

CS-5630 / CS-6630 Visualization for Data Science Text Visualization



Alexander Lex
alex@sci.utah.edu

Text / Language

Features of Text as representation language

abstract, general

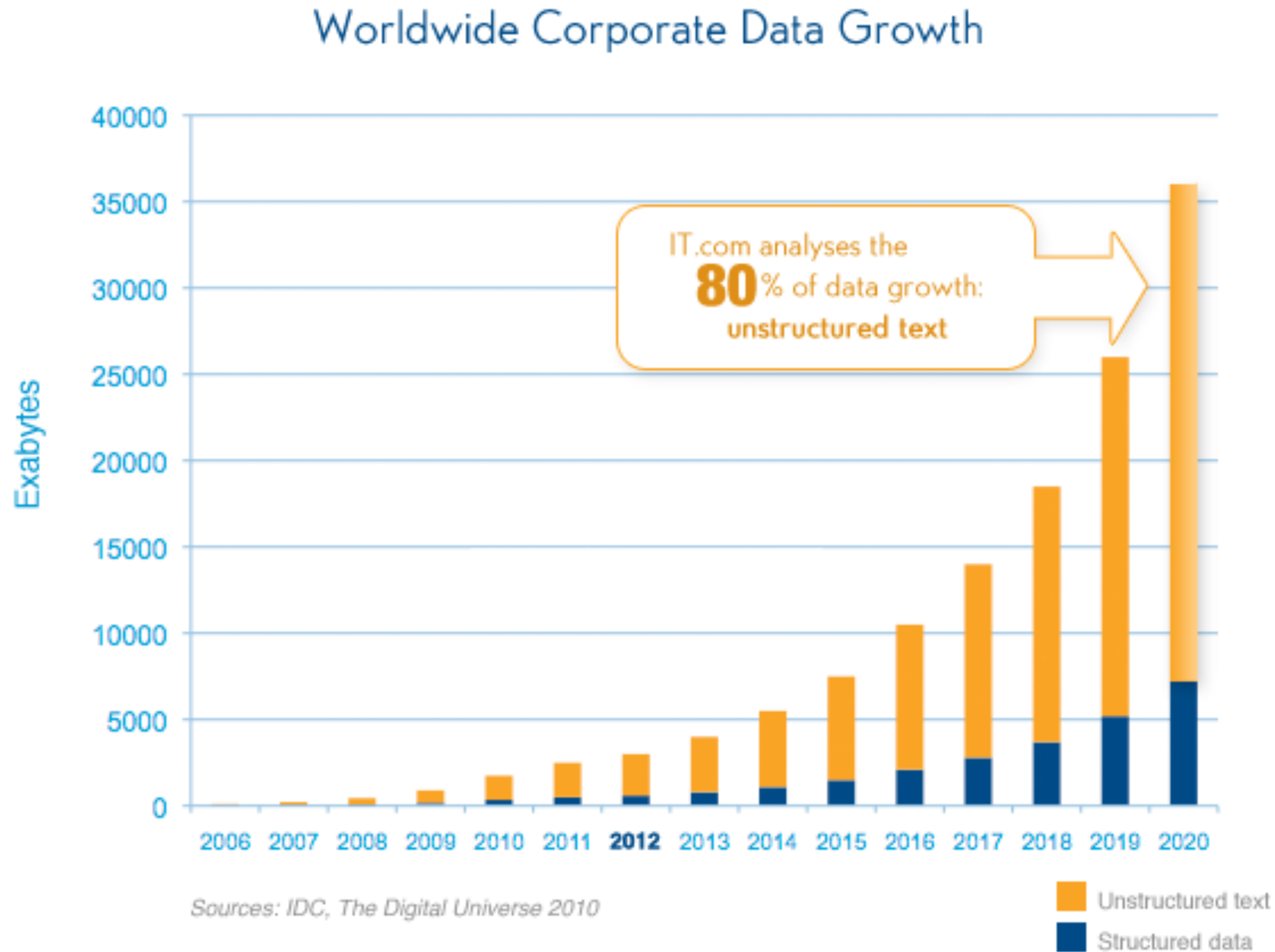
extremely expressive

different across population groups
(countries, accents, religions,...)

linear perception

semi-structured (content: grammar, words, sentences,
paragraphs,.. ; appearance: typography, calligraphy,..)

Why Visualize Text?



Design and Text

Typography:

typefaces (serif, sans-serif, **bold**, *italic*)

point size (10pt, 12pt, 24pt, 36pt..)

line length (alignment: left, right, justified)

vertical: line spacing (leading)

horizontal: spaces between groups of letters (tracking)

Kerning – space between pairs of letters

Ligatures – combining letters to a glyph

*Creating a font type is an art
that requires profound design knowledge*

A V W a
No kerning

A V W a
Kerning applied

fi → fi

fl → fl

Oscars and Typography

Wrong Movie
announced for Best
Picture

Failure of
Typography

Larger Failures in a
Complicated System



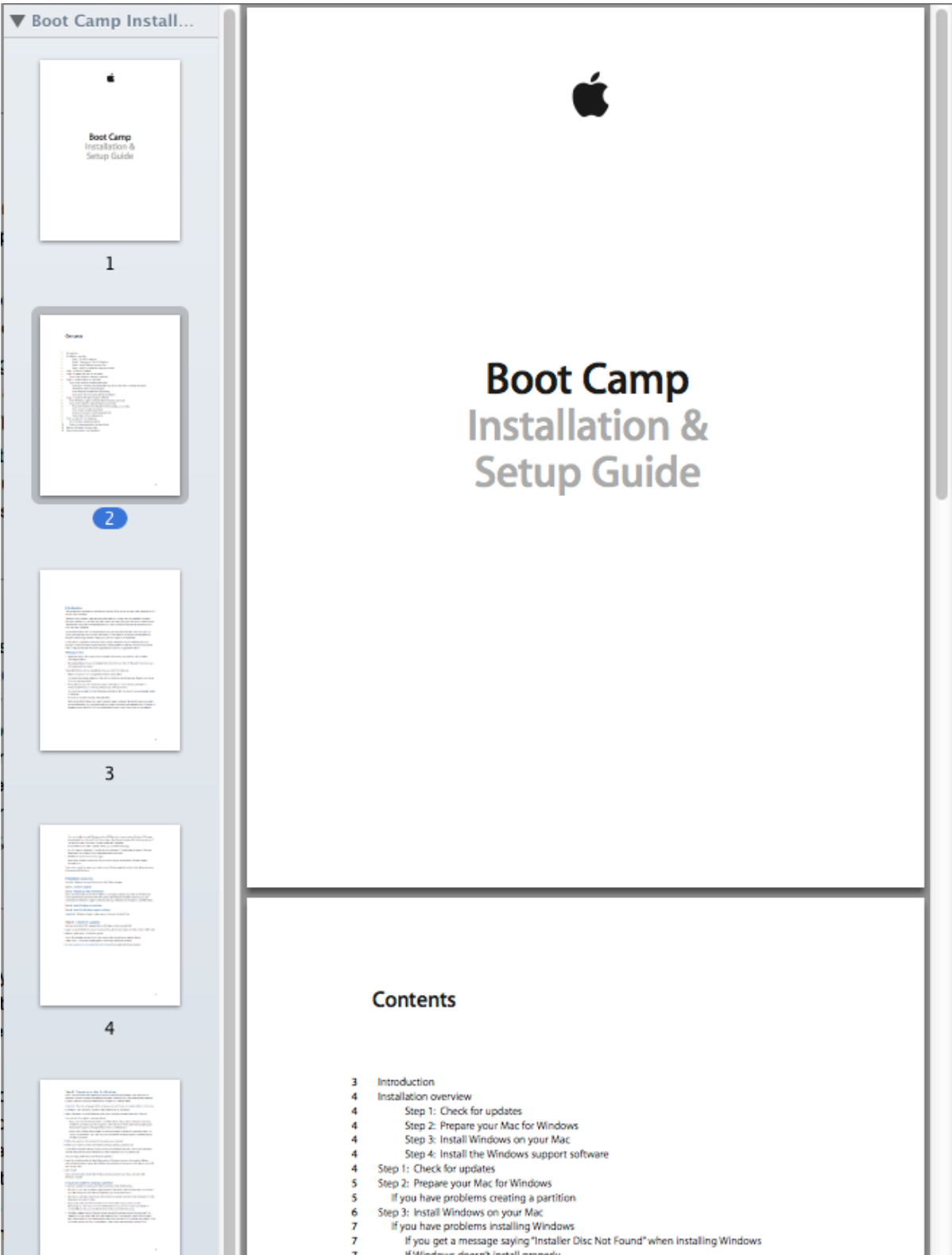
Visualization for “Raw” Text

in daily use..

