# Text and Document Visualization 1

CS 7450 - Information Visualization
October 31, 2016
John Stasko

# Learning Objectives

- Explain key challenges in visualizing a large document or body of text
- Identify and explain different techniques for representing words and concepts in a document
  - Word cloud, Wordle, Parallel tag cloud, SeeSoft, WordTree, PhraseNet, SentenTree, TextArc
- Understand the positives and limitations of word clouds and Wordles
- Describe SeeSoft-style miniature visual representations
- Explain what word concordance is
- Describe how WordTree representation works
- Identify and explain the techniques:
  - Word cloud, Wordle, Parallel tag cloud, SeeSoft, WordTree, PhraseNet, SentenTree, TextArc

1

# 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?

# 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

# 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
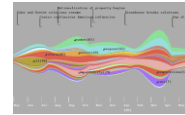
3

# 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



Today

Next
time

# Information Retrieval

- Can InfoVis help IR?

- Assume there is some active search or query
  – Show results visually
  – Show how query terms relate to results
  – …

# 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.

# Search Visualization



http://www.kartoo.com

Defunct

# 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 to-ward 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, show-ing here the 388 unique documents amongst the top 100 documents retrieved by Google, Teoma, AltaVista, Lycis 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

# To Learn More



Marti Hearst's Book

Chapter 10

http://searchuserinterfaces.com/book/

6

# 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

# One Text Visualization



Uses:
Layout
Font
Style
Color
…

# Design Challenge

- How would you visualize one of the recent presidential debates?

- Ideas?

# Tasks

- What kinds of questions or tasks would someone want to do with such a visualization?

8

# Word Counts

# More Word Counting



http://www.wordcount.org

# 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

# History

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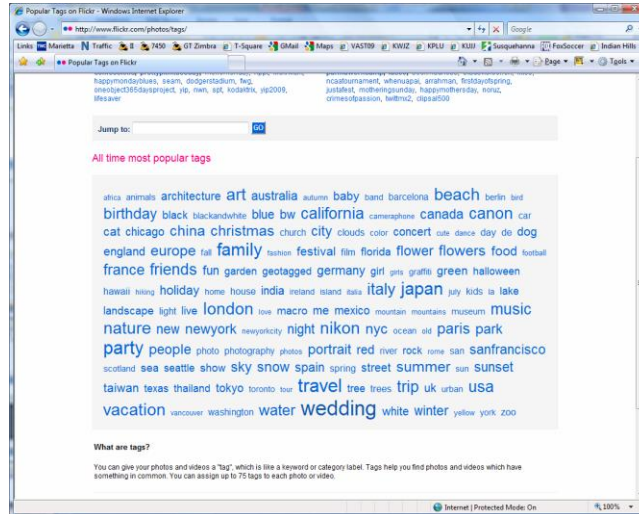- Fortune's '01 Money Makes the World Go Round

