# VISUALIZATION

**Visualizing Text**

# Overview

- Introduction
- Visualizing search results
- Visualizing documents
- Visualizing document collections

# Overview

- Introduction
- Visualizing search results
- Visualizing documents
- Visualizing document collections

# 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
  - Email
- We often have to make critical decisions based on our understanding of documents
- However, few people have enough time to read everything,
- It is more urgent to scan, understand, operate, and navigate the enormous corpus of documents, and thus to efficiently acquire useful information and knowledge.
- What can information visualization provide to help users in understanding and gathering information from text and document collections?

# Text Visualization – Task and 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.

# Related Topic - IR

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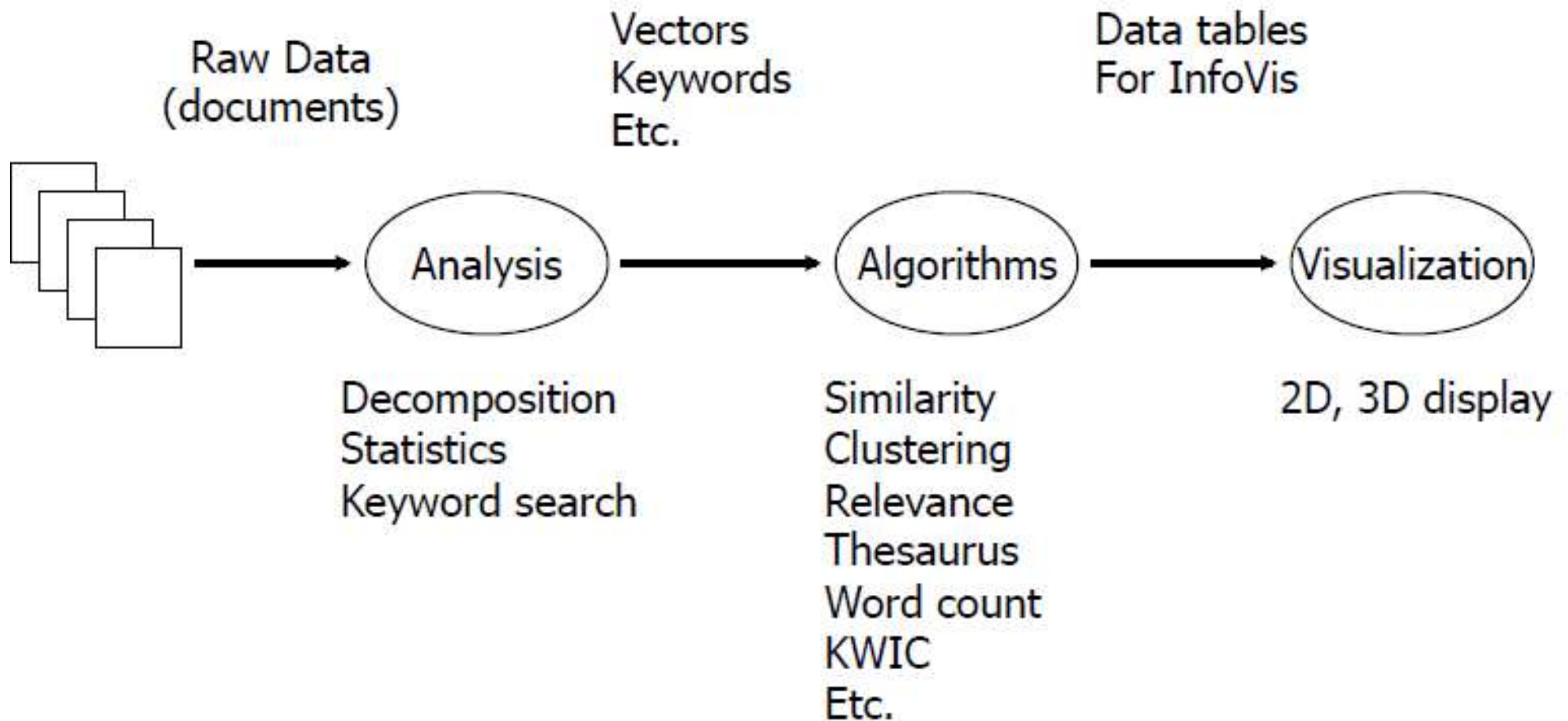  - More of a browsing task than a search one

# Related Topic - Sensemaking

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

# Challenge

- Text is nominal data
  - Does not seem to map to geometric/graphical presentation as easily as ordinal and quantitative data
- Unstructured text does NOT have any explicit meta-data
  - Just that infinitely big collection of nominal data
  - Meta-data is sometimes extracted from raw text (Text Mining)
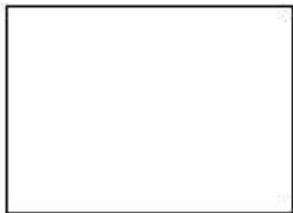- The "Raw data --> Data Table" mapping now becomes more important
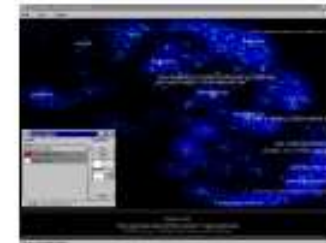
# Text Visualization - Process

Raw Data (documents)

Vectors
Keywords
Etc.

