# Text and Documents

## CS 4460 - Information Visualization

Jim Foley, some material courtesy John Stasko. Some examples from Marti Hearst, *Search User Interfaces*, Cambridge University Press, 2009

Last update: October 2016

1. You have the following information on 387 articles published in 65 issues of a newsletter for in the months of June, July and August for a given year. Each newsletter contains 5 to 7 articles

   Date published: day/month/year

   Subject category of article: Politics, Sports, Economics, Health, Science, Other

   Author of article: Walt, Mary, John, Beth

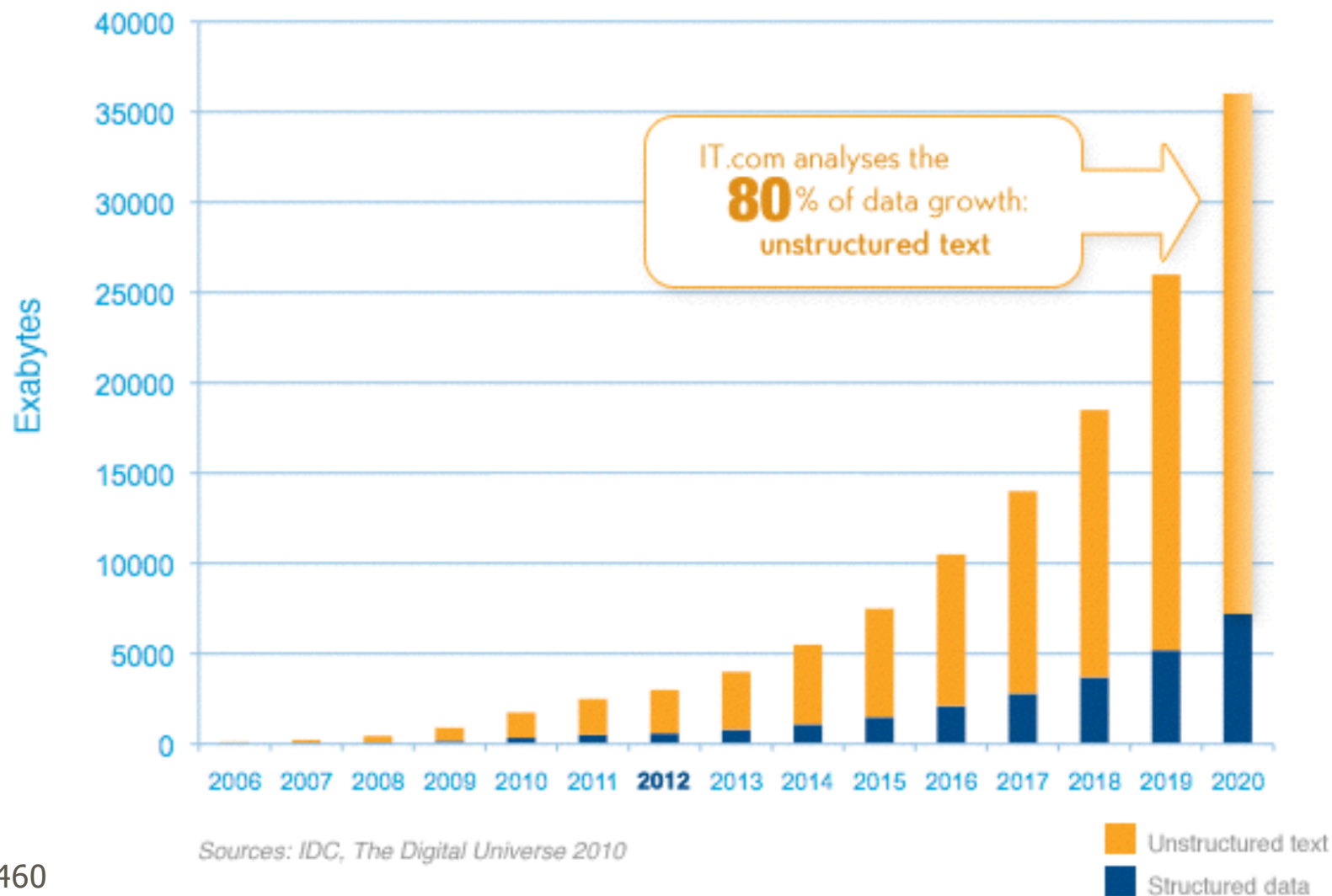   Focus of article: US, global, not known

Sketch a dashboard (several potentially linked visualizations) to convey as much of the above information as possible about the 387 articles. Ideally, you should convey all of the information.

2. There are several programs to detect plagiarism in term papers, by comparing a submitted term paper to a large collection of previously-submitted term papers.  If a possible case of plagiarism is found, we need a way to show how the submitted paper is similar to a previously-submitted paper. Sketch a way to do this.

# Who Cares about Text?

## Worldwide Corporate Data Growth

IT.com analyses the **80**% of data growth: unstructured text

Exabytes (y-axis: 0 to 40000)

Years (x-axis): 2006, 2007, 2008, 2009, 2010, 2011, **2012**, 2013, 2014, 2015, 2016, 2017, 2018, 2019, 2020

Sources: IDC, The Digital Universe 2010

Legend:
- Unstructured text (orange)
- Structured data (blue)

# The Key Question for InfoVis

- How can InfoVis help users in gathering, understanding, using information from
  - Document collections (*macro-level*)?