Viégas & Wattenberg
*interactions* '08

# Flickr Tag Cloud

# delicious Tag Cloud
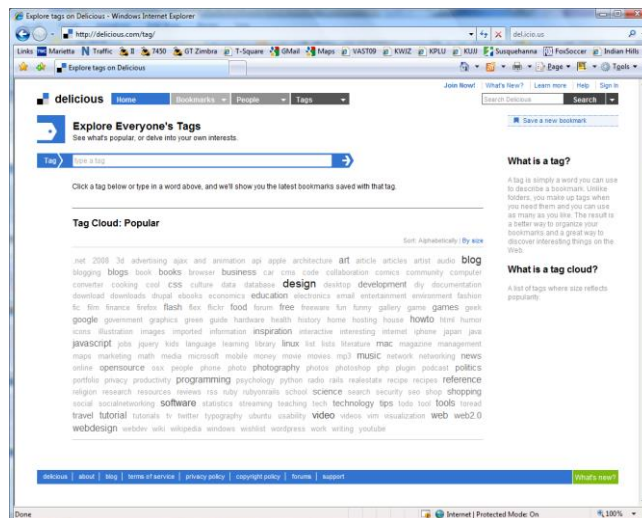
# Alternate Order

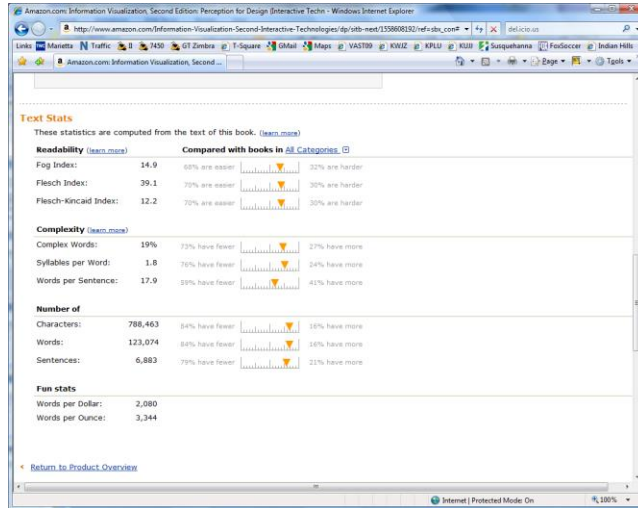# Amazon's (old) Product Concordance



Maybe now a
"word cloud"

# More (old) Info

There are other types of info about a document on Amazon

# Many Eyes Tag Cloud

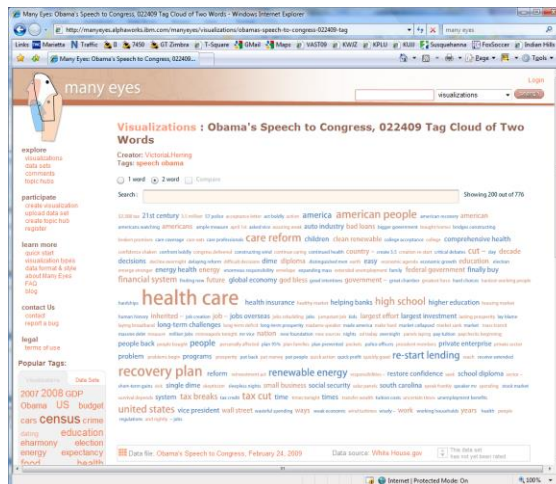Here, pairs of words are shown

13

# 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

http://www.niemanlab.org/2011/10/word-clouds-considered-harmful/

**NiemanLab**

**Word clouds considered harmful**

The New York Times senior software architect would like the newest "mullets of the Internet" to go back from whence they came.

By JACOB HARRIS   @harrisj   Oct. 13, 2011, 1:45 p.m.

In his 2003 novel *Pattern R...* named Cayce Pollard with a ... allergic to brands. Even the...

is a shoddy visualization that fails all the principles I hold dear.

> Every time I see a word cloud presented as insight, I die a little inside.

For starters, word clouds support only the crudest sorts of textual analysis, much like figuring out a protein by getting a count only of its amino acids. This can be wildly misleading; I created a word cloud of Tea Party feelings about Obama, and the two largest words were implausibly "like" and "policy," mainly because the importuned word "don't" was automatically excluded.

14

# Why So Popular?

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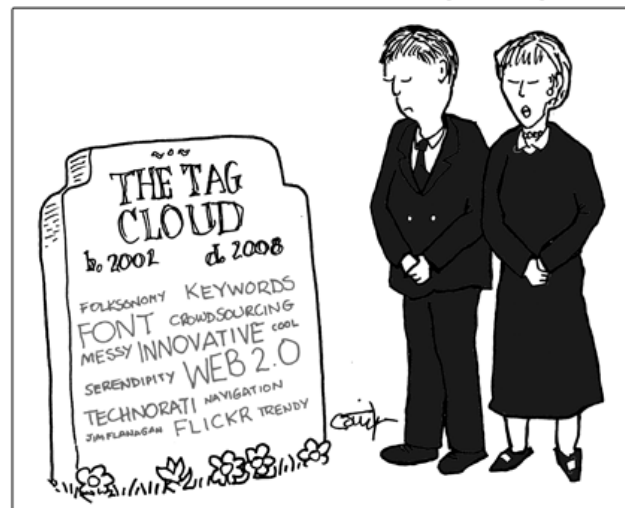- Act as individual and group mirrors
- Fun, not business-like

Hearst & Rosner
HICSS '08

http://www.socialsignal.com/system/files/images/2008-08-01-tagcloud.gif

15

# Wordle

# 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

# 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

# Fun Uses

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

# 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

# SoTU Wordles

# A Little More Order
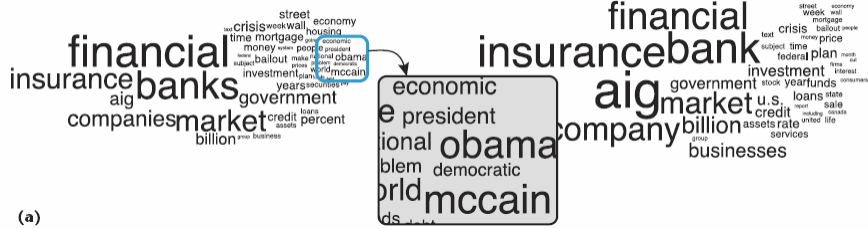


(a)