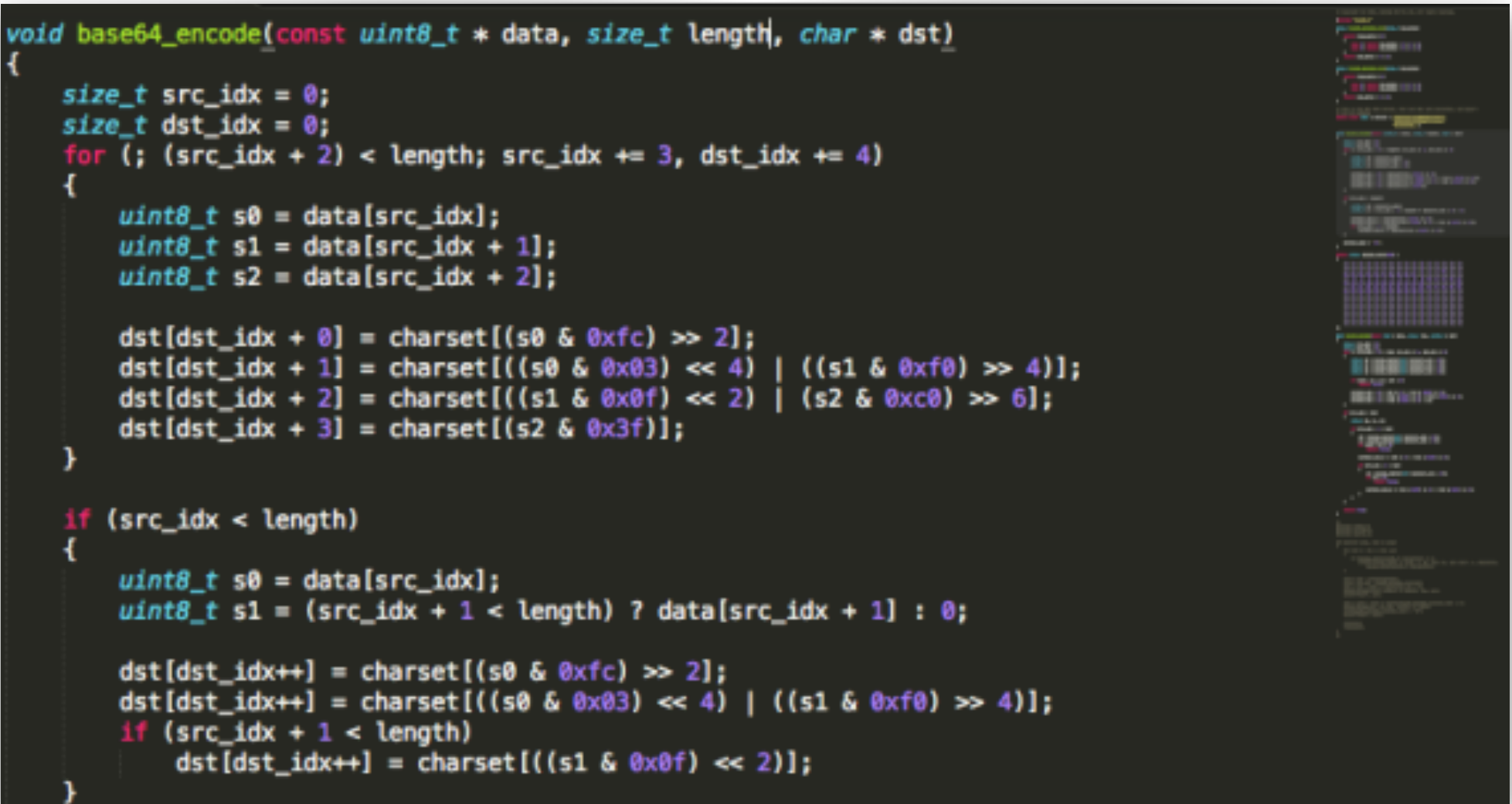
enriched text – hypertext linking (graph navigation)



overview & detail



highlighting semantics



Visualization for “Raw” Text

Document Lens

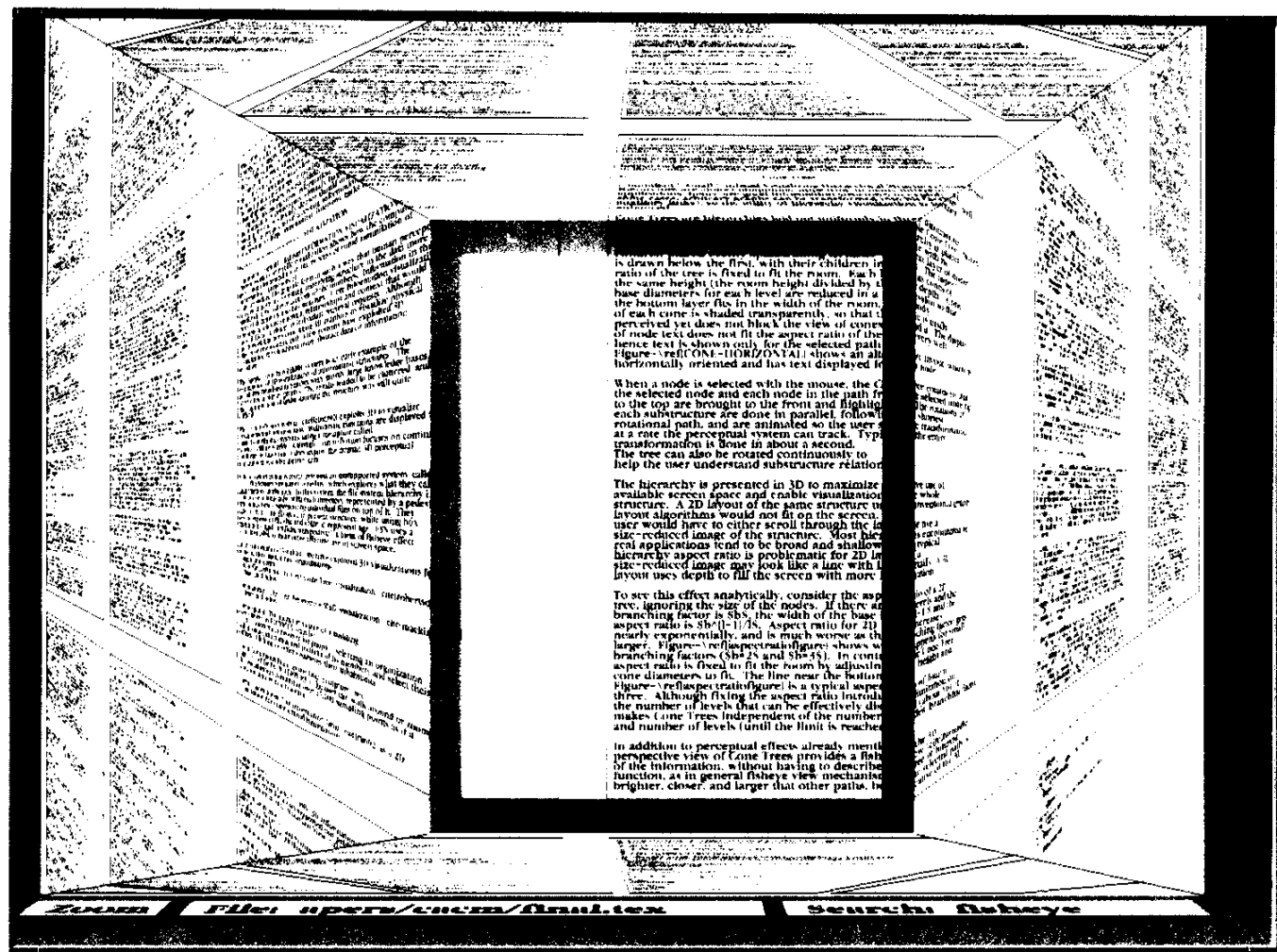


Figure 3: Document Lens with lens pulled toward the user. The resulting truncated pyramid makes text near the lens’ edges readable.

Robertson, George G., and Jock D. Mackinlay

The document lens

Proceedings of the 6th annual ACM symposium on User interface software and technology. ACM, 1993.

Document Thumbnails with Variable Text Scaling

A. Stoffel, H. Strobel, O. Deussen, D. A. Keim

Computer Graphics Forum, volume 31 issue 3 pp.

Visualizing Search Results

Eurographics Conference on Visualization (EuroVis) 2012
S. Bruckner, S. Miksch, and H. Pfister
(Guest Editors)

Volume 31 (2012), Number 3

Document Thumbnails with Variable Text Scaling

A. Stoffel and H. Strobel and O. Deussen and D. A. Keim

University of Konstanz, Germany

Abstract

Document reader applications usually offer an overview of the layout for each page as thumbnail view. Reading the text in these becomes impossible when the font size becomes very small. We improve the readability of these thumbnails using a distortion method, which retains a readable font size of interesting text while shrinking less interesting text further. In contrast to existing approaches, our method preserves the global layout of a page and is able to show context around important terms. We evaluate our technique and show application examples.

1. Motivation

The user interface of

such as Adobe Reader, consists of a detail view and one or more views for navigation within documents, such as a table of contents, and a thumbnail view providing page pre-views. In addition, most document viewers offer a keyword search functionality where the occurrence of keywords is highlighted in the detail view. However, the navigation views of document viewers (e.g., thumbnails) typically do not show the occurrence of keywords in the documents.

has to

step through all occurrences of the keyword within the detail view as scrolling the pages.

To avoid this, we propose to highlight the keywords in the thumbnail view. Using the thumbnail view reduces the and the user is pointed

pages. In addition, thumbnails can be useful for retrieval

if the users are trying

know [CvDRH99, DC02]. Due to the small size of text in thumbnails, the highlighting should in addition increase the size of the keywords and their context at first to make the text better readable and second to allow a simple dis-

ambiguation of keywords by their context. For instance, it

about “user” or “user inter-

face” keyword “user” would

to a user defined interest

The global structure of a page, namely the position of im-

is used that highlights the keywords and their context. Other applications might use a different interest function, for instance a sentiment score could be used to create thumbnails for sentiment analysis.