    Everything on the web, library, …
  - Individual (or a few) documents (*micro-level*)?

    Such as a thesaurus, or a book or speech,

    Shakespeare, Bible, Qur'an, Torah, ….

- Documents include email, all social media, news, books, magazines, …..

# Text InfoViz Examples

What is each doing?

# Some Questions about Document(s)

- Summarization: what are the main ideas in document(s)?
- Sentiment analysis: what sentiment(s) document(s) express?
- Trends: how do words/ideas/sentiment change over time in documents?
- How are words/ideas/themes distributed through one or several document?
- Grammar analysis: what grammatical patterns are used in document(s)
- How do two (or more) docs compare?

Some more micro-level, some more macro-level

# More Questions about Document(s)

- Search: which docs have certain keywords/topics?
- How similar are two or more docs? In what ways?
- How do documents fit into a larger context?
- How does one document compare to or relate to other documents?
- In what docs (where in docs) is the word "inflation" used close to (next to) the word "spending?"
- Event analysis: what event(s) are reported in document(s)

# Text Visualization Browser

A Visual Survey of Text Visualization Techniques
Provided by ISOVIS group

http://textvis.lnu.se/

About    Summary    Add entry    Othe

**Techniques displayed:**
**184**

**Search:**

**Time filter:**
2001    2013

...

**Analytic Tasks**

**Visualization Tasks**

**Data**

**Source**

Note multiple filters

© ISOVIS group 2014-2016. All rights for the tech...                    ...tive owners.
Version 1.7.5. Last updated: September 7, 2016
This website is using Google Analytics (for statistic...

# Challenge

- Text is nominal data with a hugh (infinite) cardinality
  - Does not map to visual encodings as easily as nominal, ordinal & quantitative data
- It is not the text itself that gets encoded, it is information about the text
- Must first extract information from the text
  - Exception – text meta-data (more later)
- Visualizations are of the extracted info.

# Pipeline for Text/Doc Visualization

Raw Data
(documents)

Per document
　Keyword counts
　Sentiments
　Events
　….

Geometric rep'n
of desired
information

NLP ⟶ Analysis ⟶ Visualization

Natural Language Processing
　Keywords
　Key phrases
　Grammar
　Sentiments
　Statistics
　Entity extraction
　…..

Find relevant
docs
Compare docs

# Challenge (Cont'd)

- Meta-data: data about data
- *Unstructured text* does NOT have any explicit meta-data.
  - Just that infinitely big collection of nominal data
  - NLP can extract some meta-data, such as dates
- Contrast to *structured text* of an on-line library with explicit meta-data such as
  - Author name
  - Year of publication
  - Title
  - ISBN number
  - Library of Congress number
  - Publisher name
  - Etc
- Some meta-data also nominal but lower cardinality than free text
  - Simplifies retrieval and visualization process.
- Will see  some meta-data driven examples

# Outline

- Macro-level – searching larger document collections
  - Unstructured – no meta-data
  - Structured – explicit meta-data
- Micro-level
  - Inter-document methods for smaller document collections
    How do retrieved documents relate to a query?
    How do retrieved documents relate to one another?
  - Intra-document methods
    Word usage, grammatical style, …

- With the caveat that some methods can be used in multiple ways

# Macro-Level: Large Unstructured

- LARGE does not mean entire WWW!!
- A number of systems endeavor to give a "big picture view" – the "gist" of a large collection of documents
  - Themescape
  - WebThemes
  - Galaxies
  - ThemeRiver
  - Sentiment/emotion over time

# Themescape



Height/color encode document density

Received 2016 "test of time" award

# TopicLens



No Lens | Topic Lens

Vis Papers

Hover 'Lens' over cluster, see sub-clusters

field flow vector method tensor

field method scalar

graphic applic object

system inform present

user inform interact   render imag base techniqu

interact system explor

model simul interact base view

design tool analyt user inform

interact user explor method view

surfac method mesh algorithm model

set techniqu larg present task

inform displai map space color

graph network ... le layout edg

Docs color-coded & clustered by topic

Video Next

# TopicLens Video

# ThemeRiver

# Visual Backchannel for Large-Scale Events



Fig. 1. The Visual Backchannel—here shown for Twitter posts about the event Park(ing) Day—consists of a) Topic Streams: a visualization representing topical development, b) controls for filtering and searching, c) People Spiral indicating the activity of participants, d) chronologically ordered list of posts, and e) Image Cloud displaying shared photos.

# Emotions/trends/sentiments over Time

- Theme River style
- Many small multiples
- Animation
- Sparkclouds
- Parallel wordclouds

# Emotions from tweets

- A: Twenty emotions
- B: Thickness = number of tweets analyzed
- C: Star (radar) plot of emotions

# Visualizing Emotions in Event-Related Tweets



Figure 2: The detailed view showing the Women's Gymnastics Floor Exercise Final. A - Emotion wheel showing the emotion profile of the current time interval; B - Timeline visualizing the emotion flow; C - Button to stop/resume the animation; D - Tweets of the current time interval; E - Video; F - Background with color of the dominant emotion.

