

Text and Document Visualization 1



CS 7450 - Information Visualization
October 19, 2015
John Stasko

Text is Everywhere



- We use documents as primary information artifact in our lives
- Our access to documents has grown tremendously in recent years due to networking infrastructure
 - WWW
 - Digital libraries
 - ...

Big Question



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

Fall 2015

CS 7450

3

Tasks/Goals



- What kinds of analysis questions might a person ask about text & documents?

Fall 2015

CS 7450

4

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.

Fall 2015

CS 7450

5

Related Topic - IR



- Information Retrieval
 - Active search process that brings back particular/specific items (will discuss that some today, but not always focus)
 - I think InfoVis and HCI can help some...
- InfoVis, conversely, seems to be most useful when
 - Perhaps not sure precisely what you're looking for
 - More of a browsing task than a search one

Fall 2015

CS 7450

6

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

Fall 2015

CS 7450

7

Challenge



- Text is nominal data
 - Does not seem to map to geometric/graphical presentation as easily as ordinal and quantitative data
- The “Raw data --> Data Table” mapping now becomes more important

Fall 2015

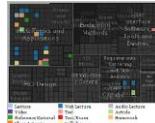
CS 7450

8

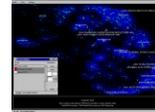
This Week's Agenda



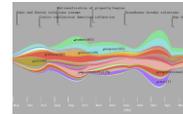
Visualization for IR
Helping search



Visualizing text
Showing words,
phrases, and
sentences



Visualizing document sets
Words, entities & sentences
Analysis metrics
Concepts & themes



Fall 2015

CS 7450

9

Information Retrieval



- Can InfoVis help IR?
- Assume there is some active search or query
 - Show results visually
 - Show how query terms relate to results
 - ...

Fall 2015

CS 7450

10

Generalize More



- How about the “holy grail” of a visual search engine?
 - Hot idea for a while
- My personal view: It’s a mistake in the general case. Text is just better for this.

Fall 2015

CS 7450

11

Search Visualization



<http://www.kartoo.com>

Defunct

Fall 2015

CS 7450

12

Sparkler



- Abstract result documents more
- Show “distance” from query in order to give user better feel for quality of match(es)
- Also shows documents in responses to multiple queries

Havre et al
InfoVis '01

Fall 2015

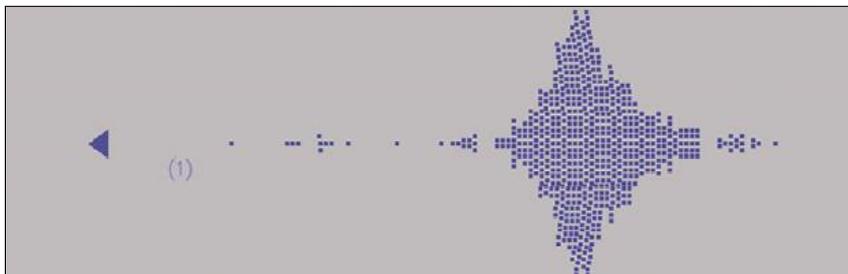
CS 7450

13

Visualizing One Query



- Triangle – query
- Square – document
- Distance between query and documents represents their relevance