Data tables
For InfoVis

Analysis → Algorithms → Visualization

Decomposition
Statistics
Keyword search

Similarity
Clustering
Relevance
Thesaurus
Word count
KWIC
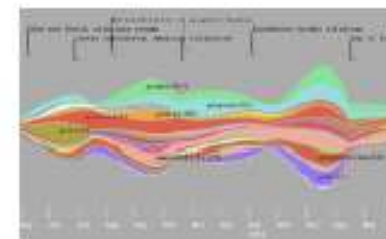Etc.

2D, 3D display

# Types of Text Visualization



**Visualization for IR**
Helping search

**Visualizing text**
Showing words, phrases, and sentences

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

# Overview

- Introduction
- **Visualizing search results**
- Visualizing documents
- Visualizing document collections

# Visualization for IR

- Can InfoVis help IR?
- Assume there is some active search or query
  - Show results visually
  - Show how query terms relate to results
  - …
- Visualizing the results of search operations is another big area in text infovis

# Search - Problems

- Query responses do not include:
  - How strong the match is
  - How frequent each term is
  - How each term is distributed in the document
  - Overlap between terms
  - Length of document
- Document ranking is opaque
- Inability to compare between results

# Tile Bars

- 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
  - Graphical representation of term distribution and overlap using Tile bars
- Simultaneously indicate:
  - Relative document length
  - Frequency of term sets in document
  - Distribution of term sets with respect to the document and each other
- TileBars: Visualization of Term Distribution Information in Full Text Information Access Marti A. Hearst   Proceedings of CHI '95, Denver, CO, May 1995

# Tile Bar



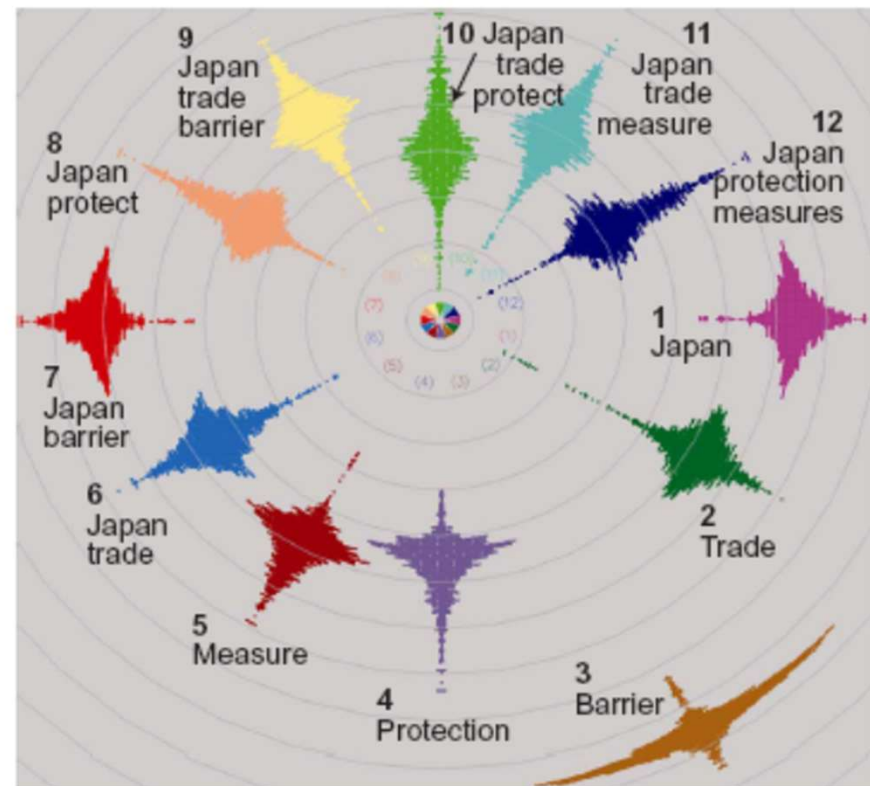Search terms

Presentation

Two search terms

Blocks indicate "chunks" of text, such as paragraphs

Blocks are darkened according to the frequency of the term in the document

# Sparkler

- Visually presenting and exploring the results for different queries on the same topic simultaneously
- Glyphs representing documents arranged along a line
  - Distance from the center indicates degree of relevance to the query.
  - When there were multiple documents with the same relevance score, they were spread out horizontally from the line, forming a visualization of the distribution of relevance scores.
- The Sparkler for the different queries arranged along a circle
  - Each query's visualization assigned a different color.
- Selecting a document in one sparkler caused its position to be highlighted in the other Sparkler visualizations.