# Trends via "Sparkclouds"



Sparkclouds: visualizing trends in tag clouds. Lee et al, IEEE TRANSACTIONS ON VISUALIZATION AND COMPUTER GRAPHICS, Nov/Dec 2010

# Trends via Parallel Wordclouds

# Pros/Cons

- Theme River style
- Many small multiples
- Animation
- Sparkclouds
- Parallel wordclouds

- Other ways to do it???

# Understanding Relevance - Macro-level – Dozens or Hundreds

- Visualize Keywords and Doc's
- Show relation of each Doc to Keywords
- "Similar" Doc's cluster together

- VIBE
- SQUID
- JIGSAW
- Google
- Veerasamy's work
- TileBars (dozen or less)

# A VIBE Visualization

laser

plasma

fusion

# SQWID: Search Query Weighted Info Display (VIBE-like)



- Keywords "pull" Doc's
  - (University, Visualization, Tools)
- Doc's can go outside convex hull of keywords (unlike some other approaches)

McCrickard and Kehoe, Visualizing Search Results using SQWID, Poster paper in Proceedings of the 6th World Wide Web Conference (WWW6), Santa Clara CA, April 1997

# Jigsaw

Related docs clustered together

# Problems with Google

- Query results given as text
- Difficult to understand:
  - Relationships between documents with the same relevance value
  - Relative relevance/irrelevance to the query
    Only by position in list of search results
  - Relationships of matching documents to components of the query statement
  - Relationships between results of multiple queries

information visualization text

About 8,290,000 results (0.48 seconds)

**Information Visualization for Text Analysis (Ch 11) | Search ...**
searchuserinterfaces.com/book/sui_ch11_**text**_analysis_**visualization**.html ▾
Full **text** content of the book Search User Interfaces, written by Marti Hearst and published by Cambridge University Press, 2009. Chapter 11: **Information** ...

**See Text in Whole New Way: Text Visualization Tools | ETC ...**
https://blogs.princeton.edu/.../see-**text**-in-whole-new... ▾   Princeton University ▾
Aug 16, 2012 - **Text visualization** adds another dimension to data mining a **text**. You can see in a simple and fast way how many words make up a **text**, what ...
You visited this page on 10/5/15.

**Images for information visualization text**                    Report images



**More images for information visualization text**

**Marti A. Hearst: Research: Information Visualization**
people.ischool.... ▾   University of California, Berkeley School of Information ▾
**Text** Visualization. Tag Clouds: Data Analysis Tool or Social Signaller?, Hearst and Rosner, HICSS 2008, Social Spaces minitrack. pdf; **Information Visualization** ...

**Information visualization - Wikipedia, the free encyclopedia**
https://en.wikipedia.org/wiki/**Information_visualization** ▾   Wikipedia ▾
The abstract data include both numerical and non-numerical data, such as **text** and geographic information. However, **information visualization** differs from ...

[PDF] **Visualization of Text Streams: A Survey**
www.fer.unizg.h... ▾   Faculty of Electrical Engineering and Computing, Uni... ▾
Keywords—**Information Visualization**, Visual Analytics, Topic. Detection and ... **text** visualization is a subfield of **information visualization**. Visual analytics is a ...

CS 4460

# Adding InfoViz to Google



(a)

(b) HotMap

(c) Concept Highlighter

Search terms

Concepts related to search terms

Items in window

HotMap and Concept Highlighter tested somewhat better. See paper for details.

Hoeber & Yang, Comparative Study of Web Search Interfaces, 2006 Conference on Web Intelligence (ACM Digital Library)

# Adding InfoViz to Google

- Another idea (due to JDF)

- Pop-up (tool-tip) word clouds over each retrieved document ☺

# Understanding Relevance - Veerasamy



Veerasamy & Belkin, *Evaluation of a tool for visualization of information retrieval results, ACM Conference on Research and Development in Information Retrieval*, Pro. 19th annual international ACM SIGIR conference (ACM Digital Library)

# Understanding Relevance



- Keyword bars show relevance for document
- "Total sum" gives overall relevance
- Experimental comparison against no InfoVis
  - No difference :-(

# Pros and Cons??

(b) HotMap

# Understanding Relevance - TileBars

- Goal
  - Minimize time and effort for deciding which documents to view in detail

- Idea
  - Show the role of the query terms in the retrieved documents, making use of document structure

Hearst
CHI '95

# TileBars

- Graphical representation of term distribution and overlap

- Simultaneously indicate:
  - Relative document length
  - Frequency of term sets in document
  - Distribution of term sets with respect to the document and each other

# TileBars Interface

Search terms

Presentation



**TileBars: Term Distribution in Information Access**

Min Overlap Span

Search   Redisplay   Reset   Exit    0   2   4   6   8   10

Min Hits                          Min Distribution (%)

Term Set 1: network        0  2  4  6  8  10    0  10  20  30  40  50

Term Set 2: image          0  2  4  6  8  10    0  10  20  30  40  50

Term Set 3:                0  2  4  6  8  10    0  10  20  30  40  50

**Documents Within Contraints**

100   "Hot technologies for 1
101   "Information age  the Smithsonians' LAN
102   "Hot T-1 stuff. (customer premises equ
103   "Comdex Fall. (1989)"
104   "MAN about town: taking the local out of l
105   "Ethernet products: you can get there from
106   "HDTV and
107   "Backing up. (tape back-up strategies)"
108   "CPC '90: gatheri
109   "DEC imaging workstations to challenge PC
110   "Paradox 3.0. (Software Review) ( one of s
111   "Xerox goes wild. (a new version of the

CS 4460

# Additional Terms => More Rows



Three search terms

Three rows per tile bar

# TileBars Interface

Color coding shows degree of relevance

Notice how doc length is treated

# TileBar Issues

- Other Issues?
  - How about scaling up to more documents and search terms
- But what does TileBars give us that Veerasamy method does not? Pros/cons?