Fall 2015

CS 7450

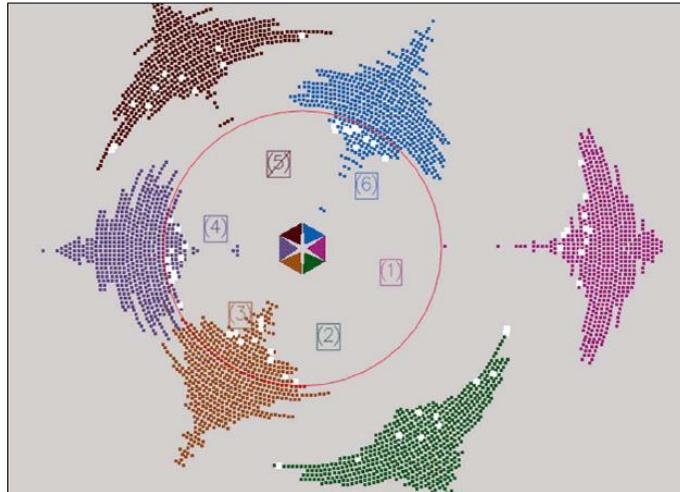
14

Visualizing Multiple Queries



Six queries here

Bullseye allows viewer to select quality results



Fall 2015

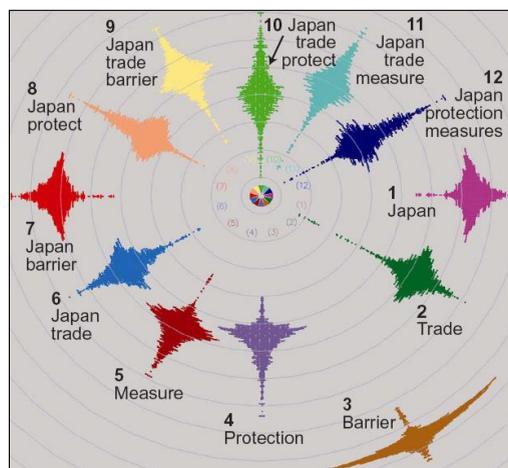
CS 7450

15

Test Example



- Text Retrieval Conference (TREC-3) test document collection
- AP news stories from June 24–30, 1990
- TREC topic: Japan Protectionist Measures
- Sparkler found 16 of 17 relevant documents

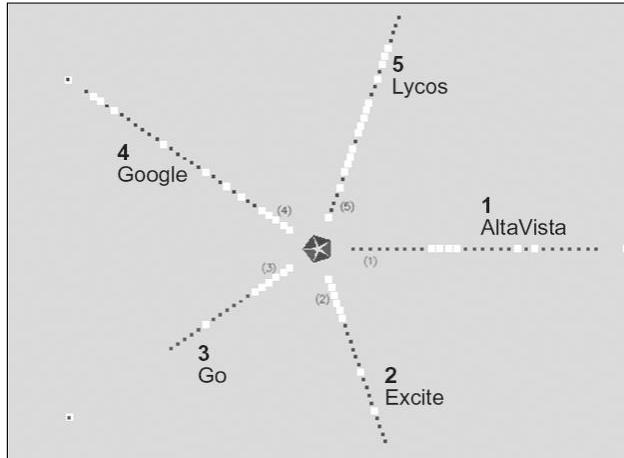


Fall 2015

CS 7450

16

Another Idea



Use it to compare search results from different search engines

Fall 2015

CS 7450

17

RankSpiral

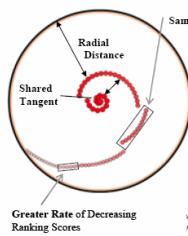
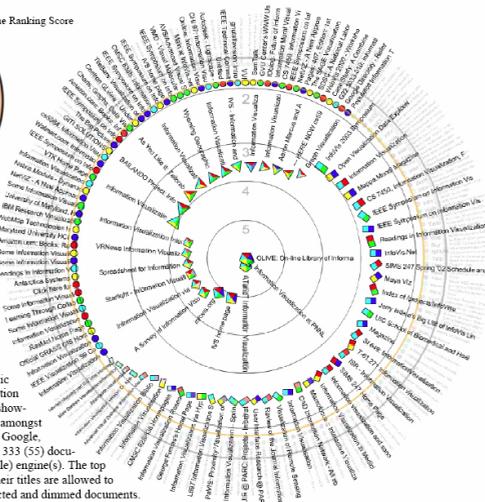


Figure 1. (Top) *RankSpiral* places consecutive document icons next to each other so that they do not overlap. Total ranking score of documents increases toward the center. Radial distance between documents that have the same angle can be used to display title fragments. (Right) shows a static *RankSpiral* that maximizes information density and minimizes occlusions, showing here the 388 unique documents amongst the top 100 documents retrieved by Google, Teoma, AltaVista, Lycos and MSN. 333 (55) documents were found by single (multiple) engine(s). The top 100+ documents are selected and their titles are allowed to extend across the remaining unselected and dimmed documents.



