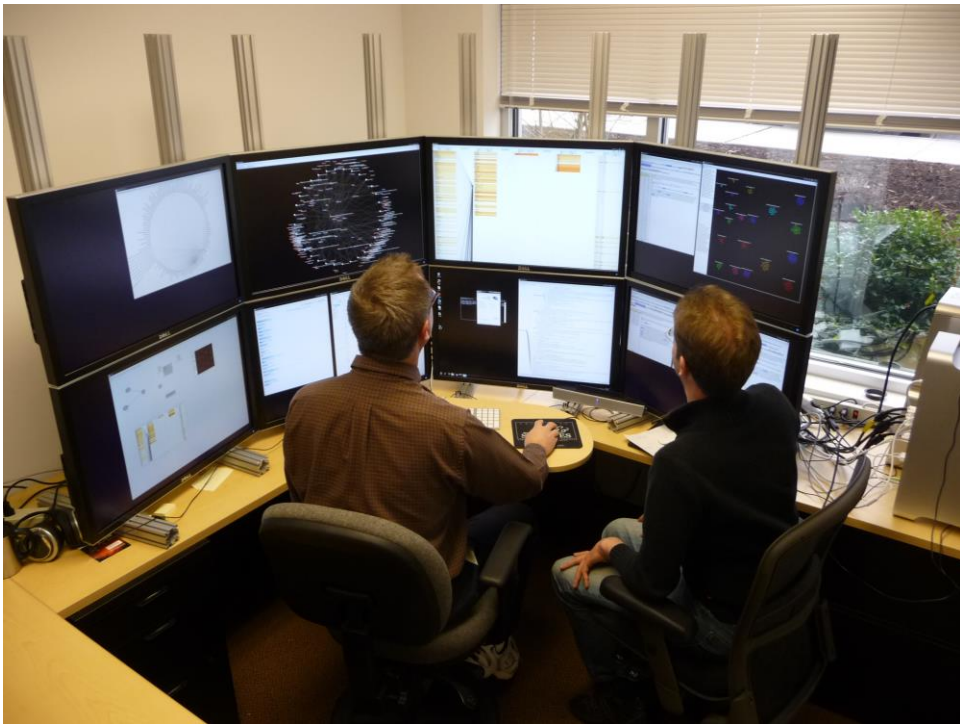
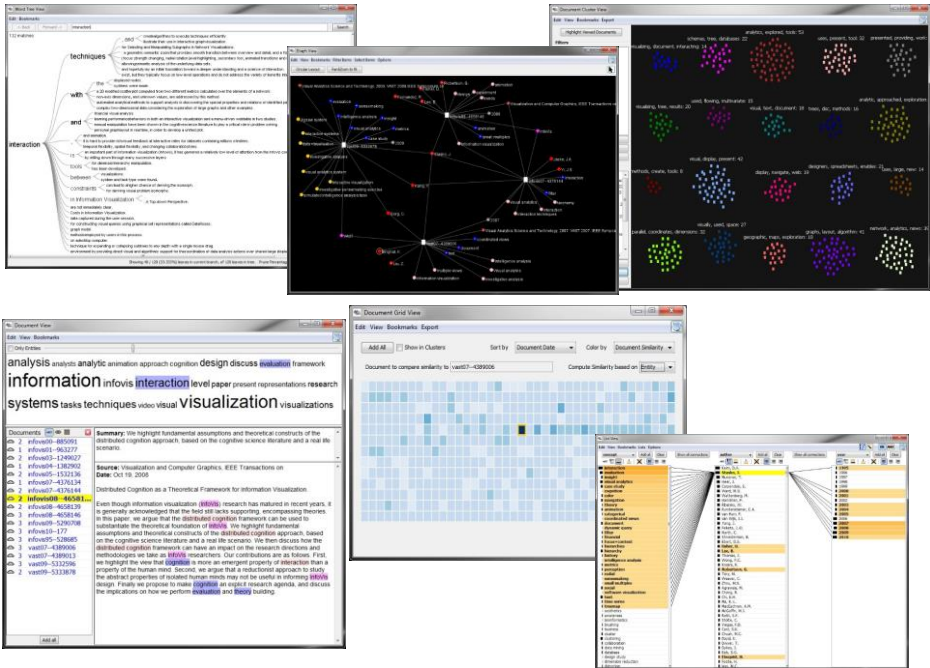
- Targeting sense-making scenarios
- Variety of visualizations ranging from word-specific, to entity connections, to document clusters
- Primary focus is on entity-document and entity-entity connection
- Search capability coupled with interactive exploration

Stasko, Görg, & Liu
Information Visualization '08
Görg et al
IEEE TVCG '13

Fall 2017

CS 4460

30



Temporal Issues



- What if documents have dates?
- Can we visualize their contents over time?