# End of Understanding Relevance

- VIBE
- SQUID
- JIGSAW
- Google (and extensions)
- Veerasamy's work
- TileBars (dozen or less)

- Pros/cons?
- How does each scale with
  - Number of documents?
  - Number of keywords?

# Making Sense of Smaller Collections

- Getting an overall sense
  - Word/Tag clouds
- Comparing multiple docs
  - NYT State of Union example
  - Washington Post State of Union example
  - New York Times Political Convention speeches
  - Plagiarism

# Tag Clouds and Word Clouds

- **Tag** Clouds represent explicit meta-data about a document or web site or picture.
  - Typically user-assigned – "crowd-sourced keywords"
- **Word** Clouds, in contrast, are generated from the words of a document or document collection or web site.
  - In some sense they are automatically generated meta-data.
  - We could call them implicit meta-data
- Ways to use are same, but we will always refer to Word Clouds, not Tag Clouds

# Alpha order / prominent in center / etc

# Same / Different Orientation



CS 4460

# Lots of Other Design Options

# Bubble Chart

## Pros/ Cons?

Bubble Charts (implemented by Wattenberg)

# Other Uses for Word Bubble Chart?

- How could we use this to find differences between
  - Authors
  - Centuries
  - Versions
- Group time ☺

# Word Pairs more Useful



Made with Many Eyes
Tag Cloud Viz
(not same as Word Cloud Viz)

# Meaningful Associations Confused

Find the country names in this cloud

FAST!!

# Alternative: "Semantic" Layout

Tags are grouped based on clustering and co-occurrence analysis – words that co-occur close to one another in the text are placed together in the cloud

ajax apple **art** article audio **blog** blogging **blogs** books business code comics community computer cool **css** culture daily del.icio.us delicious **design development** diy firefox flash flickr free freeware **fun** funny games geek **google** graphics gtd hacks hardware history **howto** html humor images **internet** **java javascript** language lifehacks **linux mac** maps media movies mp3 **music news** opensource OSX photo **photography** photos php politics productivity **programming** python rails **reference** research rss ruby science **search** security shopping social **software** tech technology tips tool **tools** toread travel **tutorial** tutorials usability video **web** web2.0 **webdesign** webdev wiki windows writing xml

lisp perl python ruby rails
database wordpress fonts wiki gtd
books writing language math **science** philosophy religion history politics
media **news blog blogs** internet technology business web2.0 rss search google
firefox accessibility usability php xml ajax **javascript** html **css** webdesign
**design web reference** howto tutorial **java programming** development **tools software** opensource free
windows **linux** unix security networking hardware apple **mac** osx
game **games fun** funny **humor art photography flash** animation comics
cinema film movies movie **video** tv
audio **music** mp3 ipod radio podcast podcasting
mobile treo psp xbox fashion shopping
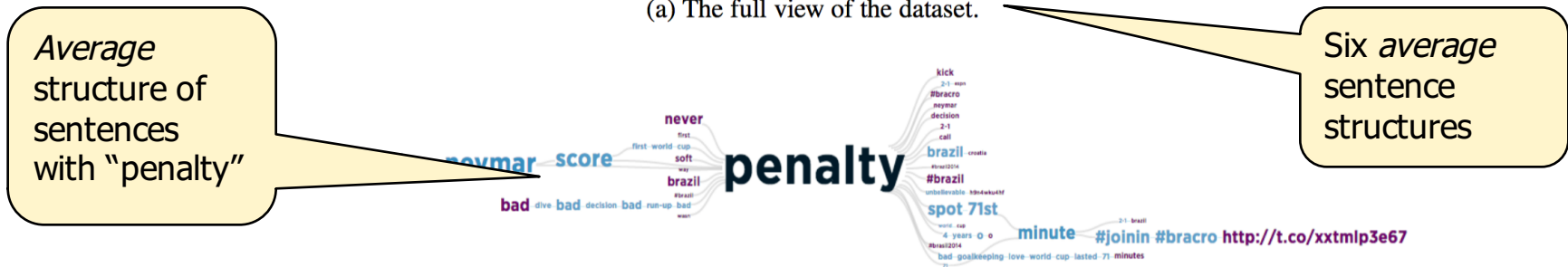travel food health marketing advertising

Hassan-Monteroa & Herrero-Solana, *Improving Tag-Clouds as Visual Information Retrieval Interfaces*, InSciT2006

CS 4460

# Semantics-preserving Word Clouds

1. Determine co-occurrence word count within sentence.

2. Separate words into clusters based on high co-occurrence count.



Semantic-Preserving Word Clouds by Seam Carving, Wu et. al, EuroVis 20111

# SentenTree



(a) The full view of the dataset.