Havre S, Hetzler E, Perrine K, Jurrus E, Miller N. Interactive visualization of multiple query results. InfoVis 2001

# Overview

- Introduction
- Visualizing search results
- **Visualizing documents**
- Visualizing document collections
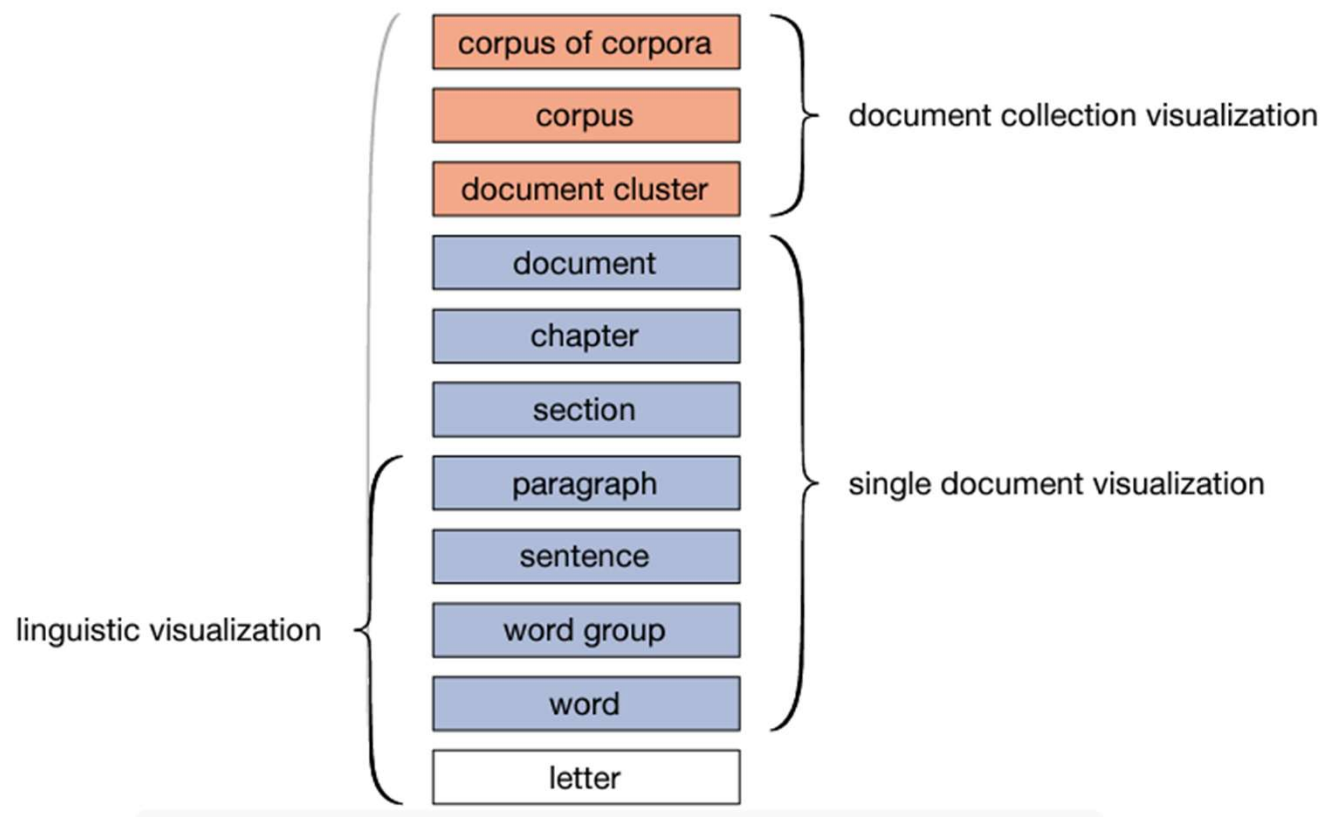
# Visualizing Single Documents

- How do we represent the words, phrases, and sentences in a document or set of documents?
  - Main goal of understanding versus search
  - Visualizing text (features) requires a transformation step: discretization, aggregation, normalization
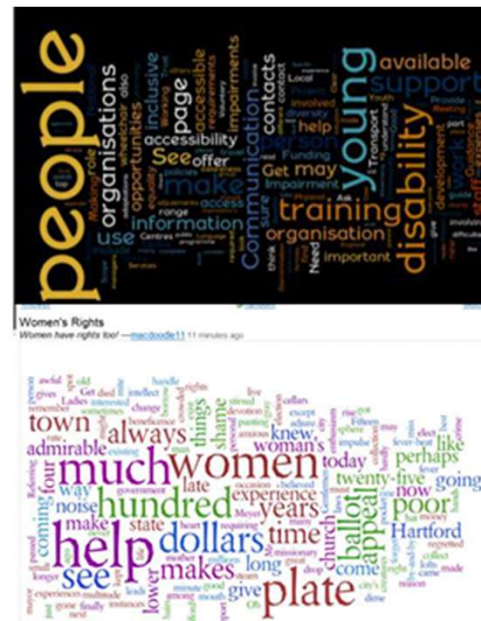


unstructured text → 4 x 't'
3 x 'u'
2 x 'r'
2 x 'e'
...

structured data

# Typical steps of processing to derive Text Features

- Large collections require pre-processing of text to extract information.
- Typical steps are:
  - Cleaning
  - Sentence splitting
  - Changing to lower case
  - Stopword removal
  - Stemming
  - PoS tagging
  - Named entity recognition
  - Deep parsing – trying to "understand" text

# Text Units Hierarchy

# Vocabulary – based Visualization

- Idea is to show word/concept importance through visual means
- Tag Clouds: Provide layouts of raw tokens, colored, and sized by the corresponding word frequency within a single document.

