

Text and Document Visualization 2



CS 7450 - Information Visualization
March 17, 2011
John Stasko

Example Tasks & Goals

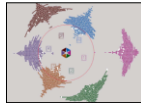


- Which documents contain text on topic XYZ?
- Which documents are of interest to me?
- Are there other documents that are similar to this one (so they are worthwhile)?
- How are different words used in a document or a document collection?
- What are the main themes and ideas in a document or a collection?
- 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.
- Find connections between documents.

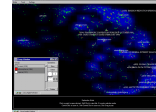
This Week's Agenda



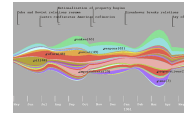
Visualization for IR
Helping search



Visualizing text
Showing words,
phrases, and
sentences



Visualizing document sets
Words & sentences
Analysis metrics
Concepts & themes



Last Time

Spring 2011

CS 7450

3

Related Topic - Sensemaking

Recall



- 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

Spring 2011

CS 7450

4

Today's Agenda



- Move to collections of documents
 - Still do words, phrases, sentences
 - Add
 - More context of documents
 - Document analysis metrics
 - Document meta-data
 - Document entities
 - Connections between documents
 - Documents concepts and themes

Spring 2011

CS 7450

5

Various Document Metrics



- Goals?
- 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

Keim & Oelke
VAST '07

Spring 2011

CS 7450

6

Vis

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

Do partial overlap in blocks for a smoother appearance

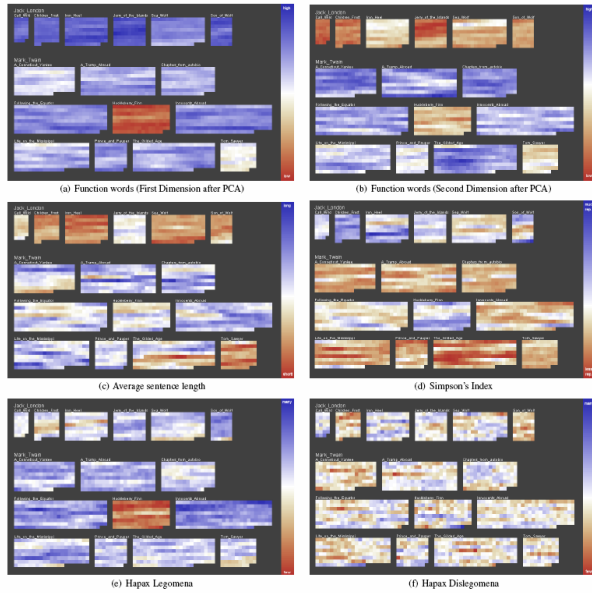


Figure 2: Fingerprints of books of Mark Twain and Jack London. Different measures for authorship attribution are tested. If a measure is able to discriminate between the two authors, the visualizations of the books that are written by the same author will equal each other more than the visualizations 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.

The Bible

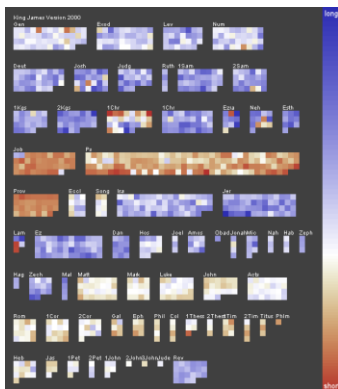


Figure 4: Visual Fingerprint of the Bible. Each pixel represents one chapter of the bible and color is mapped to the average verse length. Interesting characteristics such as the generally shorter verses of the poetry books, the inhomogeneity of the 1. Book of Chronicles or the difference between the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

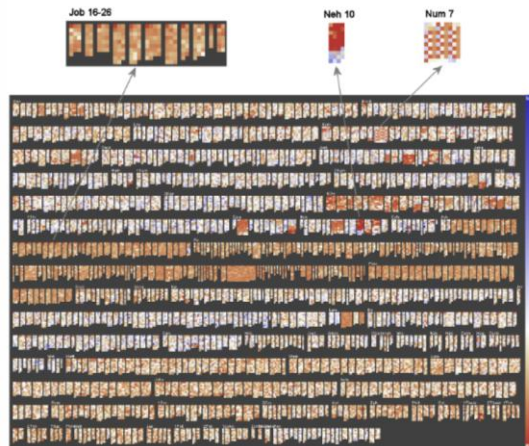


Figure 5: Visual Fingerprint of the Bible. More detailed view on the bible in which each pixel represents a single verse and verses are grouped to chapters. Color is again mapped to verse length. The detailed view reveals some interesting patterns that are camouflaged in the averaged version of fig.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