*Average* structure of sentences with "penalty"

Six *average* sentence structures

(b) A zoomed-in view focused on *penalty* after the viewer has clicked on that word.

Fig. 8: A SentenTree visualization of tweets commenting on the third goal of the opening game of the World Cup. This dataset contains 132,599 tweets (75,930 unique tweets).

## Think of this as a really smart word cloud!

CS 4460

Hu, Wongsuphasawat and Stasko, Visualizing Social Media Content with SentenTree, IEEE Transactions on Visualization and Computer Graphics, January 2017.

# Document Comparison: NYTimes



Putting it together:  Ben Werschkul of the *New York Times (DEAD LINK ☺)*
http://www.nytimes.com/ref/washington/20070123_STATEOFUNION.html?initialWord=iraq

# Word Counts – Republicans, Democrats



## At the National Conventions, the Words They Used

A comparison of how often speakers at the two presidential nominating conventions used different words and phrases, based on an analysis of transcripts from the Federal News Service.

Add word or phrase  +

← Words favored by Democrats    Number of mentions per 25,000 spoken words by ● Democrats and ● Republicans    Words favored by Republicans →

**AUTO** Democrats credited President Obama with the recovery of the auto industry after the 2009 bailout, while Republicans left the topic unmentioned.

**WOMEN** Democrats used the word much more frequently, primarily in reference to women's health and equal pay.

**BUSINESS** Republicans were more likely to talk about businesses, emphasizing Mr. Romney's private-sector experience and plans to improve the economy.

**UNEMPLOYMENT** Many Republican speakers brought up the still-high unemployment rate and the number of Americans who remain jobless, while Democrats largely avoided the topic.

Democrats mentioned Romney
92 times per 25,000 words

Republicans mentioned Romney
109 times per 25,000 words

# Document Comparison: Washington Post
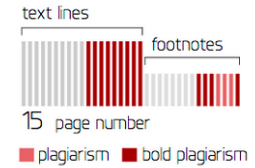
## State of the Union Speeches, 1900-2016



What year? ☹

http://www.washingtonpost.com/graphics/politics/2015-sotu/language/

# Plagiarism Example

Plagiarism in the PhD thesis
of Karl-Theodor Freiherr zu Guttenberg

# Visualization of Alleged Plagiarism

Types of plagiarism at bottom (c), a list of difflines (b) (glyph-based visualization of the finding spots (d)) in the center, overview on left (a). Copy-and-pasted passages marked in red (e). A finding spot can be opened for side-by-side comparison of suspicious and original text fragments. Overview reveals the distribution of finding spots across the document (g) and relationship to sources (h). Overview supports brushing and selection to define subset of finding spots displayed in diffline view.

# Digging into Document Details

- A **concordance** is an alphabetical list of the principal words used in a book or body of work, with their immediate contexts (sometimes called KWIC or Key Word in Context)

- A **frequency list** is a sorted list of words together with their frequencies

- Word Trees

- Grammatical Structures

# Concordance & Frequency List Together



From www.concordancesoftware.co.uk

# Concordance: 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* `08

# Word Trees (Words in context)

*Shift-click to make that word the root.*



someday go to college. And 17 years later I did go to college. But I naively chose

I — would

walk the 7 miles across town every Sunday night to get one good meal a week at

have never dropped in on this calligraphy class, and personal computers might

the final adoption papers. She relented a few months later when my parents promised that I would someday go to college.

And 17 years later I did go to college. But I naively chose a college that was almost as expensive as Stanford, and all of my working-class parents' savings were being spent on my college tuition. After six months, I could see the value in it. I had no idea what I wanted to do with my life and no idea how college was going to help me figure it out. And here I was spending all of the money my parents had saved their entire life. So I decided to drop out and trust that it would all work out OK. It was pretty scary at the time, but looking back it was one of the best decisions I ever made. The minute I dropped out I could stop taking the required classes that didn't interest me, and begin dropping in on the ones that looked interesting.

It wasn't all romantic. I didn't have a dorm room, so I slept on

https://www.jasondavies.com/wordtree/?source=steve-jobs-commencement.txt&prefix=I&reverse=0

# Multiple Words

- How about sequences or pairs of words?
- Are there good ways to present them?