# Tag Cloud/Wordle

- Tightly packed words Idea is to show word/concept importance through visual means

## Tag Clouds

- Provide layouts of raw tokens, colored, and sized by the corresponding word frequency within a single document.

- Words displayed alphabetically or most frequent words in middle

## Wordle

- Wordle adds a layout algorithm to allow users to modify the font, color, or configuration.

# Tag Cloud - Creation

In principle, the font size of a tag in a tag cloud is determined by its incidence. For a word cloud of categories like weblogs, frequency, for example, corresponds to the number of weblog entries that are assigned to a category. For smaller frequencies one can specify font sizes directly, from one to whatever the maximum font size. For larger values, a scaling should be made. In a linear normalization, the weight $t_i$ of a descriptor is mapped to a size scale of 1 through $f$, where $t_{min}$ and $t_{max}$ are specifying the range of available weights.

$$s_i = \left\lceil \frac{f_{max} \cdot (t_i - t_{min})}{t_{max} - t_{min}} \right\rceil \text{ for } t_i > t_{min}; \text{ else } s_i = 1$$

- $s_i$: display fontsize
- $f_{max}$: max. fontsize
- $t_i$: count
- $t_{min}$: min. count
- $t_{max}$: max. count

Implementations of tag clouds also include text parsing and filtering out unhelpful tags such as common words, numbers, and punctuation.

# Word Cloud - Evaluation

- 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
- Word Clouds are really more overview-style visualizations
  - Don't really support queries, searches, drill-down
- However very 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
  - Easy to understand

# Semantic/Context Word Clouds

# Beyond Individual Words

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

# Semantic Graph



- SemanticGraphs (Rusu et al InfoVis 2009) is a visualization based on the semantic representation of a document in the form of a semantic graph.
- The document is summarized with the semantic graph and the list of extracted triplets.
  - It extracts subject–verb–object for each sentence by the parse tree.
  - Then, it links the triplets to their corresponding entity, which needs to resolve pronominal anaphors as well as to attach the associate WordNet synset.

# Word Tree

- The WordTree visualization provides the representation of both word frequency and context.
- Size is used to represent frequency of the term or phrase.
- The root of the tree is a user-selected word or phrase
- The branches represent the contexts in which the word or phrase is used in the document.
- Users can click on a branch, choose a different search term or re-center the tree.
- Clicking on phrase makes it the focus
- Shows context of a word or words
  - Follow word with all the phrases that follow it

- *Wattenberg & Viégas TVCG (InfoVis) '08*

# Example - Word Tree



this nation will rise up and live out the true meaning of its creed: We hold these truths to be self-evident.

on the red hills of Georgia the sons of former slaves and the sons of former slave owners will be able to sit down together at

one day — even the state of Mississippi, a state sweltering with the heat of injustice, sweltering with the heat of oppression, will be transformed into

that

i have a dream — down in Alabama, with its vicious racists, with its governor having his lips dripping with the words of interposition and nutificaion - one day

every valley shall be exalted, and every hill and mountain shall be made low, the rough places will be made plain, and the

my four little children will one day live in a nation where they will not be judged by the color of their skin but by the

cown in Alabama, wih its vicious racists, with its governor having his lips dripping with the words of interposition and nutification - one day

today. i have a dream that one day

every valley shall be exalted, and every hill and mountain shall be made low, the rough places will be made plain, and the

WordTree shows all occurrences of 'I have a dream' in Martin Luther King's historical speech

# Phrase Nets

- A phrase net examine unstructured text documents by displaying a graph whose nodes are words and whose edges indicate that two words are linked by a user-specified relation.
  - These relations may be defined either at the syntactic or lexical level; different relations often produce very different perspectives on the same text.
  - Taken together, these perspectives often provide an illuminating visual overview of the key concepts and relations in a document
- Presents pairs of terms from 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

# Example – Phrase Nets



Matching the pattern "X of Y" to compare old and new testaments. Israel takes the central place in the Old Testament and God acts at he the main pattern in the New Testament.

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

# Visualizing different versions of a document

- History Flow is designed to show changes between multiple document versions on Wikipedia..

- It also reveals some complex patterns of cooperation and confliction, such as vandalism, anonymity versus named authorship, etc.

- Each version of the document is represented by a line with length which is related to the length of its text

- The lines are ordered by date.

- Sections of each line are colored differently according to the different authors.