Follow-On Work



- Focus on readability metrics of documents
- Multiple measures of readability
 - Provide quantitative measures
- Features used:
 - Word length
 - Vocabulary complexity
 - Nominal forms
 - Sentence length
 - Sentence structure complexity

Oelke & Keim
VAST '10

Visualization & Metrics



		Voc. Difficulty	Word Length	Nominal Forms	Sent. Length	Compl. Sent. Struc.
(a)	The intention of TileBars [9] is to provide a compact but yet meaningful representation of Information Retrieval results, whereas the FeatureLens technique, presented in [5], was designed to explore interesting text patterns which are suggested by the system, find meaningful co-occurrences of them, and identify their temporal evolution.	Blue	Red	Blue	Red	Red
(b)	This includes aspects like ensuring contextual coherency, avoiding unknown vocabulary and difficult grammatical structures.	Red	Red	White	Blue	Blue

Figure 5: Two example sentences whose overall readability score is about the same. The detail view reveals the different reasons why the sentences are difficult to read.

Uses heatmap style vis (blue-readable, red-unreadable)

Interface

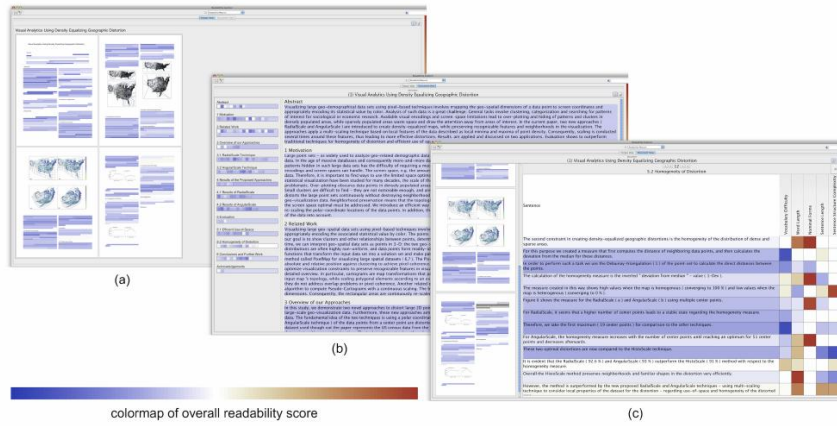


Figure 3: Screenshot of the VisRA tool on 3 different aggregation levels. (a) Corpus View (b) Block View (c) Detail View. To display single features, the colormap is generated as described in section 3.4 and figure 2.

Their Paper (Before & After)

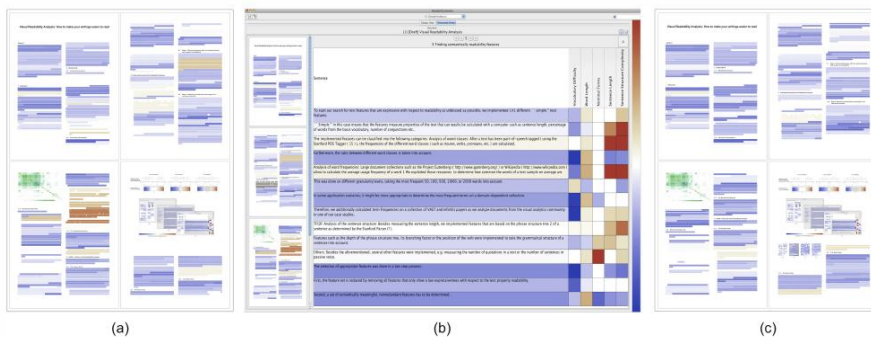


Figure 6: Revision of our own paper. (a) The first four pages of the paper as structure thumbnails before the revision. (b) Detail view for one of the sections. (c) Structure thumbnails of the same pages after the revision.

Comment from the Talk



- In academic papers, you want your abstract to be really readable
- Would be cool to compare rejected papers to accepted papers

Overviews of Documents

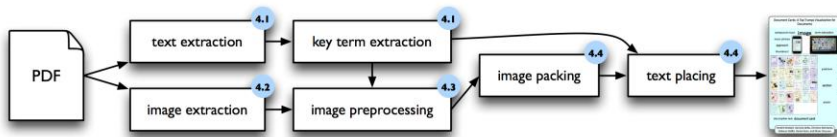


- Can we provide a quick browsing, overview UI, maybe especially useful for small screens?