# Phrase Nets

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

van Ham et al
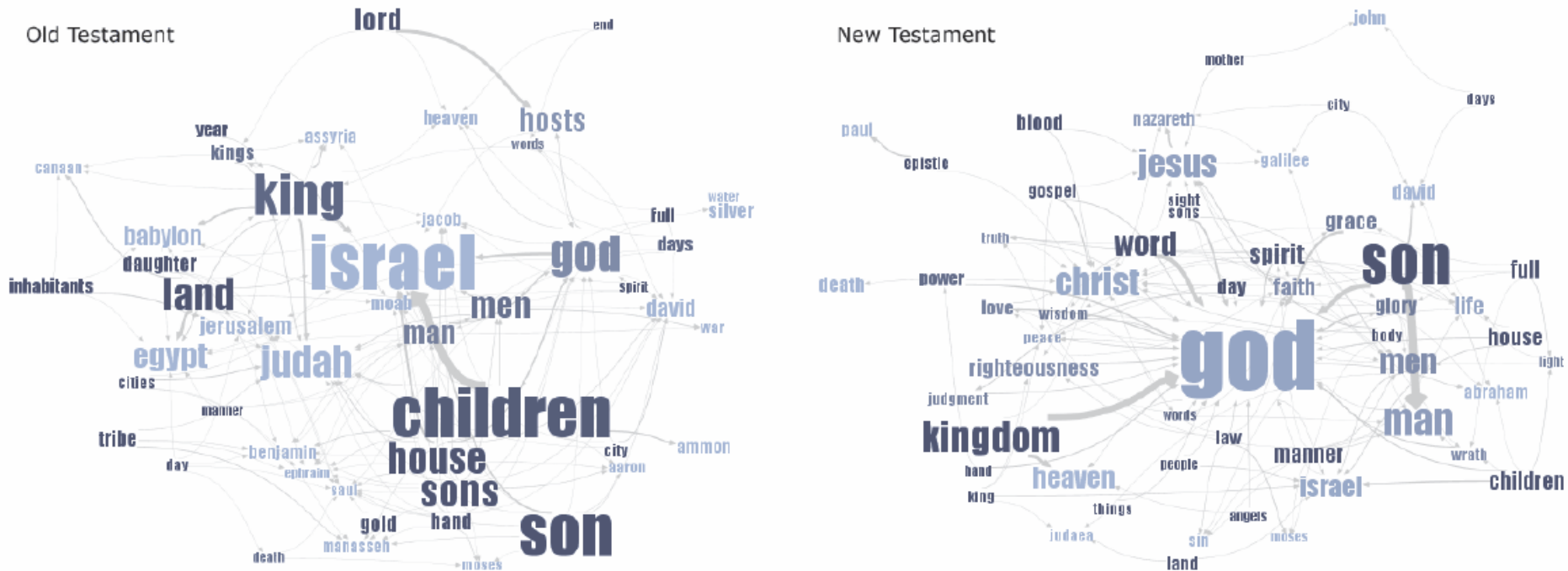*TVCG* '09

# Phrase Net 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.

# Phrase Net 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.

# Structured Document

- Lots of meta-data

- "Sig Dig"example
- FacetMaps
- PaperLens
- ResultMaps

# News Story Trends



2016 Summer Recap: A visualization of the noteworthy news from Memorial Day to Labor Day.

# Structured Info Spaces: FacetMaps

Facets: aka attributes

Displays attributes of a "case"

Drill down to one or a few cases (papers in this example)

http://research.microsoft.com/en-us/um/redmond/groups/cue/facetlens/FacetLens-Video-CHI.wmv

# Structured Info Spaces: PaperLens



a) Popularity of topic

b) Selected authors

c) Author list

d) Degrees of separation of links

e) Paper list

f) Year-by-year top ten cited papers/ authors – can be sorted by topic

Video on next slide

# PaperLens -  Conclusion

- Shows all the data at once but not all the relationships
- Leverages tightly coupled views to show correlations
- Effective in answering questions regarding:
  - Patterns such as frequency of authors and papers cited
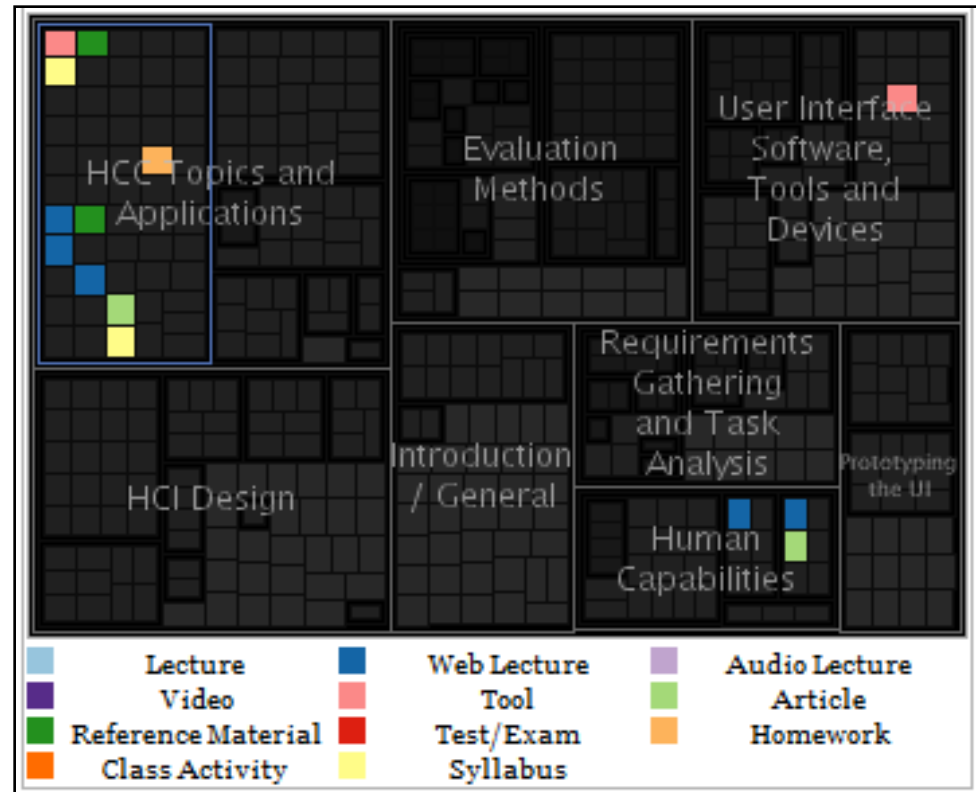  - Themes are clearly visible
  - Trends such as number of papers published in a topic area over time
  - Correlations between authors, topics and citations
- User study
  - Average task performed in less than 20 seconds
  - Participants provided correct answers to tasks 97% of the time
- Scaling challenge
  - Fish-eye technique not as effective with larger number of papers
  - Over-lapping highlights
  - No-more than four authors can be color-coded in the overview