- Link sections of text that keep the same between adjacent version by drawing shaded connections

Viegas FB, Wattenberg M, Dave K.
Studying cooperation and conflict between authors with history flo visualizations
*Proceedings of the SIGCHI 2004*

# Example – History Flow



- Version history for the Wikipedia entry *Treaty of Trianon*.
- The left side reveals different authors that are colored differently.
- The center part shows the visualization;
- In the right side, a text view closely linked with the visualization shows the detailed content.
- Users can locate on the visualization by moving a set of crosshairs
- The text view shows the corresponding version and position

# Overview

- Introduction
- Visualizing search results
- Visualizing documents
- **Visualizing document collections**

# Comparing Multiple Documents

❖ 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

# Parallel Tag Cloud



The visualization technique combines graphical elements from parallel coordinates and traditional tag clouds to provide rich overviews of a document collection (C. Collins, F. B. Viegas and M. Wattenberg, "Parallel Tag Clouds to explore and analyze faceted text corpora," *2009 IEEE Symposium on Visual Analytics Science and Technology*)

# Themail



A visualization that portrays relationships using the interaction histories preserved in email archives. Using the content of exchanged messages, it shows the words that characterize one's correspondence with an individual and how they change over the period of the relationship.

Viegas et al CHI 2006

# Document Cards

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



Strobelt et al
*TVCG* (InfoVis) '09

# 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

# Example

# Zooming In

# Jigsaw

- Jigsaw is an interactive visualization for document exploration and sense-making
- Variety of visualizations ranging from word-specific, to entity connections, to document clusters
- Primary focus is on entity-document and entity-entity connection
- Entities could be  people, places, dates, organizations, and so on
- Search capability coupled with interactive exploration

- *Stasko, Görg, & Liu*
- Jigsaw: supporting investigative analysis through interactive visualization
- *InfoVis* 2008

# Jisgaw- Document View

# Jisgaw- Graph View



Shows connections between documents and entities in a form of node-link diagram, documents are represented by white rectangles and entities are represented by circles colored differently

# Jisgaw - Document Cluster View

# Visualization of Document Themes

- Someone may not have time to read them all documents
  - Someone just wants to understand the main topics in them

- Look for sets documents that all have common theme
  - Closely related to each other, but different from rest

- The main goal is to discover one or more specific topics and to reflect the relationships among various topics

- Various visualizations developed by *Pacific Northwest National Laboratory*
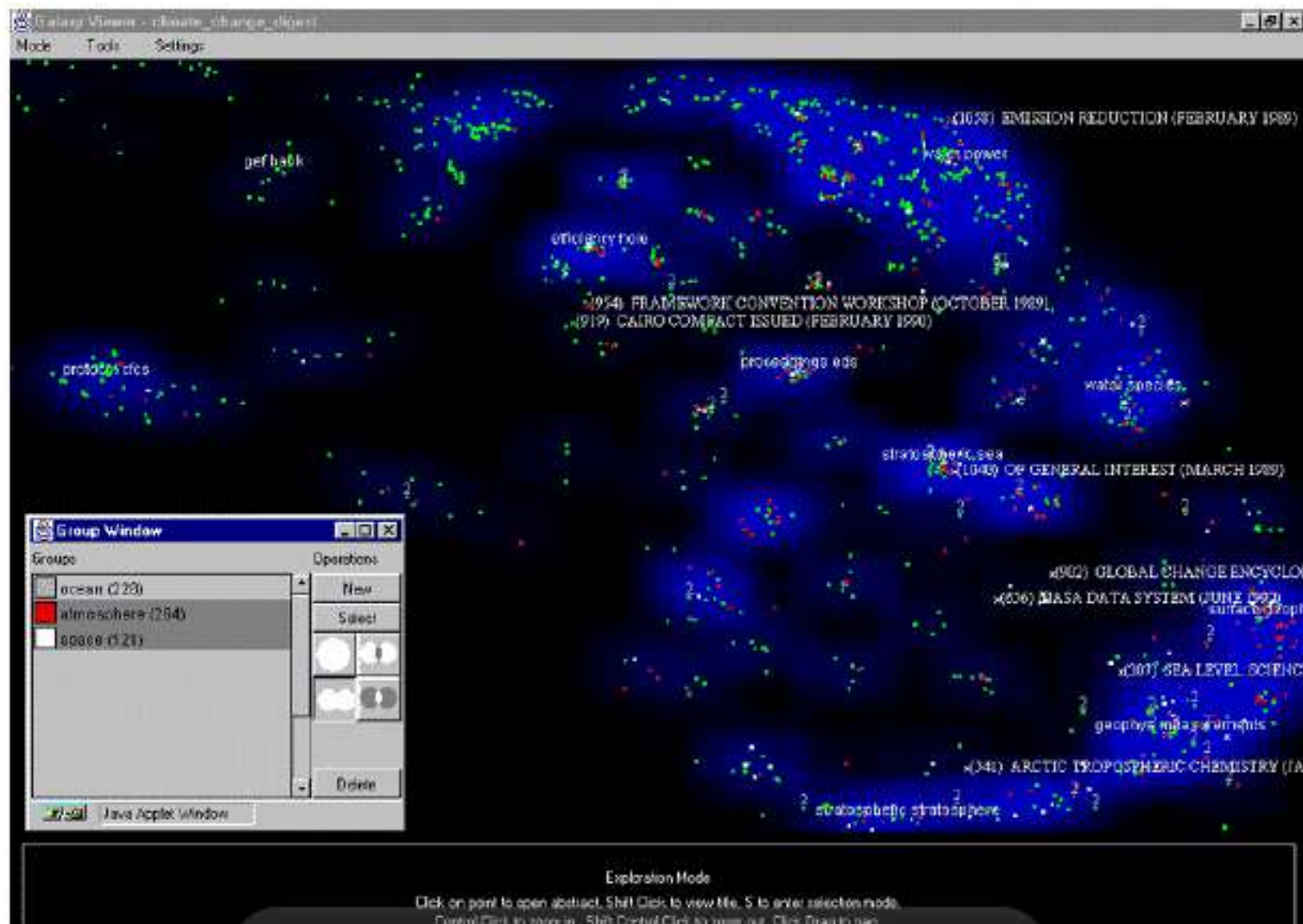
# ThemeView