Color represents different search engines

Spoerri
InfoVis '04 poster

Fall 2015

CS 7450

18

Transition 1



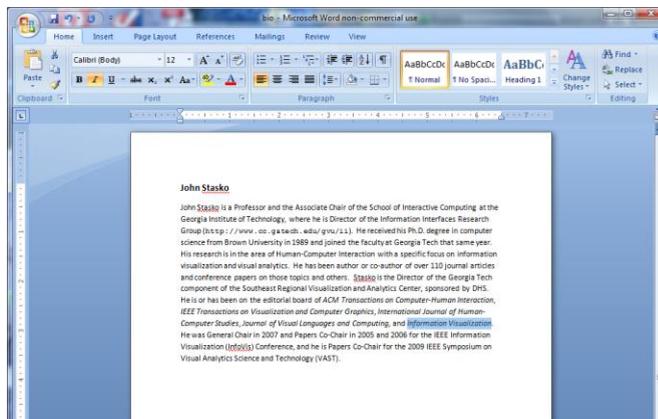
- OK, let's move up beyond just search/IR
- How do we represent the words, phrases, and sentences in a document or set of documents?
 - Main goal of *understanding* versus search

Fall 2015

CS 7450

21

One Text Visualization



Uses:
Layout
Font
Style
Color
...

Fall 2015

CS 7450

22

Tag/Word Clouds



- Currently very “hot” in research community
- Have proven to be very popular on web
- Idea is to show word/concept importance through visual means
 - Tags: User-specified metadata (descriptors) about something
 - Sometimes generalized to just reflect word frequencies

Fall 2015

CS 7450

25

History



- 90-year old Soviet Constructivism
- Milgram’s ‘76 experiment to have people label landmarks in Paris
- Flanagan’s ‘97 “Search referral Zeitgeist”
- Fortune’s ‘01 Money Makes the World Go Round