Document Cards



- Compact visual representation of a document
- Show key terms and important images



Strobelt et al
TVCG (InfoVis) '09

Spring 2011

CS 7450

15

Representation



Layout algorithm searches for empty space rectangles to put things

Spring 2011

CS 7450

16

Interaction



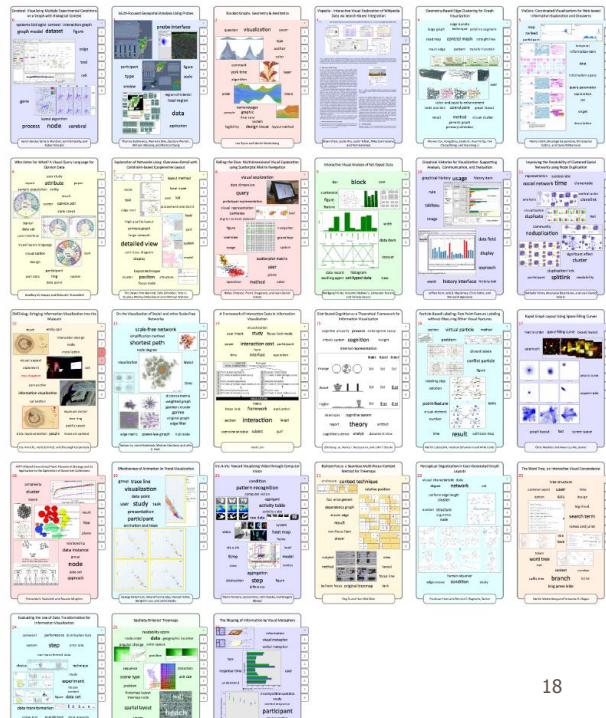
- Hover over non-image space shows abstract in tooltip
- Hover over image and see caption as tooltip
- Click on page number to get full page
- Click on image goes to page containing it
- Clicking on a term highlights it in overview and all tooltips

Spring 2011

CS 7450

17

InfoVis '08
Proceedings

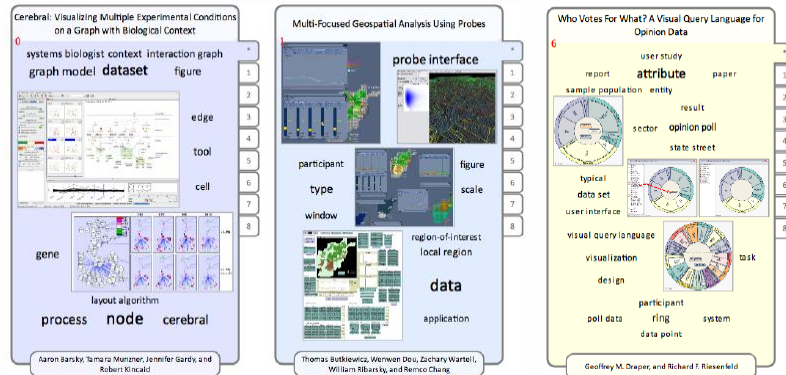


Spring 2011

18



Zooming In



Spring 2011

CS 7450

19

PaperLens



- Focus on academic papers
- Visualize doc metadata such as author, keywords, date, ...
- Multiple tightly-coupled views
- Analytics questions
- Effective in answering questions regarding:
 - Patterns such as frequency of authors and papers cited
 - Themes
 - Trends such as number of papers published in a topic area over time
 - Correlations between authors, topics and citations

Lee et al
CHI '05 Short

Spring 2011

CS 7450

20

PaperLens

Video



- 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

Spring 2011

CS 7450

21

NetLens

Kang et al
Information Visualization '07



Figure 1 NetLens has two symmetric windows. The left is for Content (papers) and the right for Actors (authors). Each side is further divided into panels; overview at the top, filters on the right, and lists at the bottom. Here, the Content side has two lists to reflect papers and their citations or references, and the lists on the Actor side show authors and their co-authors, respectively. The paper overview panel shows the distribution of papers (in logarithmic scale) over time, grouped by topics. Users can see which topics have their number of papers increase or decrease over 22 years. On the right side, the overview of the authors shows the distribution of countries of origin in logarithmic scale.

Spring 2011

CS 7450

22

More Document Info



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

Spring 2011

CS 7450

23

Jigsaw



- 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

Spring 2011

CS 7450

24

Document View



The screenshot shows a 'Document View' window. At the top, a word cloud displays terms like 'analysis', 'information', 'interaction', and 'visualization'. Below it is a 'Doc List' with document IDs. A selected document is shown with a summary and a snippet of text. Annotations with arrows point to the word cloud (labeled 'Wordcloud overview'), the document list (labeled 'Doc List'), the document summary (labeled 'Document summary'), and the text snippet (labeled 'Selected document's text with entities identified').