2. Related Work

Three different techniques are currently used for handling document overview and navigation: abstraction from the document with pixel based representations, thumbnails with different highlighting techniques, and semantic zooming.

A common pixel based technique is TileBars [Hea95], which visualizes the length of documents and the distribution of search terms within these documents with a rectangular pixel-based visualization. Byrd [Byr99] combines the scrollbar of the document view with a pixel visualization of

allowing the user to scroll

rence of the terms. Both techniques do not show the context

and a user has to

from the detail view in

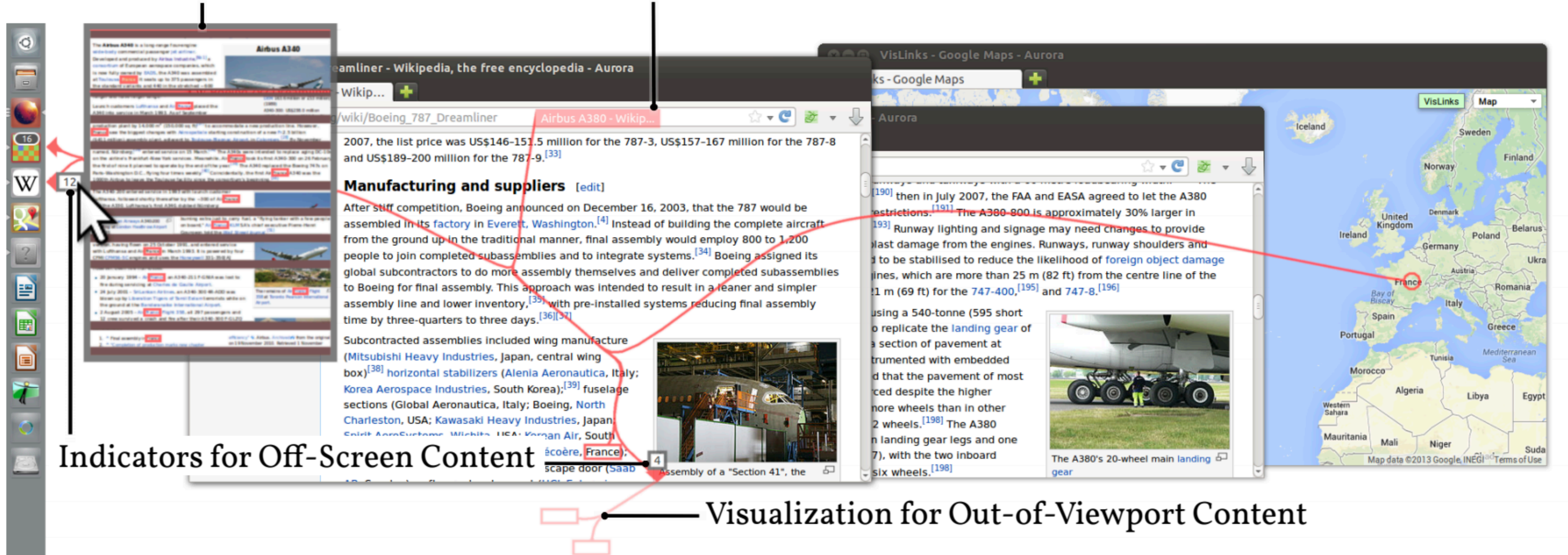
order to access the context of the search terms.

Thumbnails, small version of the document or page, are commonly used for overview and navigation. The space-filling thumbnail approach of Cockburn et al. [CGA06] avoids scrolling in the overview of a document, by positioning the thumbnails of all pages on a grid on the screen and resizing the thumbnails to fit the window size. Suh et al. [SWRG02] combined the thumbnails with popouts, which highlight search terms by rendering them in a readable size with a semi-transparently colored background above

Visualizing Hidden Text

Smart Preview for Off-Screen Content

Indicators for Occluded Windows



Working with Text

unstructured text



4 x 't'
3 x 'u'
2 x 'r'
2 x 'e'

...

structured data

Structured Text Features

simple counts (bag of words)
used for similarity measures

	princess	dragon	castle
doc1	1	1	1
doc2	0	0	1

Processing to Derive Features

Typical steps are:

- cleaning (regular expressions)

- sentence splitting

- change to lower case

- stopword removal (most frequent words in a language)

- stemming

- POS tagging (part of speech)

- noun chunking

- NER (name entity recognition)

- deep parsing - try to “understand” text.

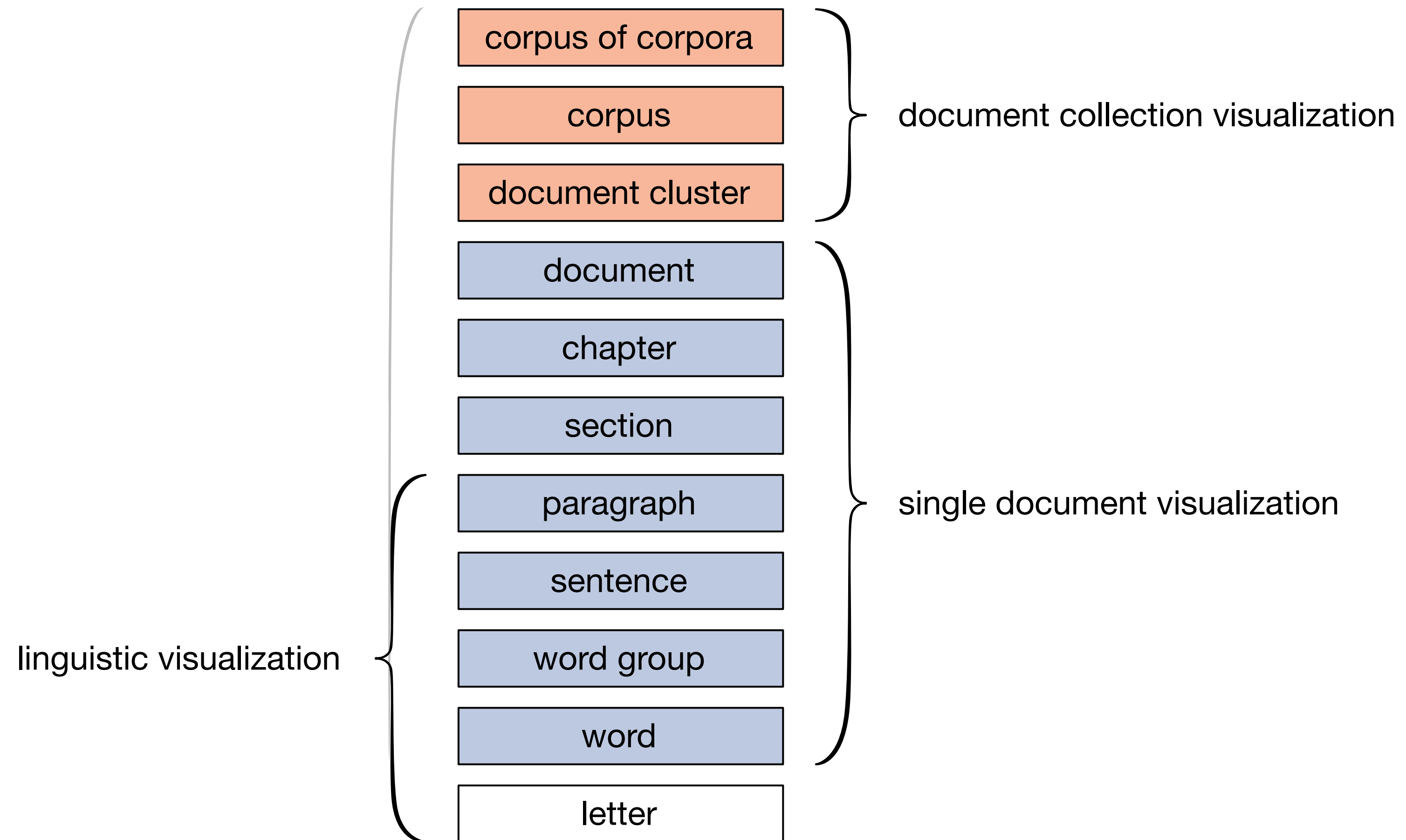
Text features are complicated

Toilet out of order. Please use floor below.

One morning I shot an elephant in my pajamas. How he got in my pajamas, I don't know.

Did you ever hear the story about the blind carpenter who picked up his hammer and saw?

Text Units Hierarchy



Types of Text Visualizations

Document Visualization

Corpus Visualization

Visualization for NLP

Creativity Support

Document Visualization

Wordle

Frequency-based

words that occur often are large

Can vary font type,
size, color, etc.