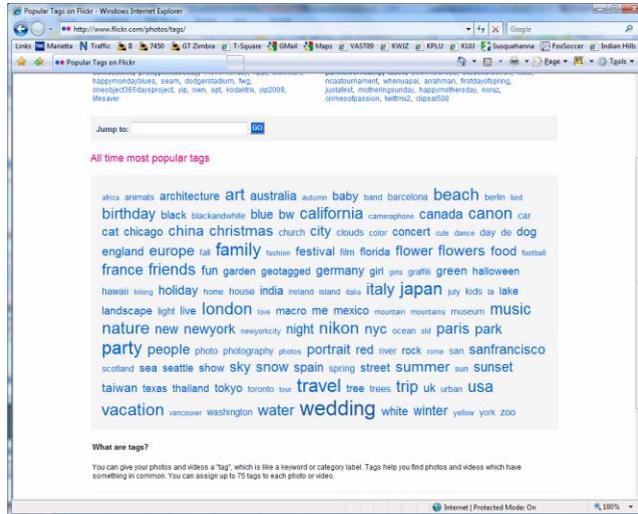
Viégas & Wattenberg
interactions ‘08

Fall 2015

CS 7450

26

Flickr Tag Cloud

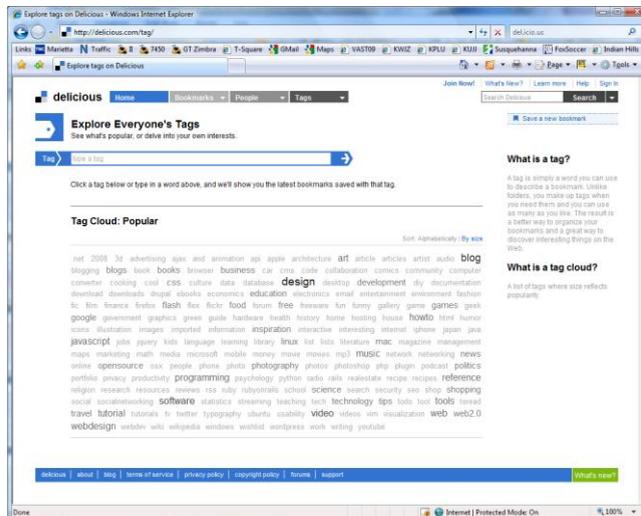


Fall 2015

CS 7450

27

delicious Tag Cloud

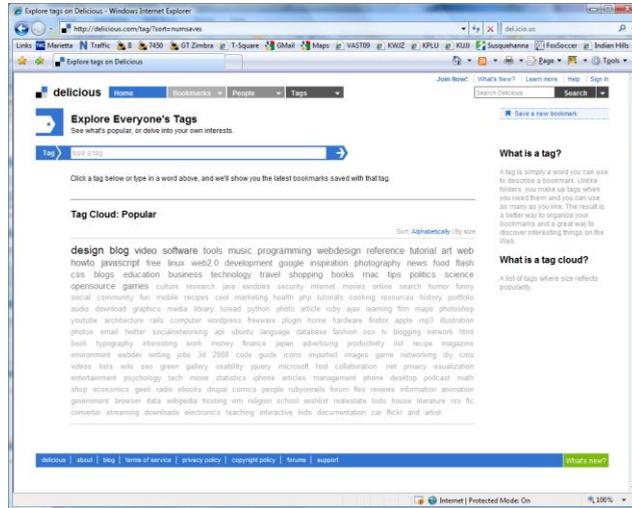


Fall 2015

CS 7450

28

Alternate Order

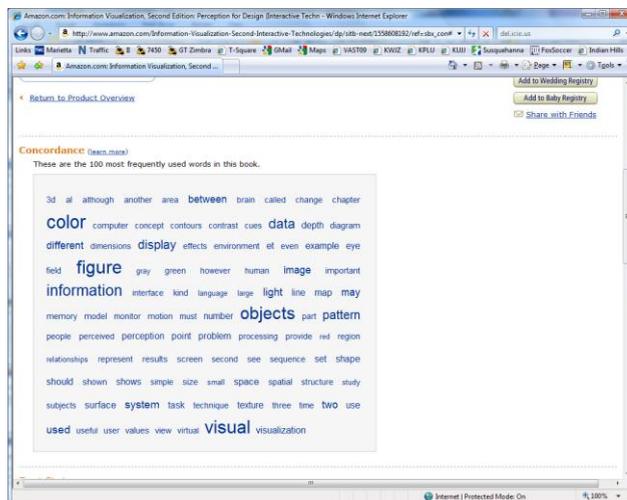


Fall 2015

CS 7450

29

Amazon's Product Concordance



Maybe now a "word cloud"

Fall 2015

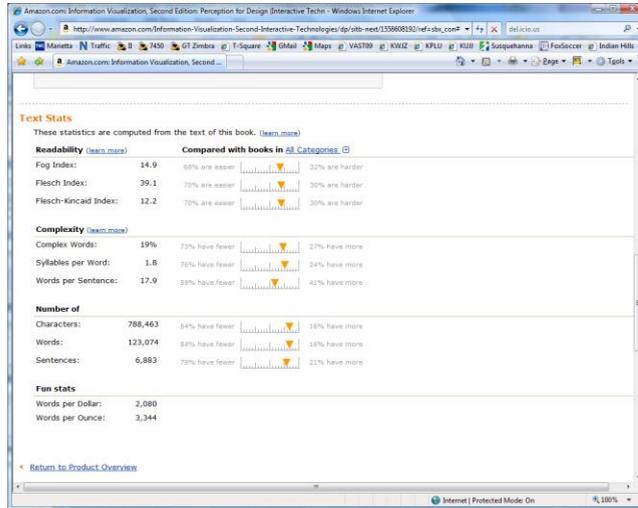
CS 7450

30

Sidenote



There are other types of info about a document on Amazon



Fall 2015

CS 7450

31

Many Eyes Tag Cloud



Here, pairs of words are shown



Fall 2015

CS 7450

32

Problems



- Actually not a great visualization. Why?
 - Hard to find a particular word
 - Long words get increased visual emphasis
 - Font sizes are hard to compare
 - Alphabetical ordering not ideal for many tasks
- Studies have even shown they underperform

Gruen et al
CHI '06

Fall 2015

CS 7450

33

Why So Popular?