Spring 2011

CS 7450

25

List View

Entities listed by type



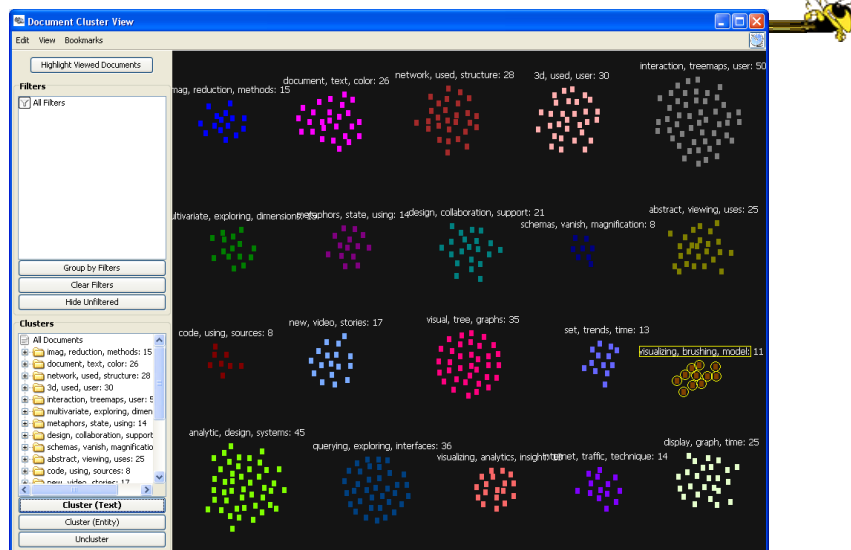
The screenshot shows a 'List View' window. On the left is a hierarchical tree of concepts such as 'interaction', 'evaluation', 'insight', 'visual analytics', 'case study', 'cognition', 'color', 'navigation', 'animation', 'categorical', 'document', 'dynamic query', 'filter', 'focus+context', 'hierarchy', 'intelligence analysis', 'metrics', 'perception', 'social', 'software visualization', 'text', 'theory', 'time series', 'treemap', 'graph', and 'high-dimensional data'. Lines connect these concepts to three central lists: 'author', 'year', and 'conference'. The 'author' list includes names like 'Spence, B.', 'Springue, D.W.', 'Stasko, J.', 'Steed, C.A.', 'Stien, C.', 'Stodolinger, K.', 'Stothel, A.', 'Stothel, C.', 'Storey, M.-A.D.', 'Shneider, T.', 'Strayer, D.', 'Stroholtz, H.', 'Stroholtz, P.J.', 'Studley, P.', 'Stulan, F.', 'Sturtebeck, E.P.', 'Shurtz, D.', 'Su, H.', 'Sudianto, A.', 'Suh, B.', 'Sullivan, T.', 'Suma, E.', 'Summers, K.L.', 'Sunder, J.', 'Swain, J.E.', 'Swindells, C.', 'Srinani, N.', 'Takekuma, Y.', 'TSA, A.', 'Tabbot, J.', 'Tan, D.S.', 'Tan, R.', 'Tanasse, T.', 'Tandon, S.', 'Tang, D.', 'Tamm, E.', 'Tatui, A.', and 'Tavanti, M.'. The 'year' list shows years from 1995 to 2009. The 'conference' list shows 'Infovis' and 'VAST'.