Uses stop-word removal



Wordle vs Tag Cloud

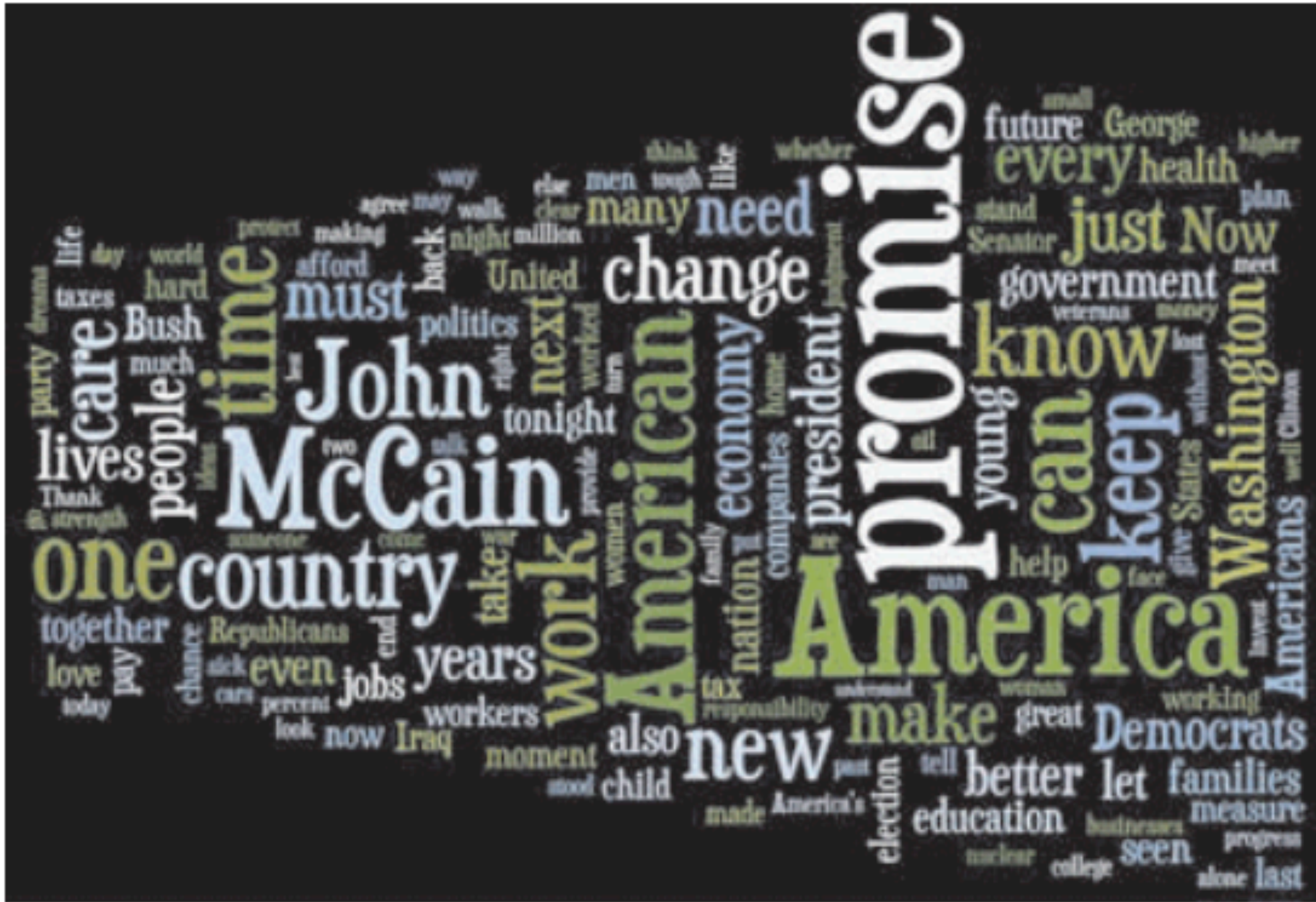


Fig 2: Wordle vs. Tag Cloud of Barack Obama's speech at the Democratic Convention in 2008.

Opinion

Use Wordle if you want something evocative.

Don't use Tag Cloud! (Looks bad, not very useful)

Use structured approach instead

- Top keywords with counts

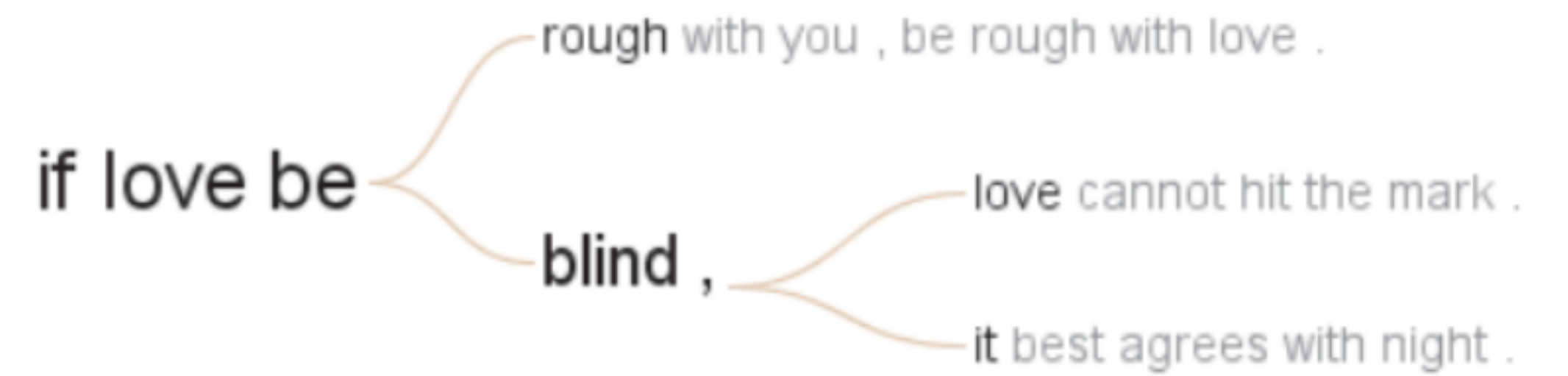
- Maybe group by topics

Word Tree

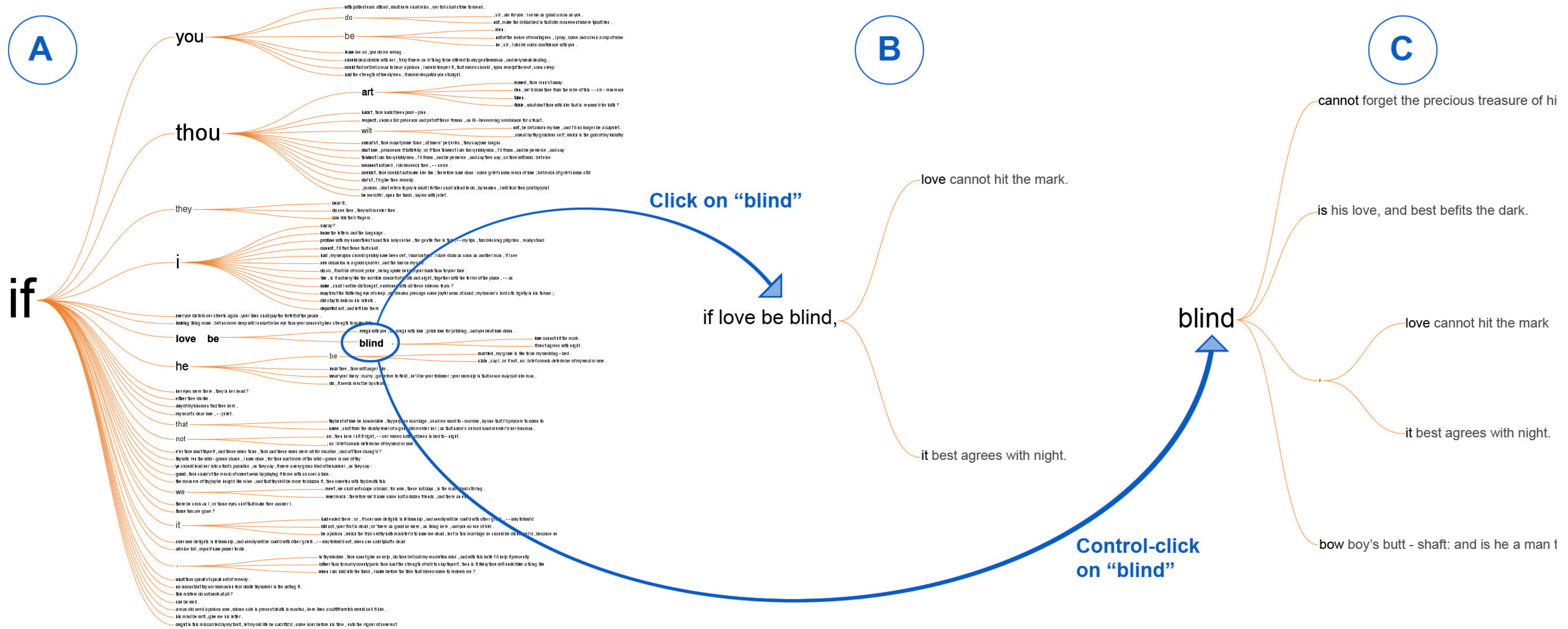
Text

if love be rough with you , be rough with love .
if love be blind , love cannot hit the mark .
if love be blind , it best agrees with night .

WordTree



Search for “if” in romeo & Juliet



The word tree, an interactive visual concordance

M Wattenberg, FB Viégas

Visualization and Computer Graphics, IEEE Transactions on 14 (6), 1221-1228

Phrase Net for Bible and “begat”

PhraseNets

1 *You create the word sequence filter:*

WORD1 and **WORD2**

2 *Many Eyes finds this word relationship in Jane Austen's text:*

Her manners were pronounced to be very bad indeed,
a mixture of **pride and impertinence**; she had no
conversation, no stile, no taste, no beauty.