- Serve as social signifiers that provide a friendly atmosphere that provide a point of entry into a complex site
- Act as individual and group mirrors
- Fun, not business-like

Hearst & Rosner
HICSS '08

Fall 2015

CS 7450

34

Wordle



- 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

Fall 2015

CS 7450

37

Layout Algorithm



- Details not published
- Idea:
 - sort words by weight, decreasing order
 - for each word w
 - $w.position := makeInitialPosition(w);$
 - while w intersects other words:
 - $updatePosition(w);$
 - Init position randomly chosen according to distribution for target shape
 - Update position moves out radially

Fall 2015

CS 7450

38

Fun Uses



- Political speeches
- Songs and poems
- Love letters (for “boyfriend points”)
- Wedding vows
- Course syllabi
- Teaching writing
- Gifts

Fall 2015

CS 7450

39

2-day Survey in Jan. 09



- 2/3 respondents were women
- Interest came from design, visual appeal, beauty
- Why preferred over word clouds:
 - Emotional impact
 - Attention-keeping visuals
 - Organic, non-linear
- Fair percentage didn’t know what size signified

Fall 2015

CS 7450

40

Wordle Characteristics



- Layout, words are automatic
- If you had some control, what would you like to change or alter?

Fall 2015

CS 7450

43

Mani-Wordle



- Start with nice default algorithm
- Give user more control over design
 - Alter color (within a palette)
 - Pin words, redo the rest
 - Move and rotate words
 - Smooth animation and collision detection for tracking changes

Koh et al
TVCG (InfoVis) '10

Fall 2015

CS 7450

44

Multiple Documents?



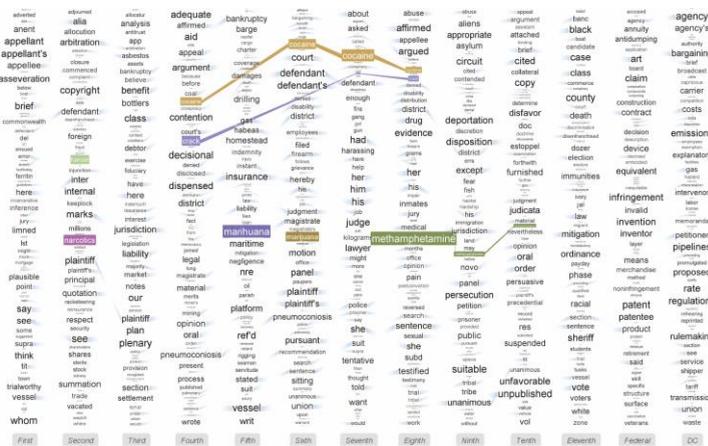
- How to show word frequencies across multiple related documents?

Fall 2015

CS 7450

47

Parallel Tag Clouds



Video

Different circuit courts

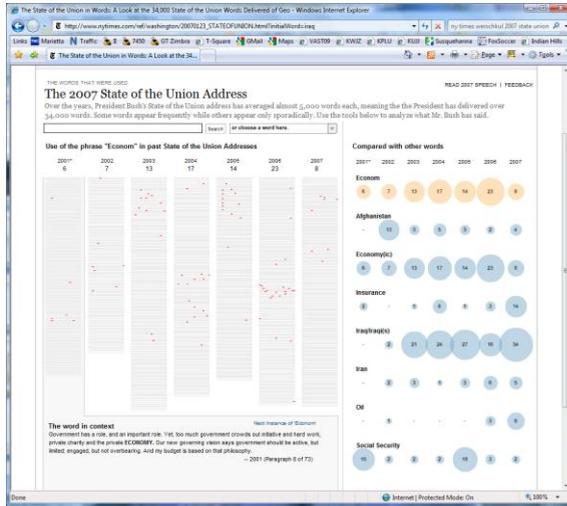
Collins et al
VAST '09

Fall 2015

CS 7450

48

Overview & Timeline



State of the Union Addresses

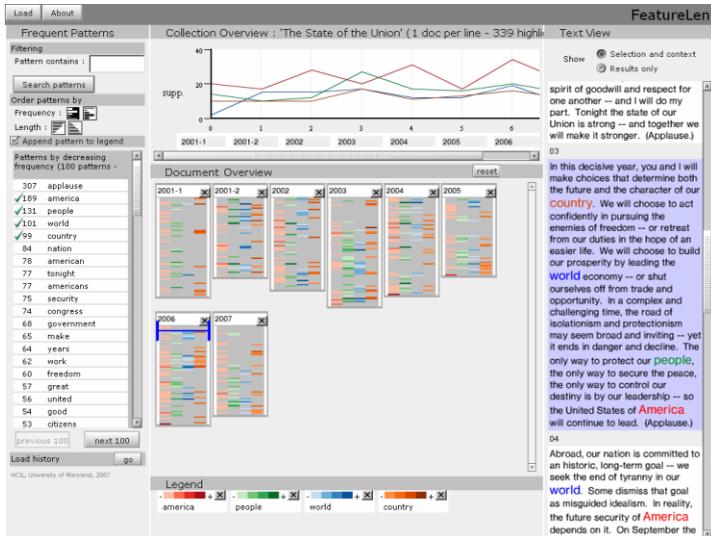
http://www.nytimes.com/ref/washington/20070123_STATEOFUNION.html?initialWord=iraq