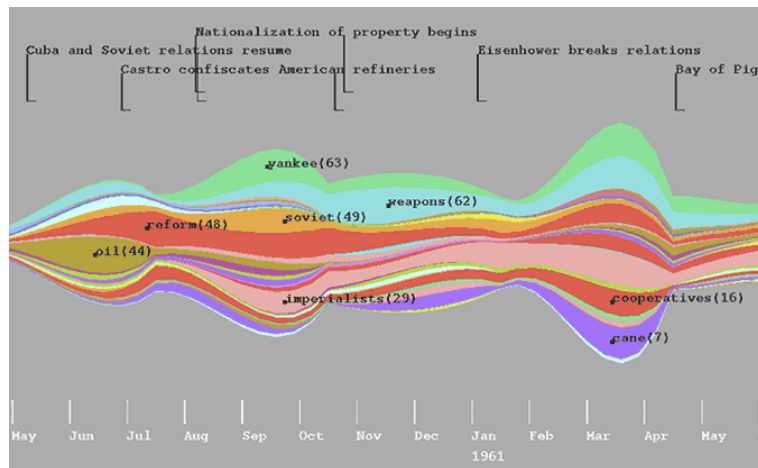
Fall 2017

CS 4460

33

ThemeRiver

Each stream is a word or a theme



Fall 2017

CS 4460

Havre, Hetzler, & Nowell
InfoVis '00

34

Up to Higher Level



- How do we present the contents, semantics, themes, etc of the documents
 - Someone may not have time to read them all
 - Someone just wants to understand them
- Who cares?
 - Researchers, fraud investigators, CIA, news reporters

Fall 2017

CS 4460

35

Vector Space Analysis



- How does one compare the similarity of two documents?
- One model
 - Make list of each unique word in document
 - Throw out common words (a, an, the, ...)
 - Make different forms the same (bake, bakes, baked)
 - Store count of how many times each word appeared
 - Alphabetize, make into a vector

Fall 2017

CS 4460

36

Vector Space Analysis



- Model (continued)
 - Want to see how closely two vectors go in same direction, inner product
 - Can get similarity of each document to every other one
- Some similarities to how search engines work

Fall 2017

CS 4460

37

Wiggle



- Not all terms or words are equally useful
- Often apply TFIDF
 - Term frequency, inverse document frequency
- Weight of a word goes up if it appears often in a document, but not often in the collection

Fall 2017

CS 4460

38

Document Collection Maps



- Use document similarity idea
- Use mass-spring graph-like algorithm for clustering similar documents together and moving dissimilar documents far apart

Fall 2017

CS 4460

39

Work at PNNL

<http://www.pnl.gov/infoviz>



- Group has developed a number of visualization techniques for document collections
 - Galaxies
 - Themescapes
 - ThemeRiver
 - ...