In the ThemeView visualization, the topics or themes within a set of documents are shown as a relief map of natural terrain. The mountains in the ThemeView™ indicate dominant themes. The height of the peaks indicates the relative strengths of the topics in the document set. Similar themes appear close together, while unrelated themes are separated by larger distances. ThemeView provides a visual overview of the major topics contained in a set of documents. Combined with its exploration tools, ThemeView permits the analyst to identify unanticipated relationships and examine changes in topics over time.



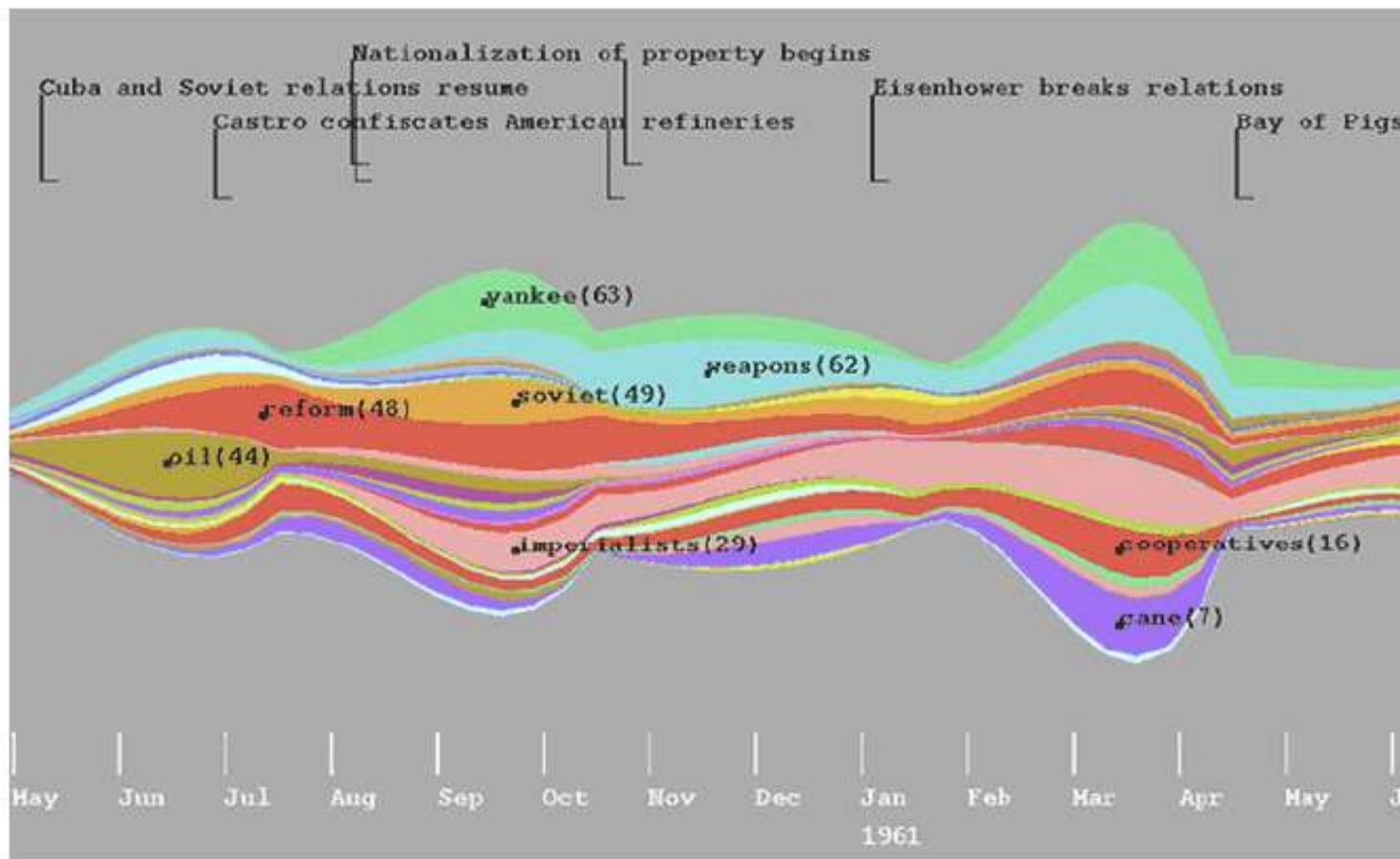https://in-spire.pnnl.gov/

# Web Theme

WebTheme provides a way to investigate and understand large volumes of textual information. It has the ability to harvest data from the World Wide Web using search terms, or by following links derived from