Fall 2015

CS 7450

FeatureLens

Video



Show patterns of words or n-grams

Don et al
CIKM '07

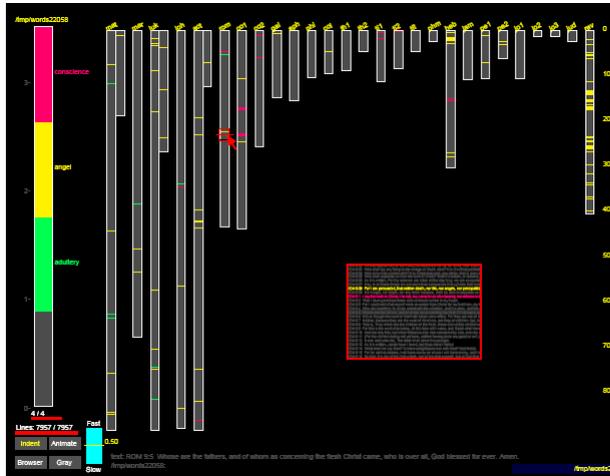
<http://www.cs.umd.edu/hcil/textvis/FeatureLens/>

Fall 2015

CS 7450

52

SeeSoft Display



Like taping text
to the wall and
walking far away

New Testament

Eick
Journal Comput. & Graph. Stats '94

Fall 2015

CS 7450

53

Beyond Individual Words



- Can we show combinations of words, phrases, and sentences?

Fall 2015

CS 7450

54

Concordance



Definition

The screenshot shows the Merriam-Webster Online Dictionary page for the word "concordance". The page includes a navigation menu on the left, a search bar at the top, and a main content area with the following information:

- concordance** (One entry found.)
- Concordance** (Sponsored Links): Find the Benefits of Concordance Software by LexisNexis. Buy Now! law.lexisnexis.com
- Main Entry:** **con·cor·dance**
- Pronunciation:** 'kan-kor-'dʌn(t)s, kən-'
- Function:** *noun*
- Eymology:** Middle English, from Anglo-French, from Medieval Latin *concordantia*, from Latin *concordant-*, *concordans*, present participle of *concordare* to agree, from *concord-*, *concor-*
- 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 2015

CS 7450

55

Concordance in Text



The screenshot shows the Concordance software interface. The main window displays a list of words and their contexts. The list is as follows:

Headword	No.	Context...	Word	...Context	Reference
HEAR	15	That my own	heart	drifts and cries, having no...	Deep Analysis
HEARD	9	By the shout of the	heart	continually at work	And the wave
HEARING	7	Nothing to adapt the skill of the	heart	to, skill	And the wave
HEARS	3	The tread, the beat of it, it is my own	heart	,	Träumerei
HEARSE	1	Because I follow it to my own	heart	,	Many famous
HEART	25	My	heart	is ticking like the sun:	I am washed i
HEART'S	2	The vague	heart	sharpened to a candid co...	The March Pa
HEART-SHAPED	1	Contract my	heart	by looking out of date.	Lines on a Yo
HEARTH	1	Having no	heart	to put aside the theft	Home is so Se
HEARTS	7	And the boy puking his	heart	out in the Gents	Essential Bea
HEARTY	1	A harbour for the	heart	against distress.	Bridge for the
HEAT	6	These I would choose my	heart	to lead	After-Dinner F
HEAT-HAZE	1	Time in his little cinema of the	heart	,	Time and Spa
HEATH	1	This petrified	heart	has taken,	A Stone Churc
HEATS	1	How should they sweep the girl clean...	heart	,	I see a girl dra
HEAVE	1	Hands that the	heart	can govern	Heaviest of fl
HEAVEN	4	For the	heart	to be loveless, and as col...	Down
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

The interface also includes a menu bar (File, Text, Search, Edit, Headwords, Contexts, View, Tools, Help), a toolbar with various icons, 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), Context sort: Asc occurrence order.