- Lee, Czerwinski, Robertson, Bederson, Understanding Research Trends in Conferences using PaperLens, Extended Abstracts of CHI 2005, pp. 1969-1972.

PaperLens Slides by Jennifer Stoll, GT

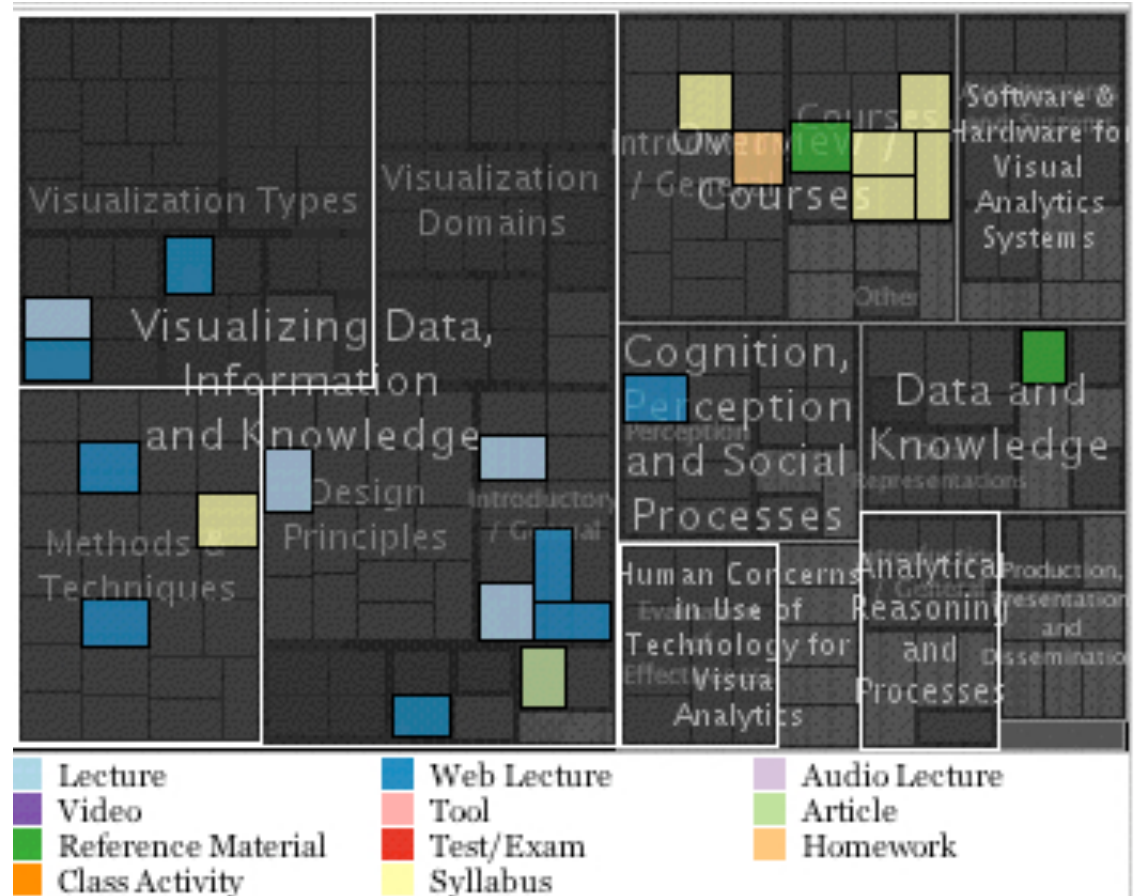CS 4460

# Structured Info Spaces: ResultMaps

- Problem - understand a HIERARCHICAL information space and how retrieved results fit into that space
  - Contrast with PaperLens
- Solution - ResultMaps
  - Based on TreeMaps
  - Highlight documents retrived via text query
  - Linkage from highlight to document in retrieval list

- BTW – do a quick critique of this InfoVis

# ResultMap

- Experimental evaluation
  - Compare with Google-style results list
  - Not much help ☹
  - Some evidence that they are subjectively preferred and help understand overall document collection structure

- (Is this InfoVis any better than previous?)