user specified URLs. Users can rapidly identify themes and concepts found among thousands of pages of text, and then further explore areas of interest.
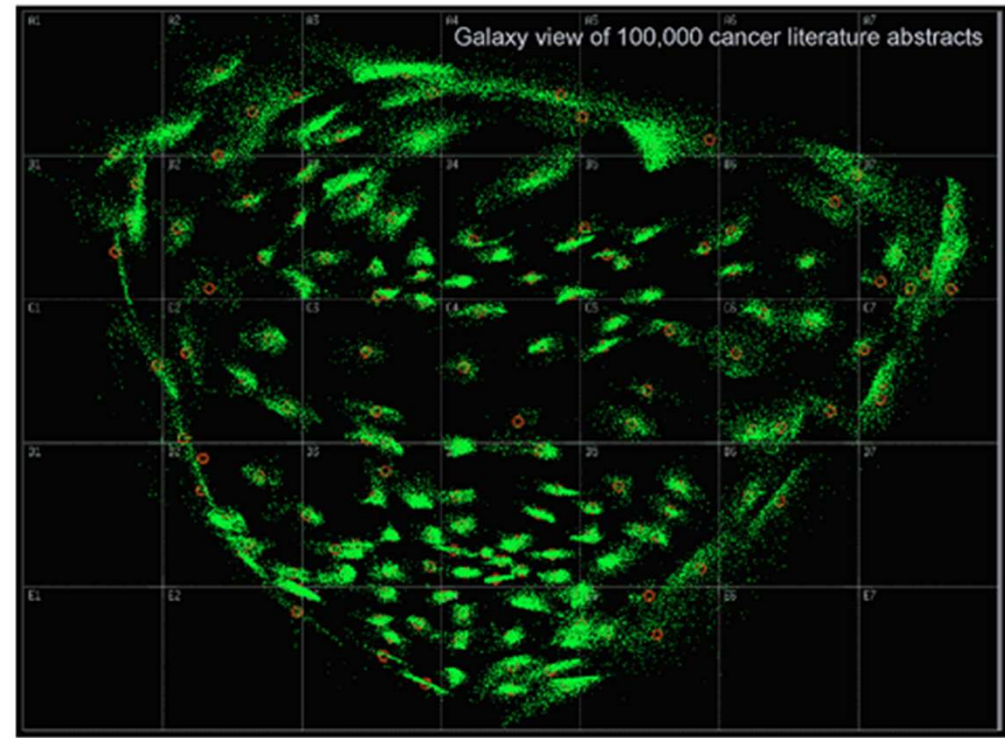
# ThemeRiver



**Temporal Variation of Themes**

Susan Havre, Elizabeth Hetzler, Paul Whitney, and Lucy Nowell. *ThemeRiver: Visualizing thematic changes in large document collections*. IEEE Transactions on Visualization and Computer Graphics, 2002.

# Galaxies

- The Galaxies visualization uses the image of stars in the night sky to represent a set of documents.

- Each document is represented by a single "docustar."

- Closely related documents cluster together while unrelated documents are separated by large distances.

- Several analytical tools are provided with Galaxies to allow users to investigate the document groupings, query the document contents, and investigate time-based trends.