<http://www.concordancesoftware.co.uk>

Fall 2015

CS 7450

56

Word Tree



Fall 2015

CS 7450

From King James Bible

57

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 2015

CS 7450

58

Another Challenge



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

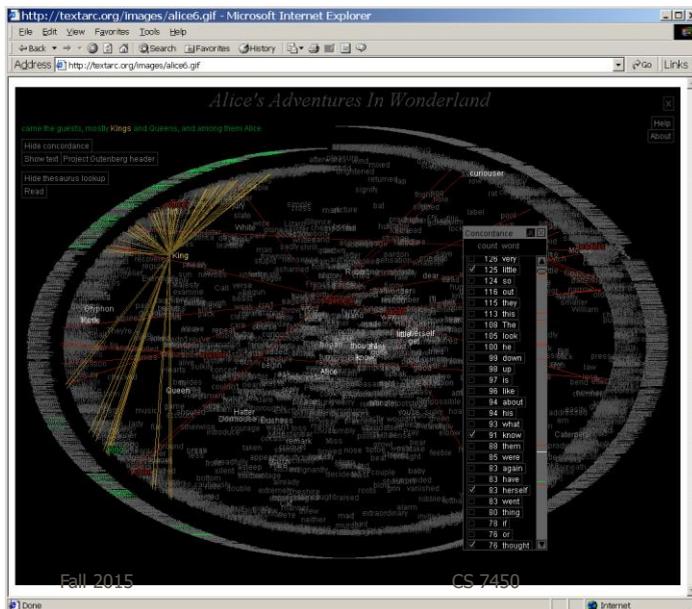
Fall 2015

CS 7450

65

TextArc

<http://textarc.org>



Sentences laid out in order of appearance

Words near to where they appear

Significant interaction

Brad Paley

66

Next Time



- More about collections of documents and showing other characteristics of documents
 - Analysis metrics
 - Entities
 - Concepts & themes

Project Design Document



- See project description & rubric
- My additions...
- Bring 2 copies
- Questions?

Survey Results



- Hours outside of class
 - <1,2,3-5,4,4,4,5,5,5,5,5,5-20,6,6,6,6-8,8,8,8*,10,10,10,11
 - *for useless meetings
- When in class, I feel
 - Interested (22/25)
 - Engaged (19/25)
 - ...
 - Confused (2/25)
 - Bored (1/25)

Fall 2015

CS 7450

69

Survey Results



- Pace?
 - About right – strong winner
 - Little too fast (4)

Fall 2015

CS 7450

70

Survey Results



- Expectations clear?
 - Yes (15)
 - More rubrics (3)
 - “I really understand the lectures given by instructor but I am so confused by assignments and final project. Assignments instruction are extremely unclear and even undefined in some cases. I really do not know what it takes to succeed in this course.”
 - “Not at all. Assignments description are really bad. I really have no idea what should I work on? what should I write? how to evaluate tools? It is not also clear how TA are doing gradings. I talked to them couple of times they gave some reasons but the reasons are really might change person to person. The same thing is happening for the final project of this course nothing is clear. If I could I would drop the course.”
- Subjectivity

Fall 2015

CS 7450

71

Survey Results



- Good?
 - “Class demos are awesome! More hands on sessions which involve all the students in the class in group activities will be awesome. And make random group assignments for these group activities. This helps us know the other people in the class. I only know 1/3rd of the students in this class.”
 - “The homeworks have been very helpful - particularly homeworks 3 and 4. D3 is difficult, but I think it is a very valuable skill to have. I also think critiquing commercial tools was a very useful exercise to think critically about the visualization decisions made. I have also really loved seeing the short demo videos. That's interesting and inspiring.”
 - “while it can be time consuming to read articles for each class, the material is very useful going into lecture. I am not certain how useful I found HW4 to be... while it was interesting to critique existing tools and systems, I would have liked to have more homework on how to design clear and effective data visualizations or perhaps spending more time figuring out what types of visualizations to use with certain data sets.”

Fall 2015

CS 7450

72

Survey Results



● Bad?