Order the words more by frequency

Cui et al
*IEEE CG&A* `10

# Semantic/Context Word Clouds



Paulovich et al
*Computer Graphics Forum* '12

Wang et al
Graphics Interface '14

Wu et al
*Computer Graphics Forum* '11

# Wordle Characteristics

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

# 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

# Video

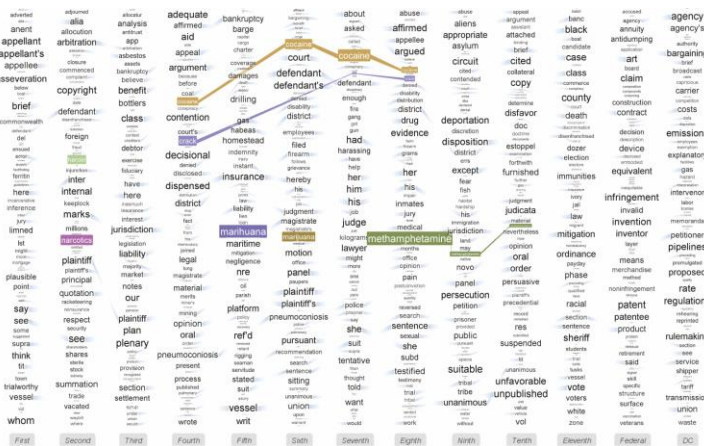# Text Analysis on Web



http://voyant-tools.org/

# Multiple Documents?

- How to show word frequencies across multiple related documents?

# Parallel Tag Clouds



Video

Different circuit courts

Collins et al
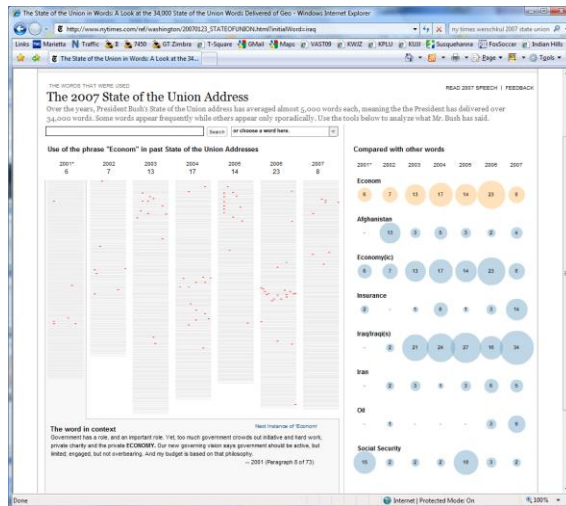VAST '09

# Analytic Support

- Note: Word Clouds and Wordles are really more overview-style visualizations
  - Don't really support queries, searches, drill-down

- How might we also support queries and search?
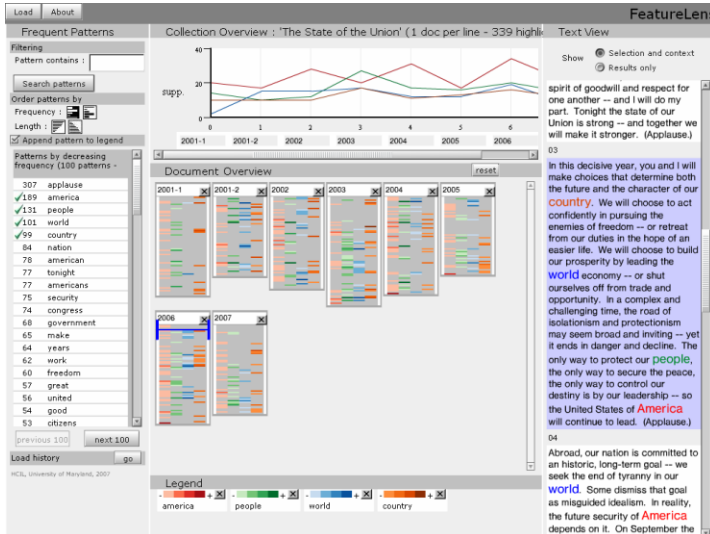
# Overview & Timeline