Wise et al
InfoVis '95

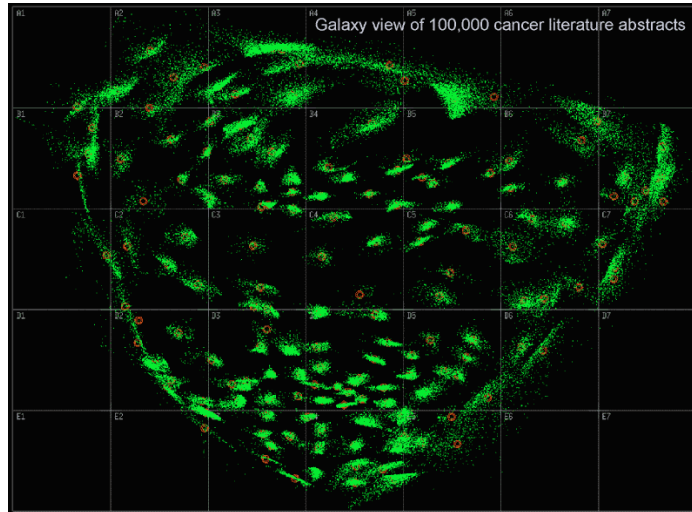
Fall 2017

CS 4460

40

Galaxies

Presentation of documents where similar ones cluster together



Fall 2017

CS 4460

41

Themespaces



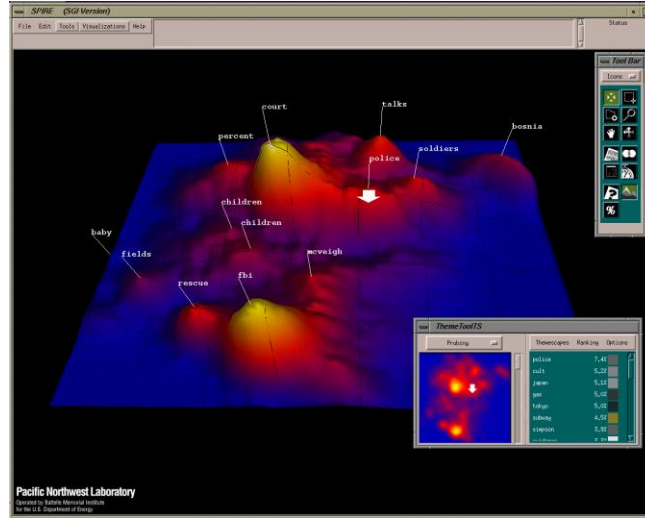
- Self-organizing maps didn't reflect density of regions all that well -- Can we improve?
- Use 3D representation, and have height represent density or number of documents in region

Fall 2017

CS 4460

42

Themescape



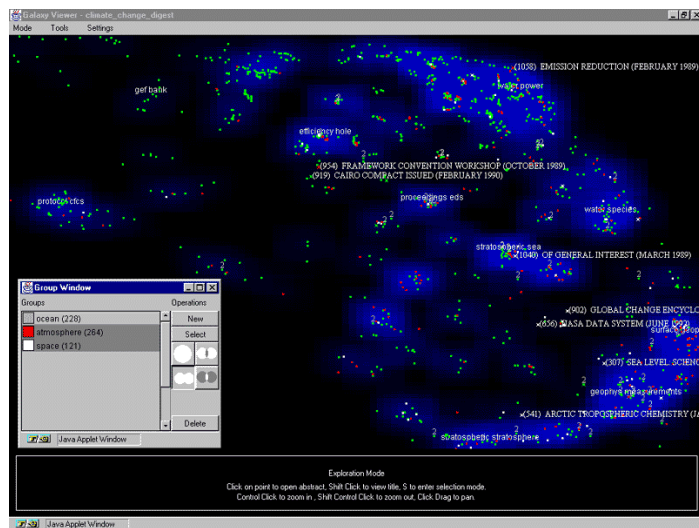
Video

Fall 2017

CS 4460

43

WebTheme



Fall 2017

CS 4460

44

Related Topic



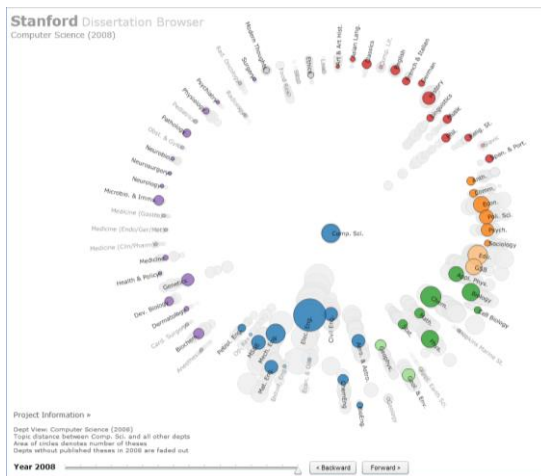
- Maps of Science
- Visualize the relationships of areas of science, emerging research disciplines, the impact of particular researchers or institutions, etc.
- Often use documents as the “input data”

Fall 2017

CS 4460

45

Stanford Diss. Browser



9,000 Stanford PhD theses

Rather than overall semantic map, you choose a focus and all update to show their relationship

Demo at

<http://nlp.stanford.edu/projects/dissertations/>

Chuang et al
CHI '12

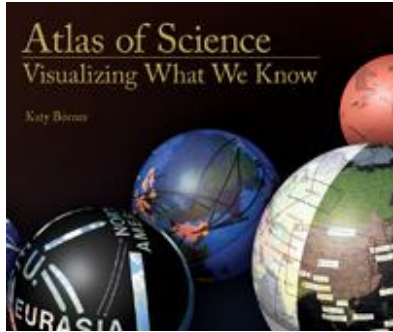
<http://nlp.stanford.edu/projects/dissertations/browser.html>