- "I hope that the questions asked to the class could be a little more difficult or thought-provoking, maybe more open-ended if that makes sense? / it's always nice to bring some kind of controversy and discussion into the issues presented, although maybe that's not best for time management."
- "This class needs more short group projects and short group activities in class with random group formation. This will keep things fresh and the students will take more interest in the course."
- "Perhaps the class can be offered every sem resulting in smaller classes. Helps with having more in class activities like design critiques, discussing visualizations etc."
- "practice with what types of visualizations to use to communicate data sets and intent."
- "more assignments which make me build visualizations"
- "Off the top of my head, maybe a bit more difficult assignments. But I understand that many in the class don't have programming experience and it might be too difficult for them."
- "More examples of past course projects. What teams did well, what types of viz's they used correctly and incorrectly. Examples are always useful."

Fall 2015

CS 7450

73

Survey Results



● Bad?

- "The group project + homeworks + readings + final exam is too much for me and others I have talked to. If you have a big group project, I feel an exam is unnecessary."
- "The course is heavy. A lot of things going on at the same time from readings, projects, hw's. It is difficult to get time to focus well on one area(eg the project). I feel there should just be an end project instead of a final exam like all of my other classes this semester."
- "Make the whole class homework and project based and more programming related. The final exam is a hindrance that must be done away with in this day and age."
- "My only criticism so far is the one reading assignment that was chapters 4-12 of the textbook. That was a bit too much at once that the information didn't "sink in" as much as it would have if it had been a bit more spread out (into two separate reading assignments, for instance). It was a bit overwhelming. Other than that, I have no criticisms."

Fall 2015

CS 7450

74

Survey Results



- **Bad?**
 - “At the beginning of the semester the instructor mentioned that this course is going to be lots of work. I said okay... I would challenge myself and learn many things. yes this course needs lots of time in a useless way. I ended up spending most of my time meeting TAs and my groupmates for this course and talking what would please the instructor and what would not. This has several reasons that I am going to explain it as bellow: First: Instructor needs to clarify his expectations for every homework or project in a very clear way. For example, in many courses system that we develop needs to have a certain functionalities. So that We try to develop that and we are clear that function must work otherwise we are not getting our grade. This is the first time that I took a class and definition of everything is summarized in a single word "interesting". Interesting design, interesting dataset. Second: Instructor mentioned the life is unfair at the beginning of the semester but he makes it even more unfair. Above 70% of the class are master HCI master students with the minimum knowledge of programming. So in our group we have three designers and two programmers. Designers all have different definition of "interesting design", and "interesting dataset". They everytime end up coming up with some imaginary designs and I need to explain them how difficult is to implement their design. What many professors in CS department are doing is that they will ask students to fill up a form explaining their skills and then the instructor and his TAs will distripute the groups in a way which is fair to everyone.”

Fall 2015

CS 7450

75

Survey Results



- **What else could you do?**
 - More reading
 - More practice
 - Meet with TAs/instructor more

Fall 2015

CS 7450

76

Survey Results



- Additional comments?

- “I do not like that there is a final exam for this type of class. This is more of a group project based class and it doesn't make much sense to have 0 tests/quizzes and then have a final at the end. Very difficult to know what to expect for an exam that is 20% of the grade.”
- “In a grad school class, I feel final exams are unnecessary when we are all trying to put a lot of work into a final project. It is archaic.”
- “The amount of papers, books and articles that this class expects the students to read on their own is staggering. Students have other difficult classes on their plate and cannot read all the papers and recommended readings. Efforts must be done to summarize those required readings so that students can read the summaries to save time.”

Survey Results



- Additional comments?

- “I hope that there will be some mechanism for reviewing group members at the end of the term project. Perhaps this was already the plan, but if not, I hope it will be possible. It becomes apparent early on which group members care about the class, the project, and their grades, and which ones can't be bothered to make an effort. I hope there will be some way to reflect on the contributions on other group members, so that those who have helped hold the group together and those who haven't done their fair share can be graded accordingly.”

Upcoming



- Text and Documents 2
 - Reading
Keim & Oelke '07

- Time Series Data
 - Reading
Aigner et al '08

Fall 2015

CS 7450

79

References



- Marti Hearst's i247 slides
- All referred to papers

Fall 2015

CS 7450

80