State of the Union Addresses

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

# FeatureLens

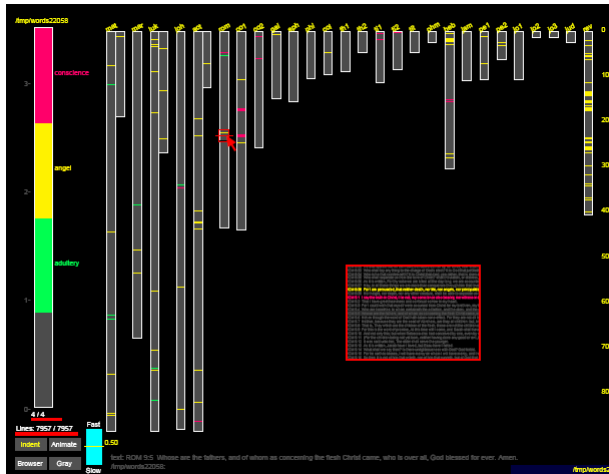Video

Show patterns
of words or
n-grams

Don et al
CIKM '07

http://www.cs.umd.edu/hcil/textvis/featurelens/

# SeeSoft Display

Like taping text
to the wall and
walking far away

New Testament

Eick
*Journal Comput. & Graph. Stats '94*

# Beyond Individual Words

- The previous techniques 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?

# Concordance



Definition

# Concordance in Text

# Word Tree

26

# 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

# Interaction

# Many Eyes' WordTree

# Phrase Nets

- Examine unstructured text documents
- Presents pairs of terms in phrases such as
  - X and Y
  - X's Y
  - X at Y
  - X (is|are|was|were) Y
- Uses special graph layout algorithm with compression and simplification

van Ham et al
*TVCG* (InfoVis) '09

# Examples



Fig 4. Matching the same pattern on different texts. Here we used the pattern "X of Y" to compare the old and new testaments. Israel takes a central place in the Old Testament, while God acts as the main pattern receiver in the New Testament.

# Examples



Fig 5. Matching different patterns on the same text. Here we analyzed Jane Austen's *Pride and Prejudice* with "X and Y" and "X at Y" respectively. The left image shows relationships between the main characters amongst others, while the right image shows relationships between locations.

# User Interface



Fig 3. The Phrase Net user interface applied to James Joyce's Portrait of the Artist as a Young Man. The user can select a predefined pattern from the list of patterns on the left or define a custom pattern in the box below. This list of patter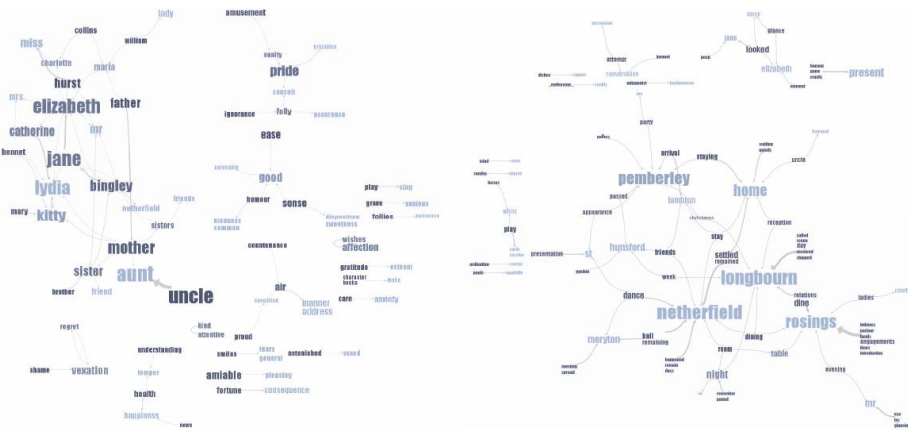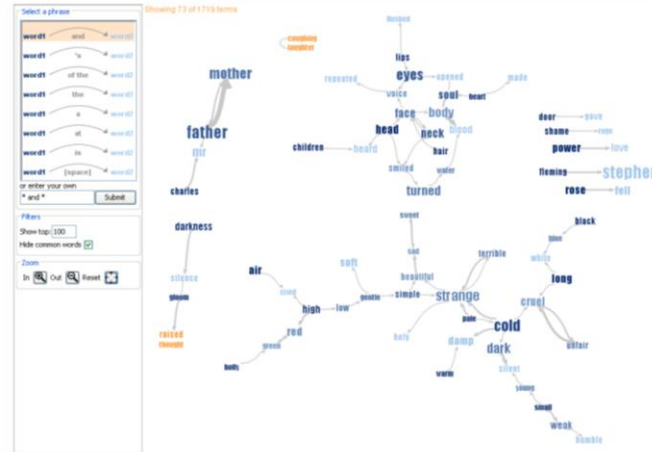ns simultaneously serves as a legend, a list of presets and an interactive training mechanism for regular expressions. Here the user has selected "...X and Y...", revealing two main clusters, one almost exclusively consisting of adjectives, the other of verbs and nouns. The highlighted clusters of terms have been aggregated by our edge compression algorithm.