3 *Many Eyes creates the word graph:*

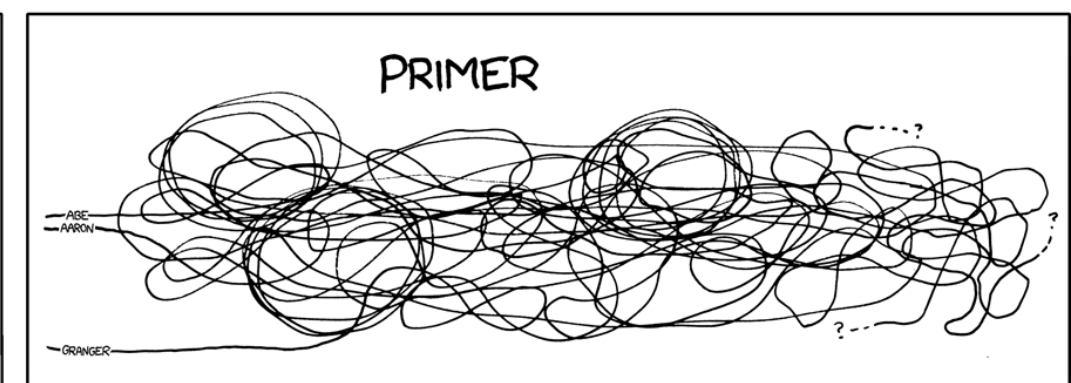
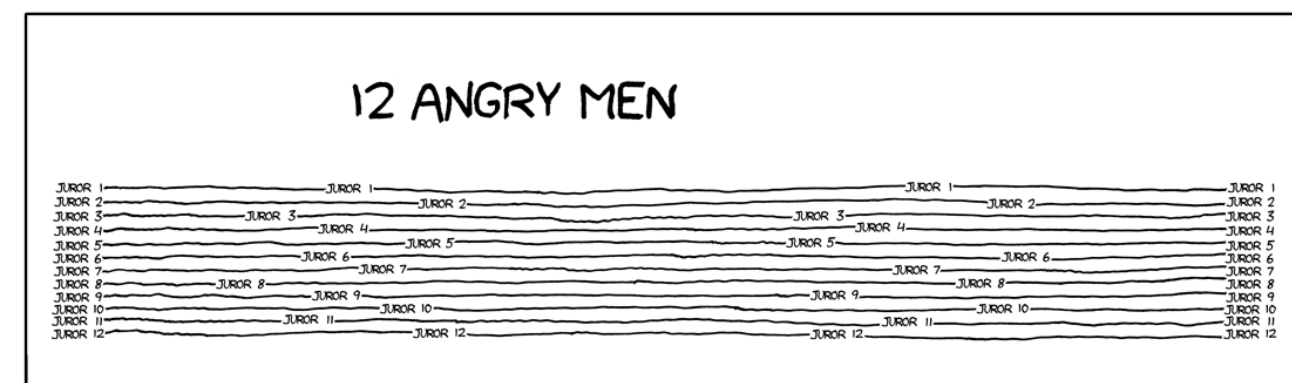
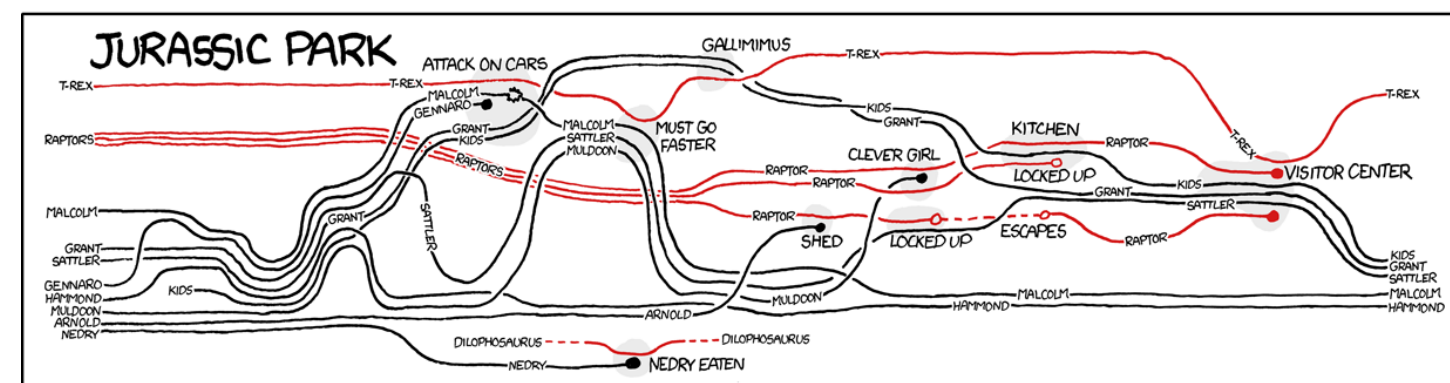
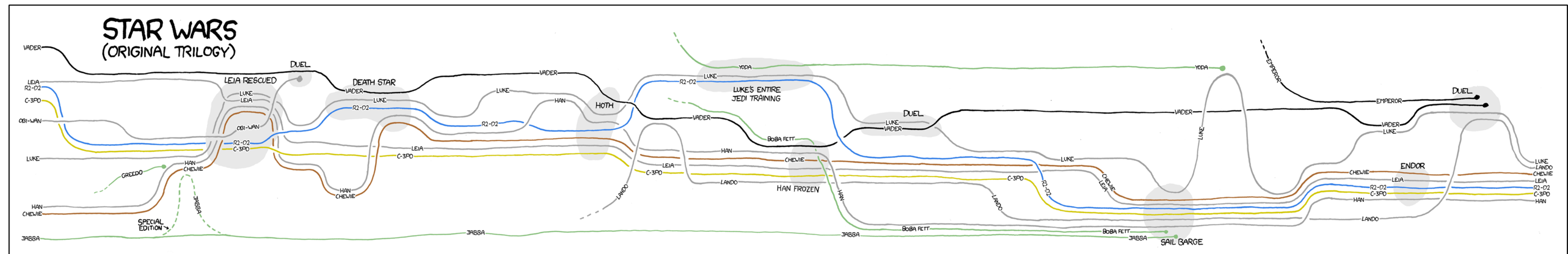
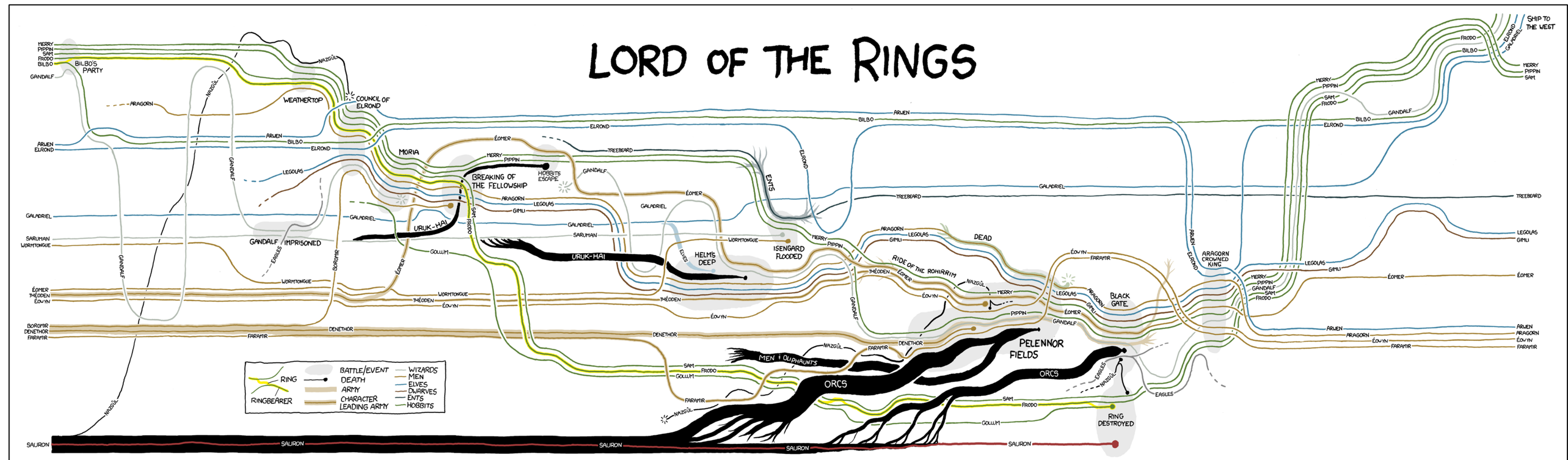
pride → **impertinence**

Frank van Ham, Martin Wattenberg, and Fernanda B. Viegas.

Mapping Text with Phrase Nets.

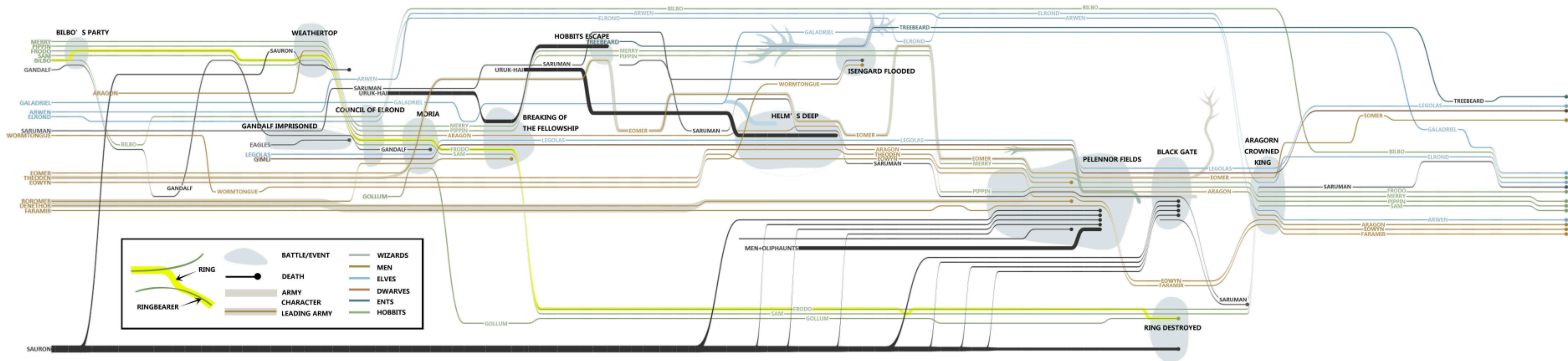
IEEE Transactions on Visualization and Computer Graphics 15, 6 (November 2009)

THESE CHARTS SHOW MOVIE CHARACTER INTERACTIONS.
THE HORIZONTAL AXIS IS TIME. THE VERTICAL GROUPING OF THE
LINES INDICATES WHICH CHARACTERS ARE TOGETHER AT A GIVEN TIME.