Spring 2011

CS 7450

26

Document Cluster View

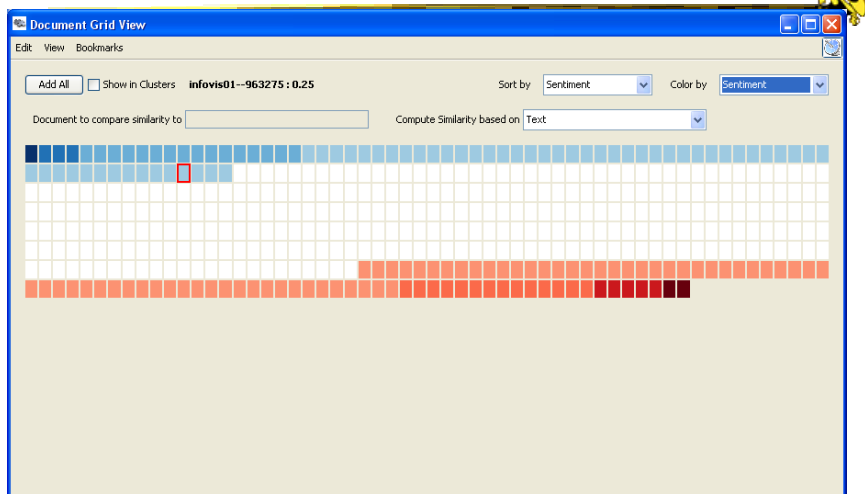


Spring 2011

CS 7450

27

Document Grid View



Here showing sentiment analysis of docs

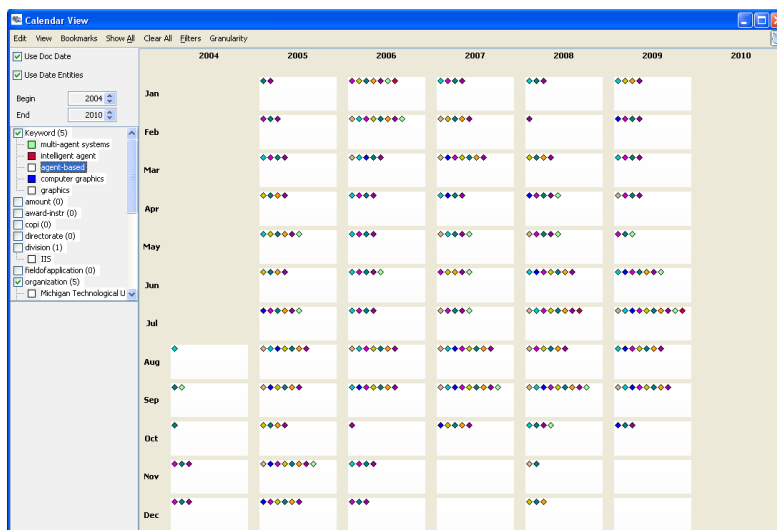
Spring 2011

CS 7450

28

Calendar View

Temporal context
of entities & docs



Spring 2011

CS 7450

Video

29

Jigsaw

- Much more to come on Visual Analytics day...

Spring 2011

CS 7450

30

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

Spring 2011

CS 7450

33

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

Spring 2011

CS 7450

34

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
 - Use a mass-spring layout algorithm to position representations of each document
- Some similarities to how search engines work

Spring 2011

CS 7450

35

Wiggle



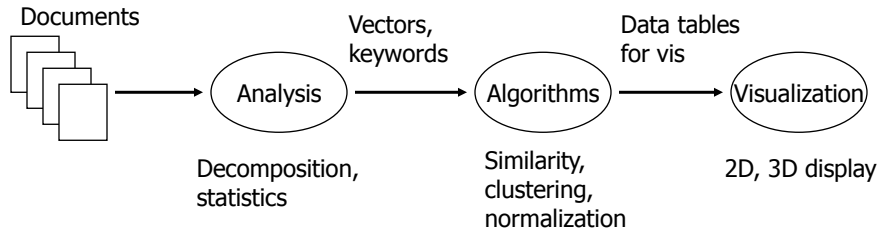
- Not all terms or words are equally useful
- Often apply TFIDF
 - Term frequency, inverse document frequency
- Weight of a word goes up if it appears often in a document, but not often in the collection

Spring 2011

CS 7450

36

Process



Spring 2011

CS 7450

37

Smart System



- Uses vector space model for documents
 - May break document into chapters and sections and deal with those as atoms
- Plot document atoms on circumference of circle
- Draw line between items if their similarity exceeds some threshold value

Salton et al
Science '95

Spring 2011

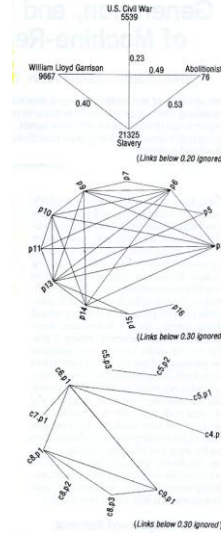
CS 7450

38

Text Relation Maps



- Label on line can indicate similarity value
- Items evenly spaced
- Doesn't give viewer idea of how big each section/document is



Spring 2011

CS 7450

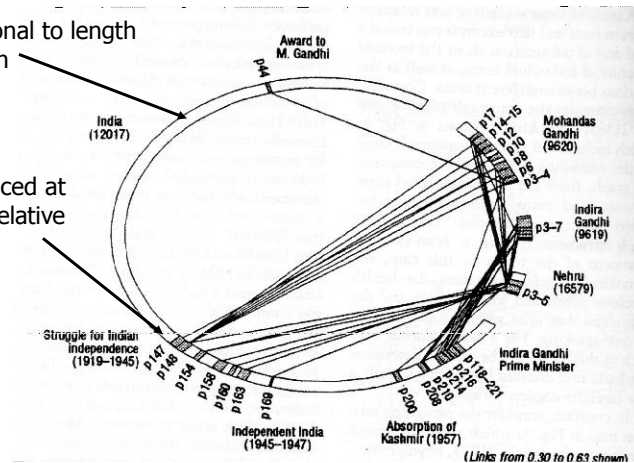
39

Improved Design



Proportional to length of section

Links placed at correct relative position



Spring 2011

CS 7450

40

Text Themes



- Look for sets of regions in a document (or sets of documents) that all have common theme
 - Closely related to each other, but different from rest
- Need to run clustering process