Fall 2017

CS 4460

46

Wonderful Book and Website



K. Börner



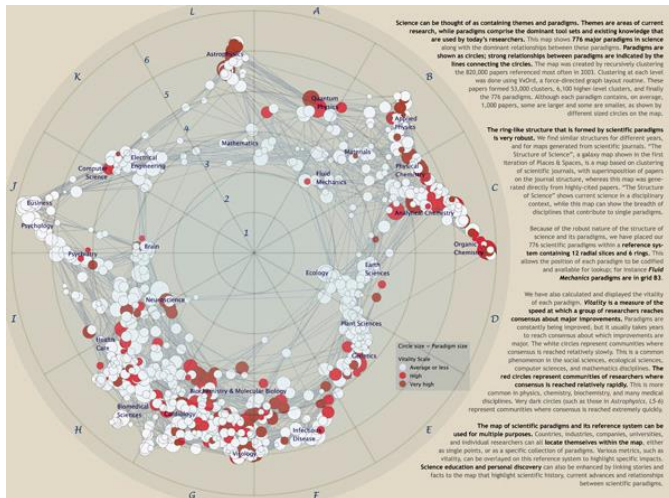
<http://scimaps.org>

Fall 2017

CS 4460

47

Some Examples



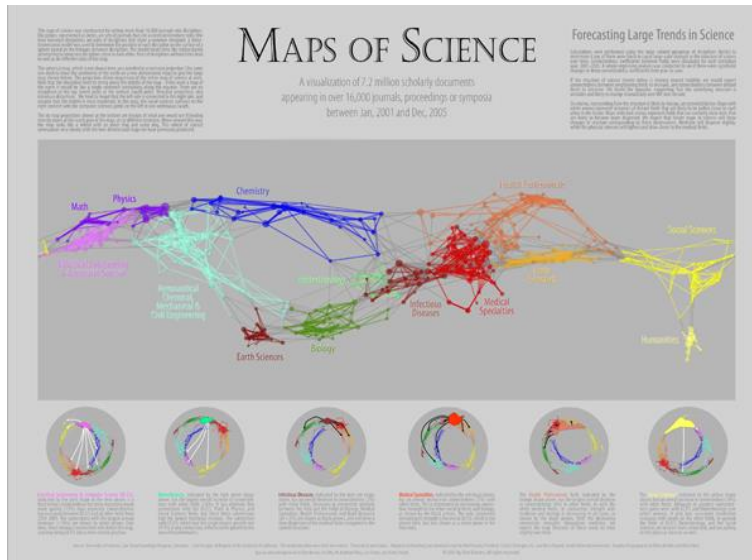
Boyack & Klavans

http://scimaps.org/maps/map_of_scientific_pa_55/

Fall 2017

CS 4460

48



Klavans & Boyack

http://scimaps.org/maps/map/maps_of_science_fore_50/

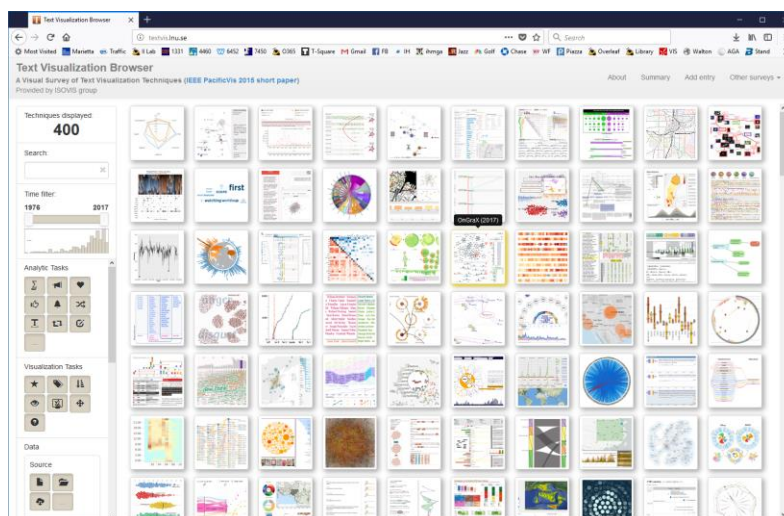
Fall 2017

CS 4460

49

<http://textvis.lnu.se>

Text Vis Browser



Fall 2017

CS 4460

50

Learning Objectives



- Explain what word concordance is & how WordTree representation works
- List different queries/tasks often needed on document collections
- List various analytic metrics often calculated on documents
- List different aspects of documents often visualized
- Explain vector space document analysis (similarity calculation, search) & TFIDF
- Describe visual representation used by and contributions of these systems
 - SentenTree, TextArc, Themail, Jigsaw, ThemeScape/IN-SPIRE, ThemeRiver

Fall 2017

CS 4460

51

Programs



- Some thoughts...

Fall 2017

CS 4460

52

Upcoming



- Lab 9 – Layout in D3 (networks)
 - Prep: Force-directed Graphs website

- Casual InfoVis
 - Prep: Watch InfoCanvas video