<https://xkcd.com/657/>

StoryFlow: Tracking the Evolution of Stories



[Liu 2013]

Visualizing Corpora

Text Corpora

Varied Goals:

Discover interesting documents

Summarize Documents

Classify Documents

Extract Facts (Intelligence Analysis)

Rich Information:

Document Metadata

Authors, date, type,

Paragraphs, figures...

Revisions, annotations, comments,

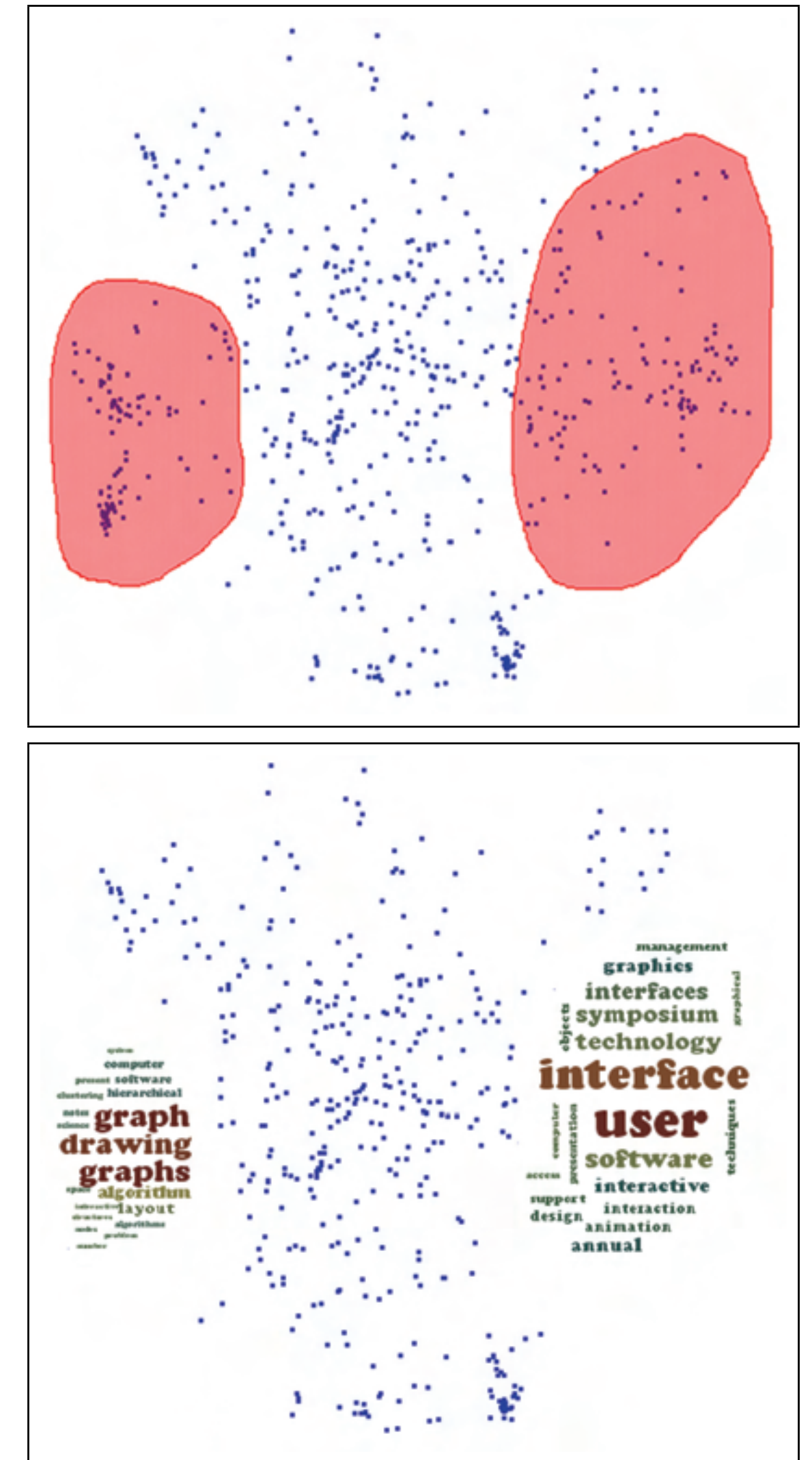
Bohemian Bookshelf



Corpora: MDS Approaches

use bag-of-word to project
documents w.r.t. text similarity
into a landscape

(only) one example

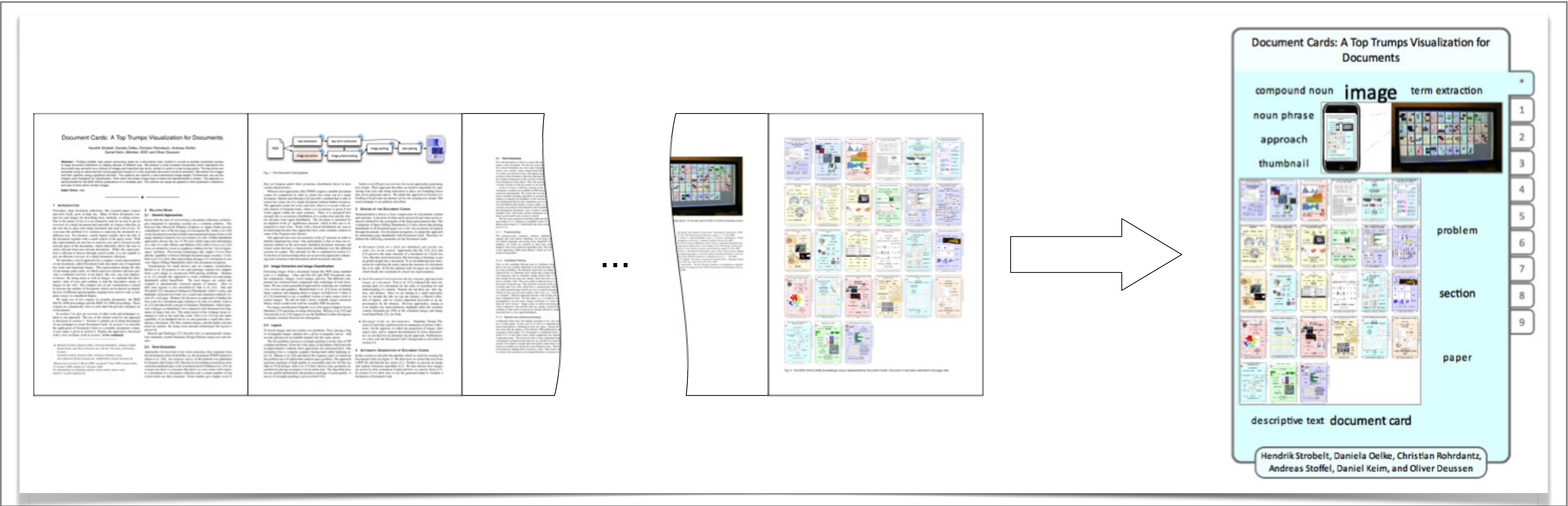


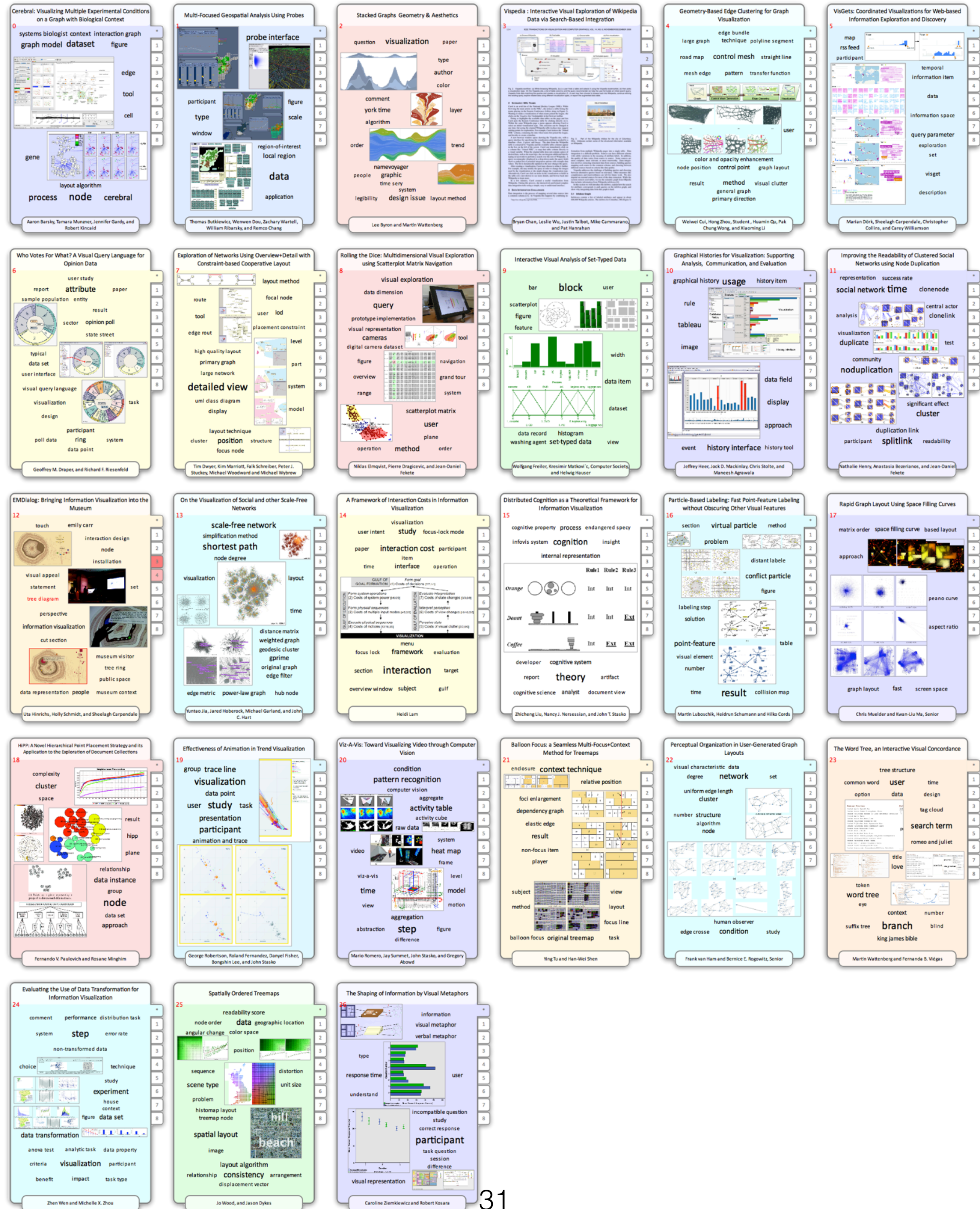