Algorithm

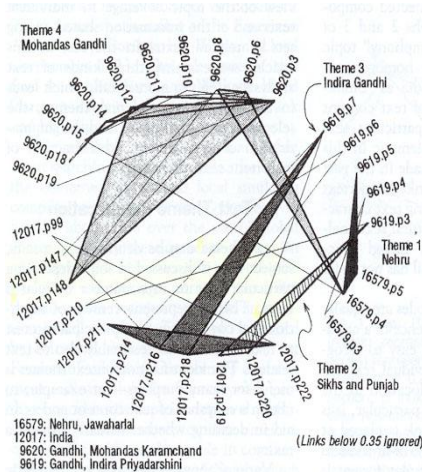


- Recognize triangles in relation maps
 - Three with edges above threshold
- Make a new vector that is centroid of 3
- Triangles merged whenever centroid vectors are sufficiently similar

Text Theme Example



- Triangles shown
- Colored in to help presentation



Spring 2011

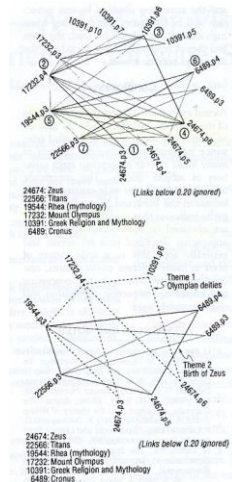
CS 7450

43

Skimming and Summarization



- Can use graph traversal to follow specific themes throughout collection
- Walk along connected edges



Spring 2011

CS 7450

44

VIBE System



- Smaller sets of documents than whole library
- Example: Set of 100 documents retrieved from a web search
- Idea is to understand contents of documents relate to each other

Olsen et al
Info Process & Mgmt '93

Spring 2011

CS 7450

45

Focus



- Points of Interest
 - Terms or keywords that are of interest to user
 - Example: cooking, pies, apples
- Want to visualize a document collection where each document's relation to points of interest is show
- Also visualize how documents are similar or different

Spring 2011

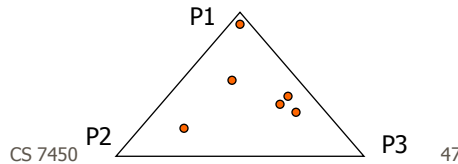
CS 7450

46

Technique



- Represent points of interest as vertices on convex polygon
- Documents are small points inside the polygon
- How close a point is to a vertex represents how strong that term is within the document

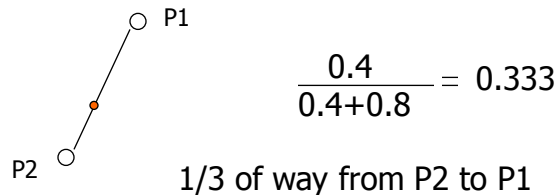


Spring 2011

Algorithm



- Example: 3 POIs
- Document (P1, P2, P3) (0.4, 0.8, 0.2)
- Take first two



Spring 2011

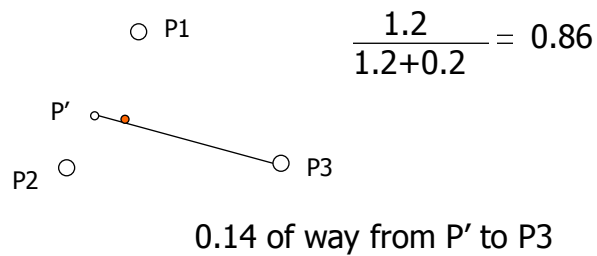
CS 7450

48

Algorithm



- Combine weight of first two 1.2 and make a new point, P'
- Do same thing for third point