# 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?

# 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

# Example



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
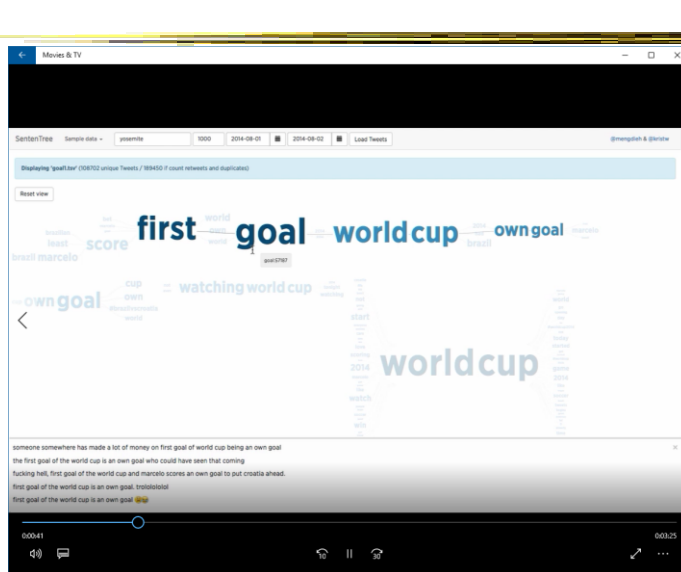
(Interaction is key & not shown here)

# Example



Tweets mentioning word "Yosemite" from Aug 1, 2014

# Video

# Another Challenge

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

# TextArc

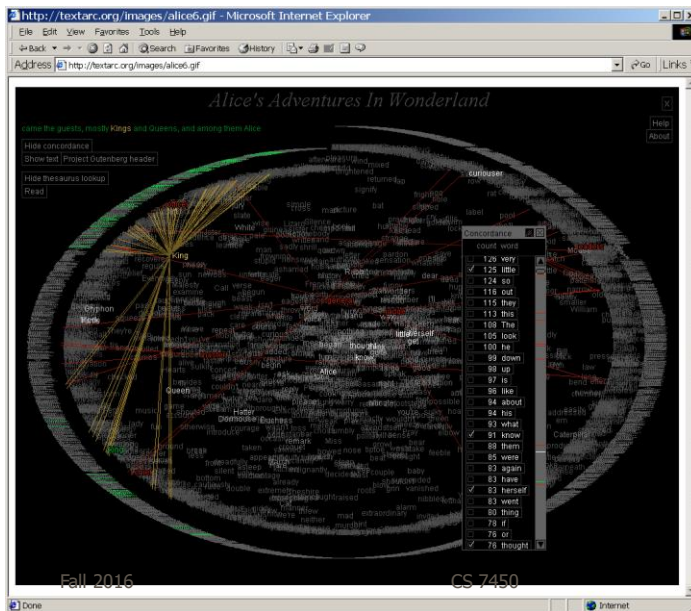http://textarc.org



Sentences laid out
in order of appearance

Words near to where
they appear

Significant interaction

Brad Paley

33

# Next Time

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

# Learning Objectives

- Explain key challenges in visualizing a large document or body of text
- Identify and explain different techniques for representing words and concepts in a document
  - Word cloud, Wordle, Parallel tag cloud, SeeSoft, WordTree, PhraseNet, SentenTree, TextArc
- Understand the positives and limitations of word clouds and Wordles
- Describe SeeSoft-style miniature visual representations
- Explain what word concordance is
- Describe how WordTree representation works
- Identify and explain the techniques:
  - Word cloud, Wordle, Parallel tag cloud, SeeSoft, WordTree, PhraseNet, SentenTree, TextArc

# Reading

- Viegas & Wattenberg '08

# Project

- Meetings with TAs
  - Once every two weeks
  - Will be logged
  - Conveys an impression

# Upcoming

- Text and Documents 2

- Hierarchical (Tree) Data 1

# References

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