Fernando V. Paulovich, Franklina M. B. Toledo, Guilherme P. Telles, Rosane Minghim, and Luis Gustavo Nonato.
Semantic Wordification of Document Collections.
Comp. Graph. Forum 31, 3pt3 (June 2012)

Figure 5: A user can interactively draw a region (polygon) containing a subset of documents of interest (top figure). Keywords are extracted from the selected document and their corresponding word cloud is built inside the user-defined region (bottom figure).

DocumentCards

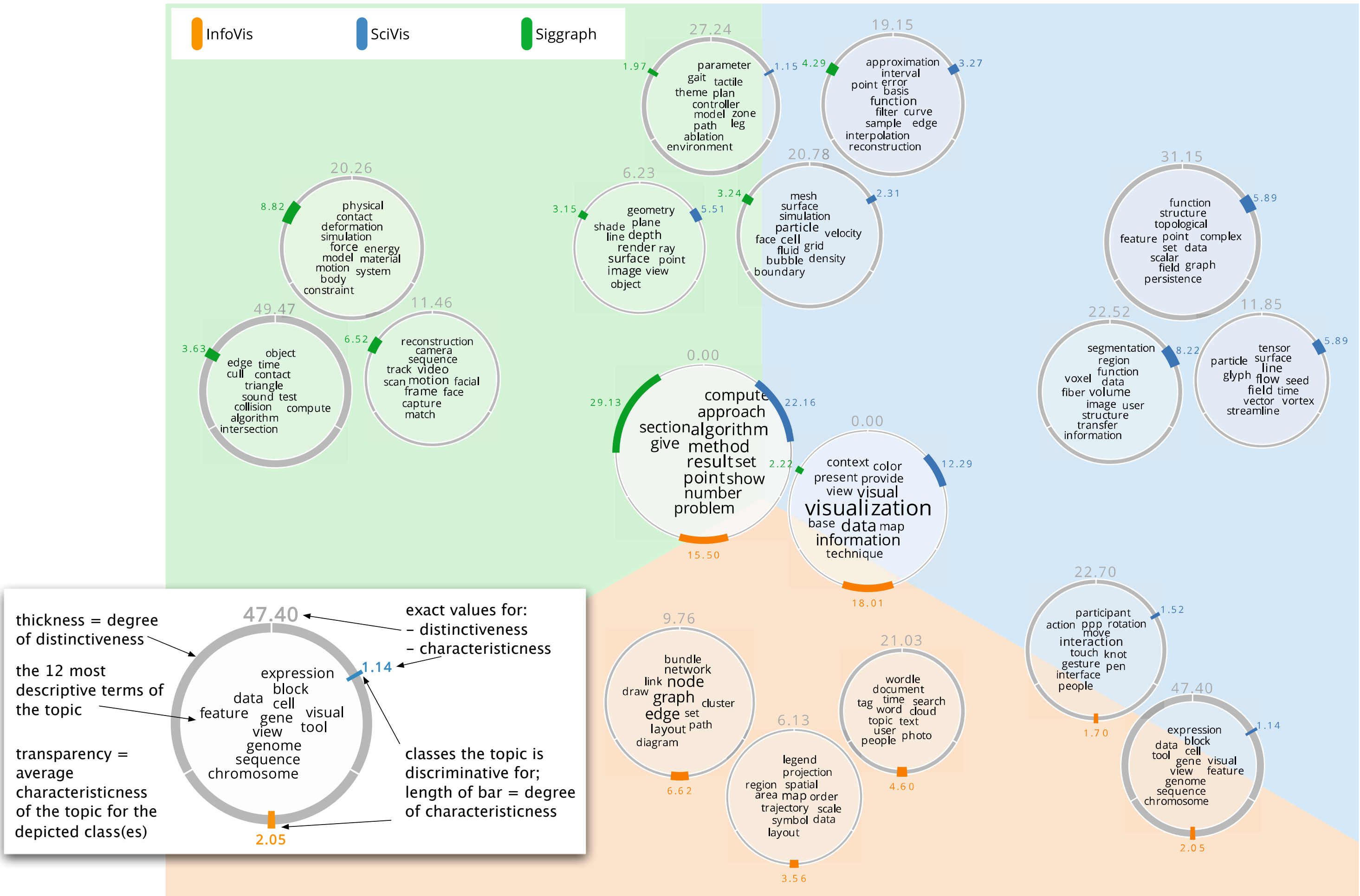
summarize scientific documents using
important terms and important figures
represent the document's content as a mix of figure and text