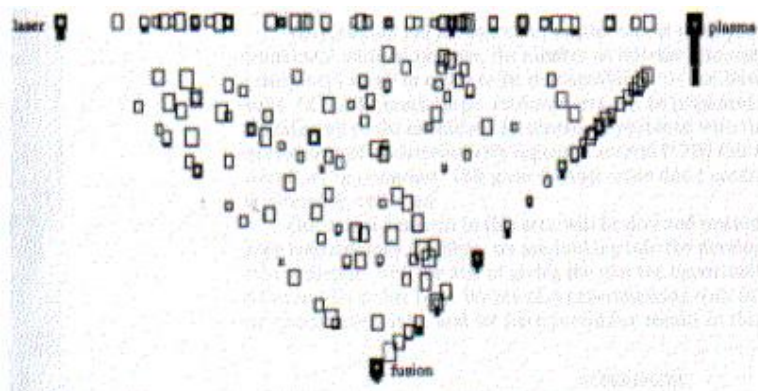


Spring 2011

CS 7450

49

Sample Visualization



Spring 2011

CS 7450

50

VIBE Pro's and Con's



- Effectively communications relationships
- Straightforward methodology and vis are easy to follow
- Can show relatively large collections
- Not showing much about a document
- Single items lose “detail” in the presentation
- Starts to break down with large number of terms

Visualizing Documents



- VIBE presented documents with respect to a finite number of special terms
- How about generalizing this?
 - Show large set of documents
 - Any important terms within the set become key landmarks
 - Not restricted to convex polygon idea

Basic Idea



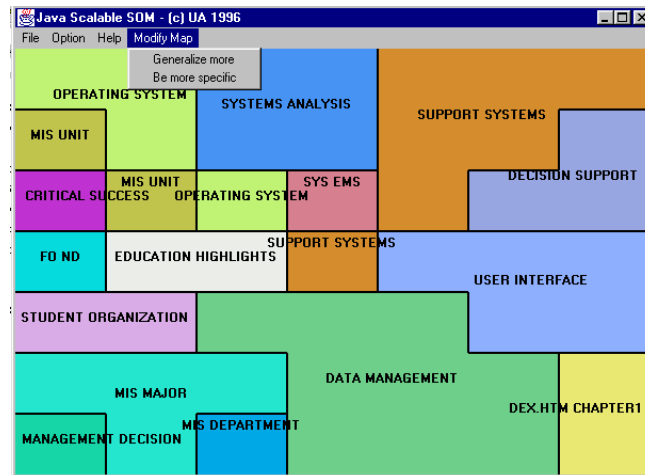
- Break each document into its words
- Two documents are “similar” if they share many words
- Use mass-spring graph-like algorithm for clustering similar documents together and dissimilar documents far apart

Kohonen’s Feature Maps



- AKA Self-Organizing Maps
- Expresses complex, non-linear relationships between high dimensional data items into simple geometric relationships on a 2-d display
- Uses neural network techniques

Map Display of SOM



Spring 2011

CS 7450

55

Map Attributes



- Different, colored areas correspond to different concepts in collection
- Size of area corresponds to its relative importance in set
- Neighboring regions indicate commonalities in concepts
- Dots in regions can represent documents