Galaxy view of 100,000 cancer literature abstracts

# Visualizing Text Streams

- Streamit, a dynamic visualization system for exploring text streams.
- Streamit is based on a dynamic force-directed simulation into which text documents are continuously inserted.
- A dynamic 2D display presents the incoming documents
- as a mass particle moving inside a 2D visualization domain,
- Users can explore documents and document clusters on the basis of keywords or topics.
- They can discover emerging patterns online by monitoring the real-time display.
- They can also examine historical data's temporal evolution through animations that play back past streams.
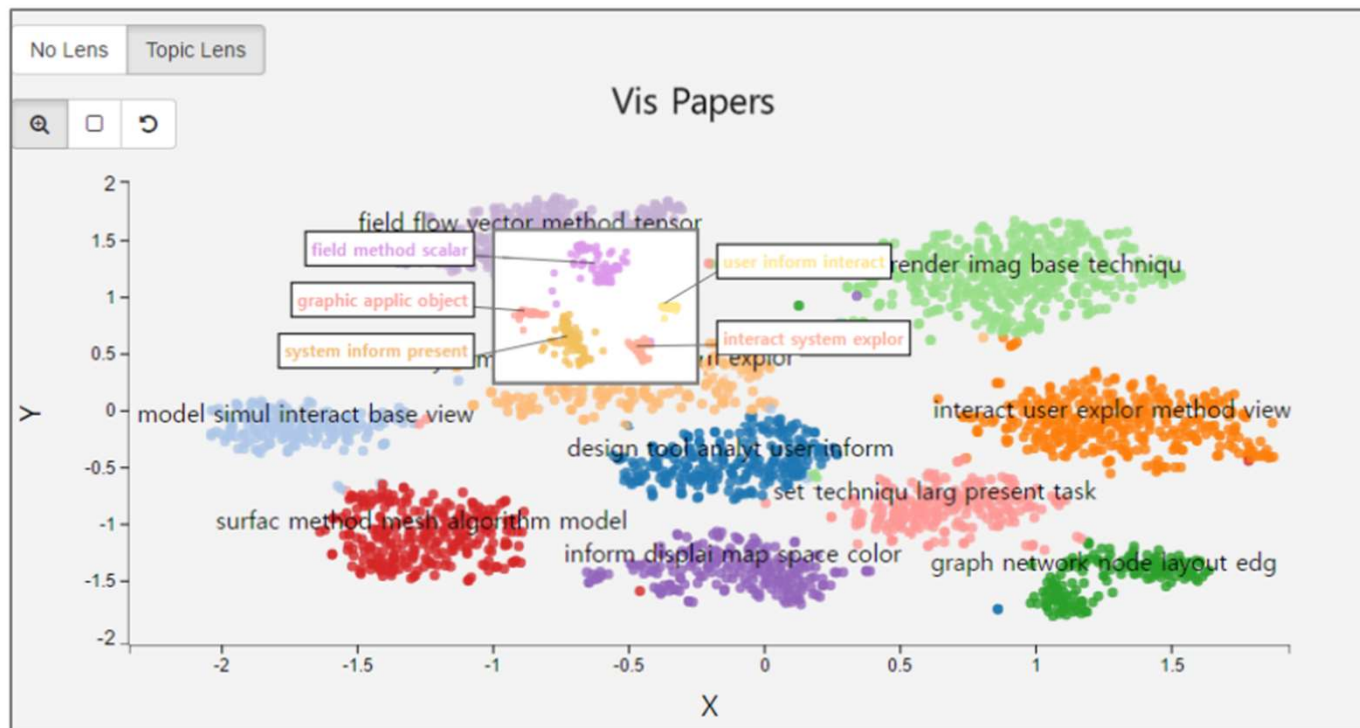
# Streamit



Streamit's interface. (a) The main window displays the particles' movement. (b) The animation control panel lets users navigate through the simulation. The (c) keyword table and (d) document table are synchronized with the particle stream to maintain an up-to-date list of the keywords and documents

Jamal Alsakran, Yang Chen, Dongning Luo, Ye Zhao, Jing Yang, Wenwen Dou, and Shixia Liu. *Real-Time Visualization of Streaming Text with a Force-Based Dynamic System*. IEEE Computer Graphics and Applications, 2012.

# Topic Lens (1/2)

TopicLens is a novel interaction technique that allows a user to dynamically explore data through a lens interface where topic modeling and the corresponding 2D embedding are efficiently computed on the fly.

# Topic Lens (2/2)

- The system initially performs topic modeling and visualizes documents as a scatterplot where the document coordinates are determined by a 2D embedding method

- To generate the 2D scatterplot of documents, a supervised version of t-distributed stochastic neighbor embedding (t-SNE) is used.

- The topic cluster memberships are color-coded.

- The representative keywords are shown in the center of each topic cluster.

- When moving the TopicLens (shown as a small rectangle), we dynamically recompute the topic model and 2D embedding in real time on those documents captured within the lens, revealing their finer-grained topical structure and their visual overview.

- The representative keywords are visualized just outside of the lens pointing to the center of each topic cluster

- Minjeong Kim, Kyeongpil Kang, Deokgun Park, Jaegul Choo, and Niklas Elmqvist. *TopicLens: Efficient Multi-Level Visual Topic Exploration of Large-Scale Document Collections*. IEEE Transactions on Visualization and Computer Graphics, 2017.

# Reading

- Search User Interfaces (Marty Hearst)
  - Chapter 10: Information Visualization for Search Interfaces
  - Chapter 11: Information Visualization for Text Analysis
  - http://searchuserinterfaces.com/book/
- Viegas FB, Wattenberg M, Feinberg J. Participatory visualization with wordle IEEE Transactions on Visualization and Computer Graphics 2009
- Carsten Görg, Zhicheng Liu, Jaeyeon Kihm, Jaegul Choo, Haesun Park, and John Stasko. *Combining Computational Analyses and Interactive Visualization for Document Exploration and Sensemaking in Jigsaw*. IEEE Transactions on Visualization and Computer Graphics, 2013.