# How to Think About all of this?

- Remember this outline?
- Macro-level – searching larger document collections
  - Unstructured – no meta-data
  - Structured – explicit meta-data
- Micro-level
  - Inter-document methods for smaller document collections
    How do retrieved documents relate to a query?
    How do retrieved documents relate to one another?
  - Intra-document methods
    Word usage, grammatical style, …

- With the caveat that some methods can be used in multiple ways

# Text and Documents Takeaways

- It's a hugh space – need to understand
  - From searching everything (WWW) to analyzing a single document
- Many opportunities for creativity
- What are user activities with Text and Documents? How can InfoVis support those activities?
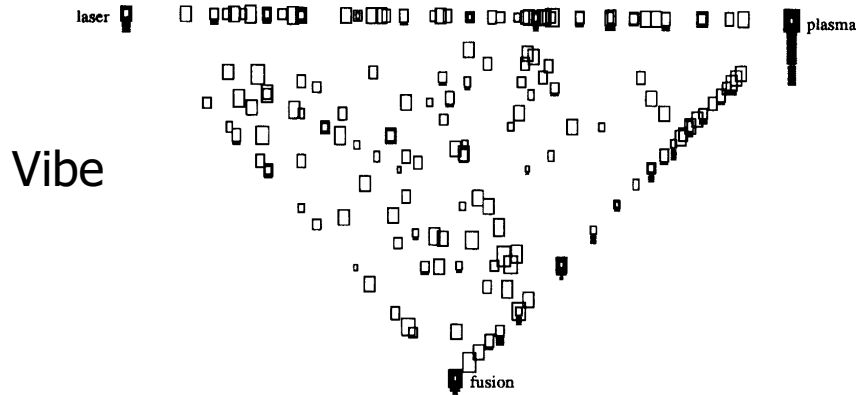
# Text and Documents Takeaways

- Which methods scale from one or a few documents to thousands of docs on up to the WWW? Why? Why not?

- How do we know which methods are good and which are not so good?

- Are there places where using InfoVis does not make sense?  What are they?
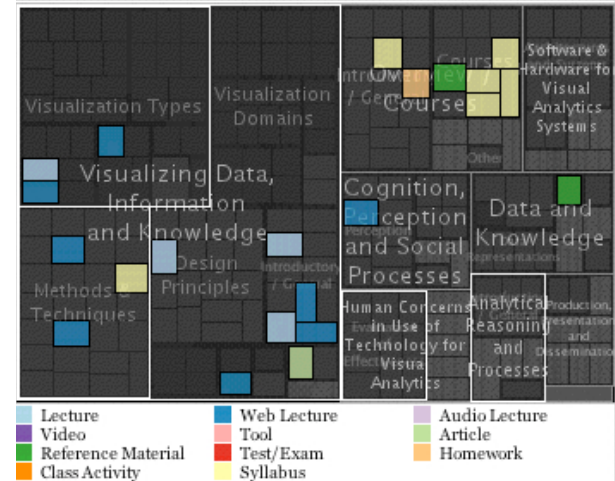
# Can Combine, Mix and Match

- From big picture overview of many docs to query-related views to detailed views of a few docs to within a single doc
- Add interactions to info presentations – the usual suspects
  - DoD
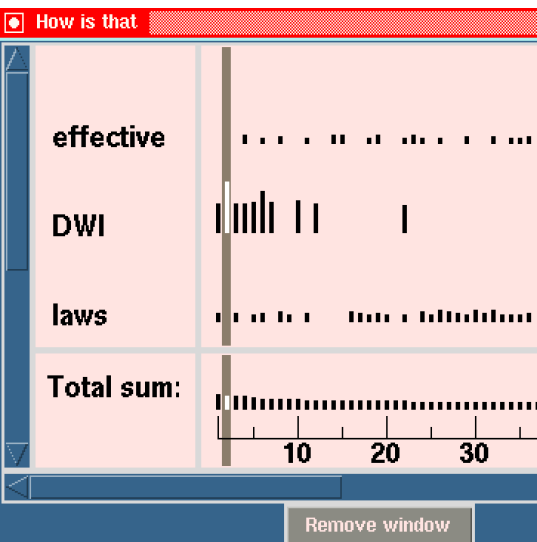  - Dynamic queries/filters
  - Animation
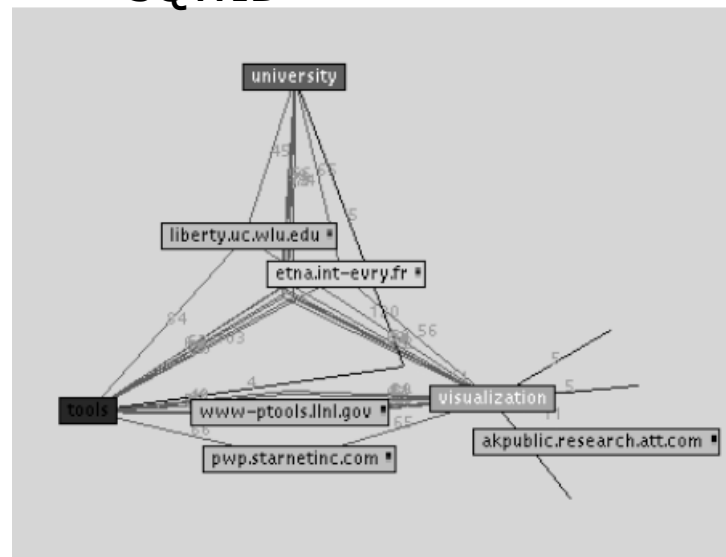  - Brush/Link

# Compare and Contrast

Vibe

ResultMap

Verasamy

SQWID

TileBars