Spring 2011

CS 7450

56



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

Wise et al
InfoVis '95

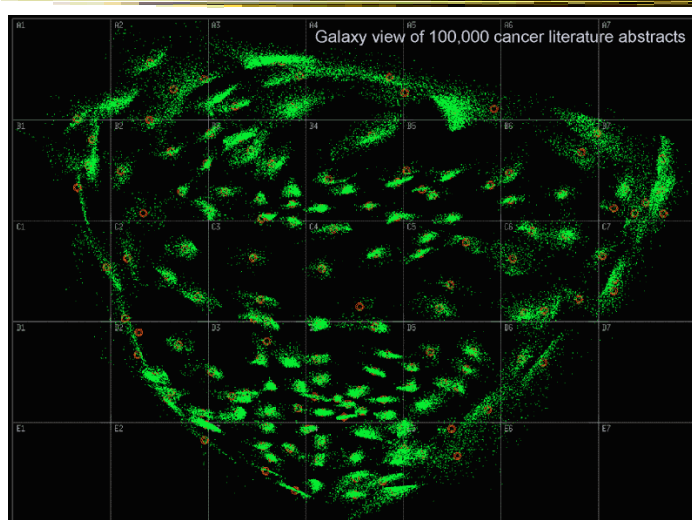
Spring 2011

CS 7450

59

Galaxies

Presentation of documents where similar ones cluster together



Spring 2011

CS 7450

60

Themescapes



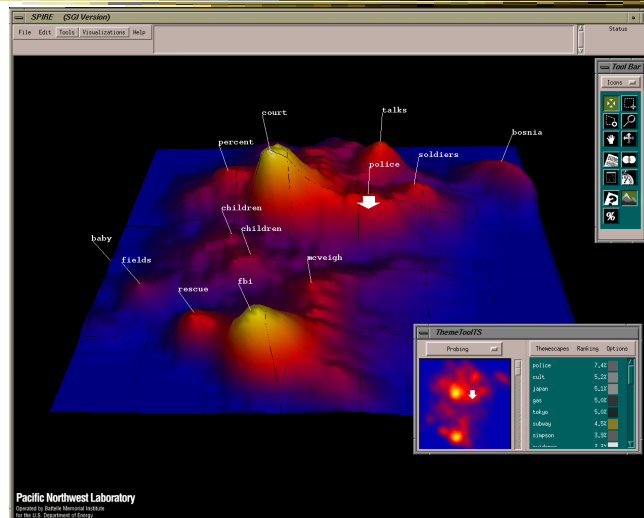
- Self-organizing maps didn't reflect density of regions all that well -- Can we improve?
- Use 3D representation, and have height represent density or number of documents in region

Spring 2011

CS 7450

61

Themescape



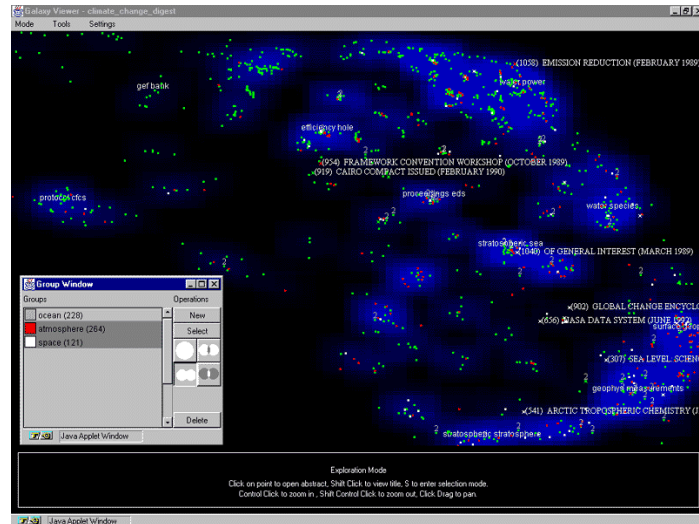
Video

Spring 2011

CS 7450

62

WebTheme



Spring 2011

CS 7450

63

Temporal Issues



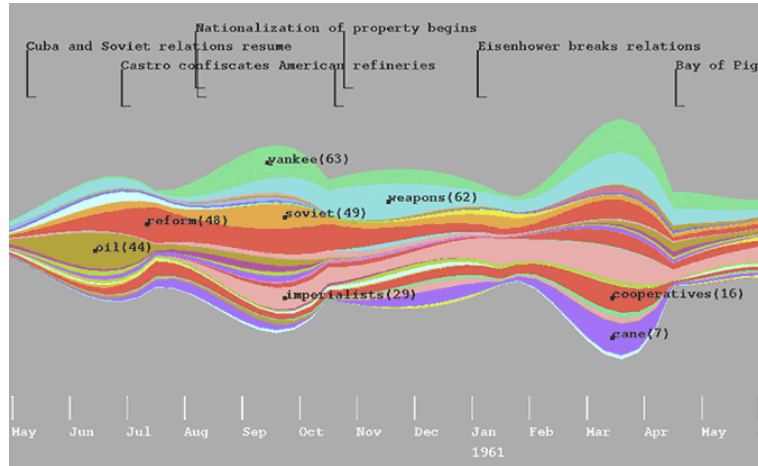
- Semantic map gives no indication of the chronology of documents
- Can we show themes and how they rise or fall over time?

Spring 2011

CS 7450

64

ThemeRiver



Havre, Hetzler, & Nowell
InfoVis '00

Spring 2011

CS 7450

65

Representation



- Time flows from left->right
- Each band/current is a topic or theme
- Width of band is "strength" of that topic in documents at that time

Spring 2011

CS 7450

66

More Information



- What's in the bands?
- Analysts may want to know about what each band is about

Spring 2011

CS 7450

67

TIARA



- Keeps basic ThemeRiver metaphor
- Embed word clouds into bands to tell more about what is in each
- Magnifier lens for getting more details
- Uses Latent Dirichlet Allocation to do text analysis and summarization

Spring 2011

CS 7450

Liu et al
CIKM '09, KDD '10, VAST '10

68

HW 7



- NodeXL reactions

Spring 2011

CS 7450

71

HW 8



- Investigative analysis
- You play the intelligence analyst
- Find the criminal plot embedded across 50 documents
- Paragraph summarizing the threat and a description of what you did

- Due Thursday 31st

Spring 2011

CS 7450

72

Upcoming

- Spring Break
- Visual Analytics 1
 - Reading
Keim et al '08
- Visual Analytics 2
 - Reading
Stasko et al '08