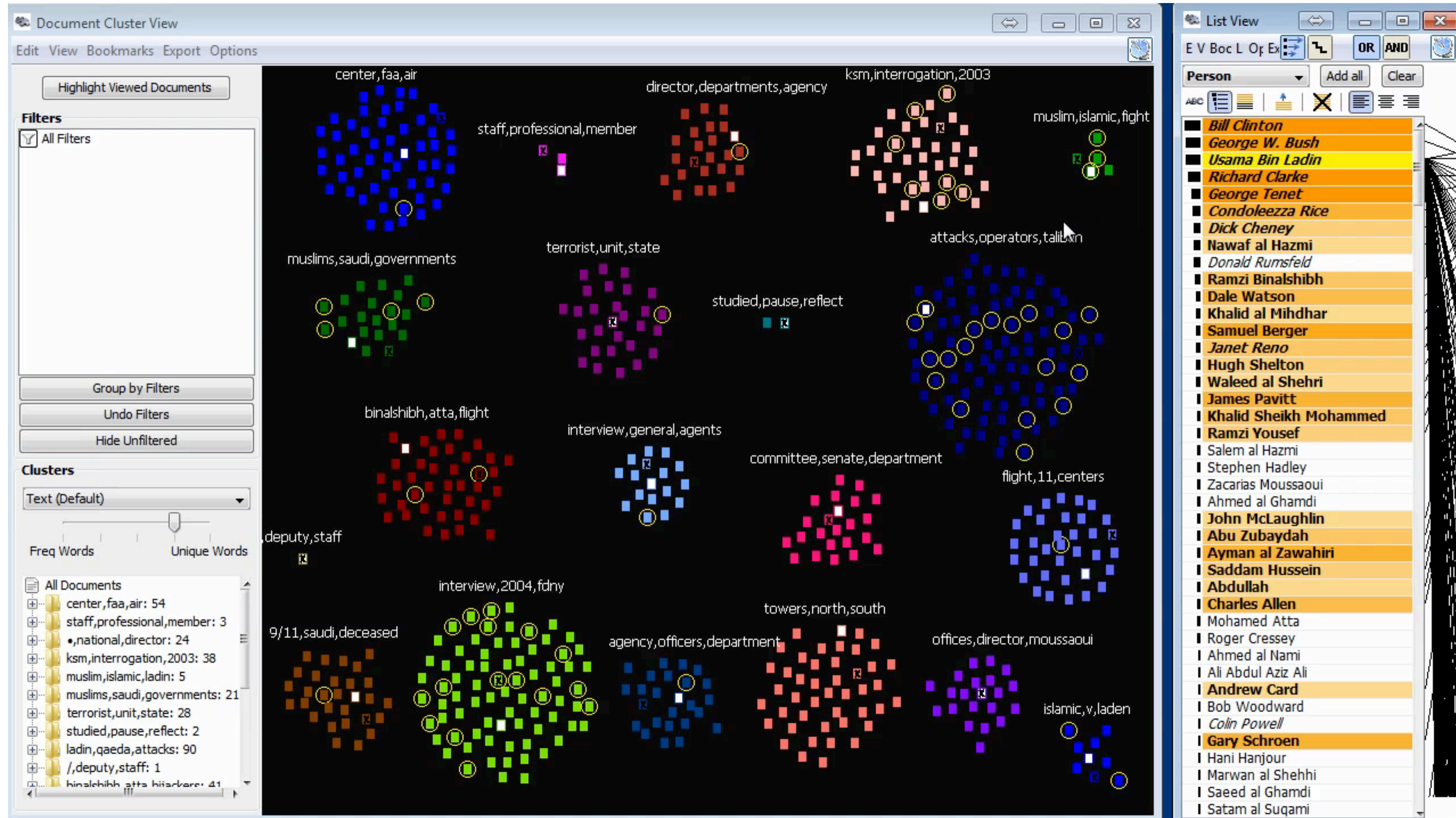
Compare Corpora

Compare topics
between text
collections

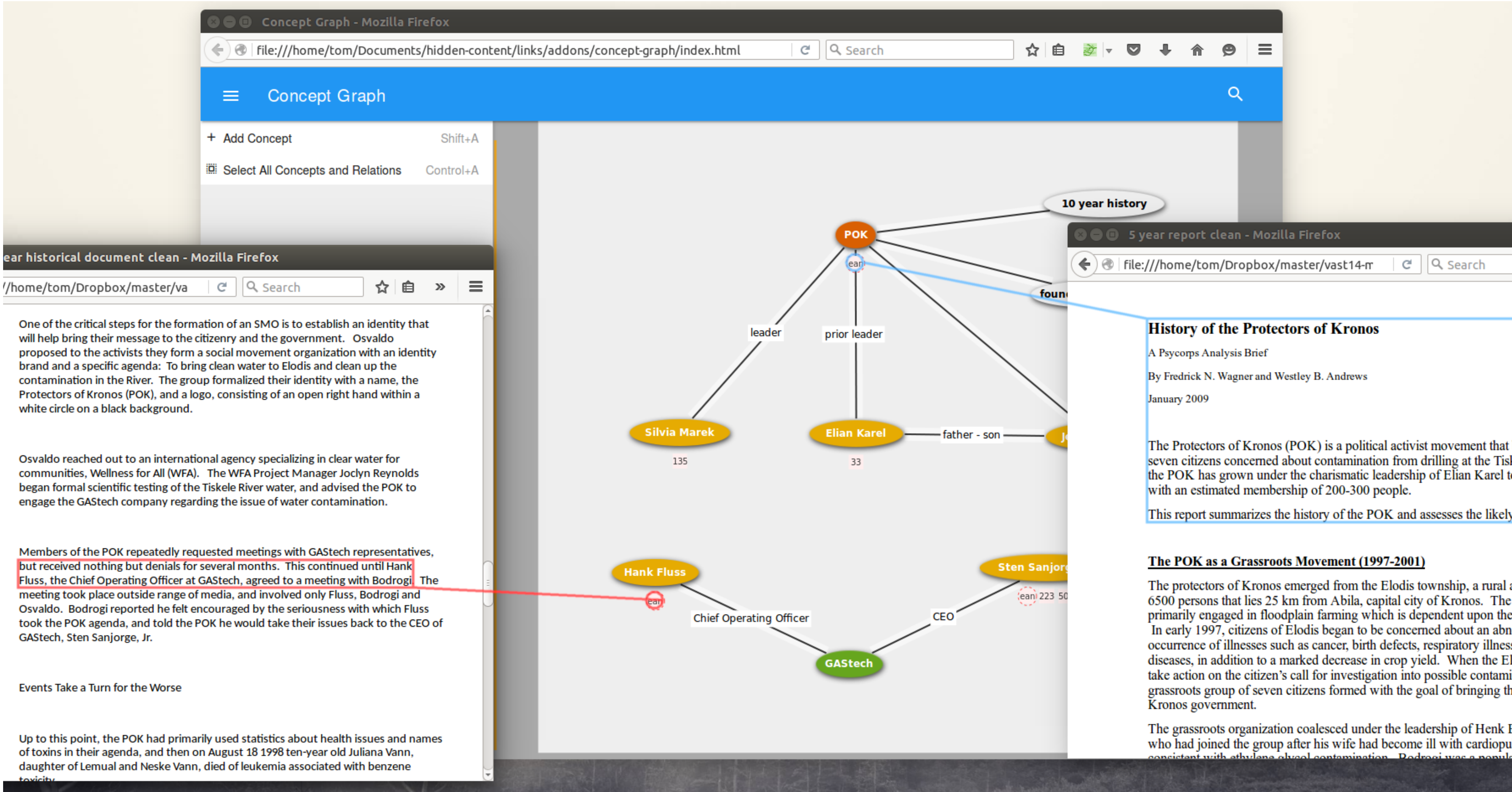


Comparative Exploration of Document Collections: a Visual Analytics Approach (<http://ditop.hs8.de>) Figure 1: Comparison of 495 papers of InfoVis, SciVis, and Siggraph (discrimination threshold = 6, number of topics = 30)
D. Oelke, H. Strobel, C. Rohrdantz, I. Gurevych, and O. Deussen

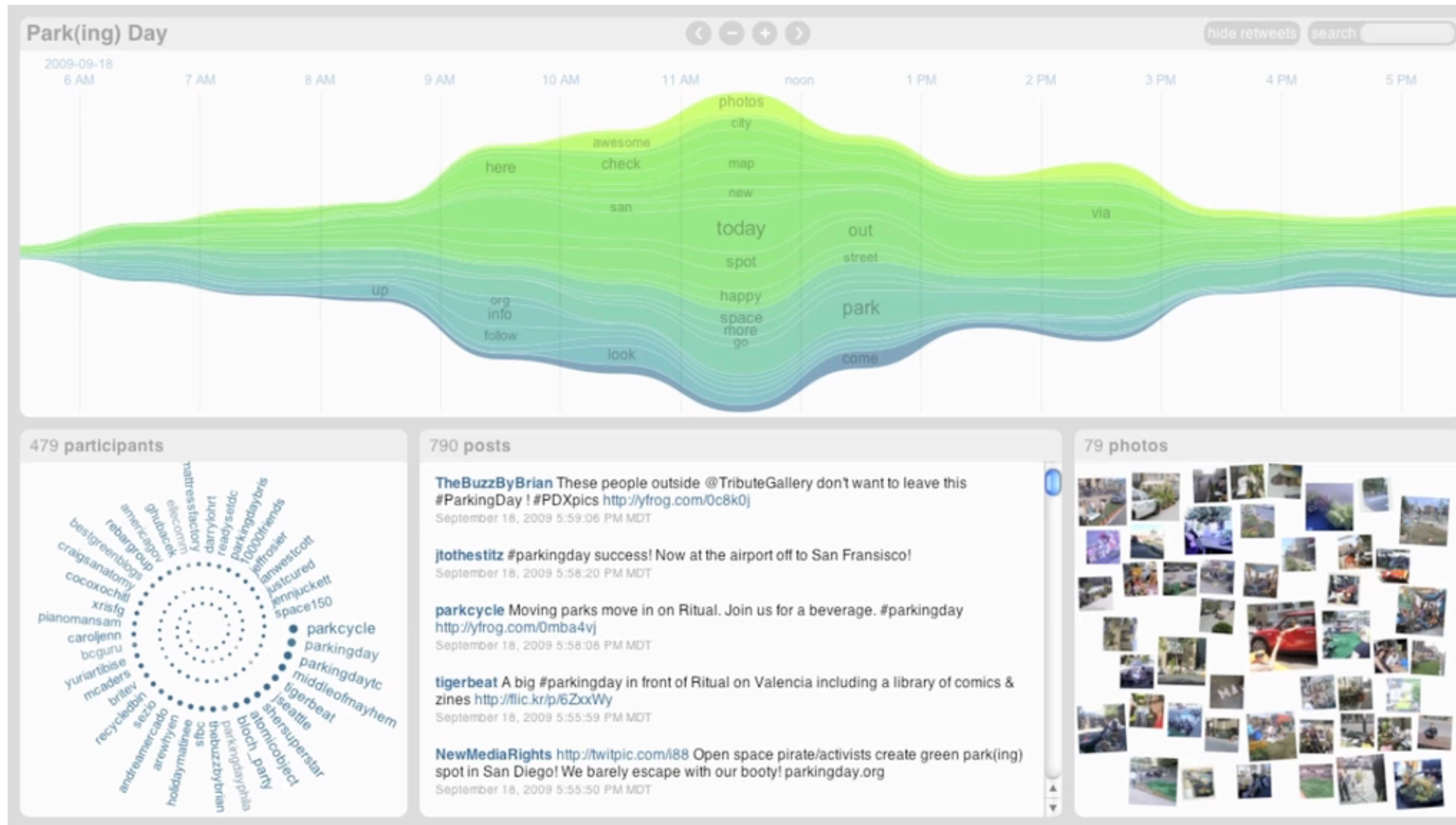
JigSaw – Intelligence Analysis



Extracting and Linking Info From Documents



Collection of Tweets



Visualization for NLP

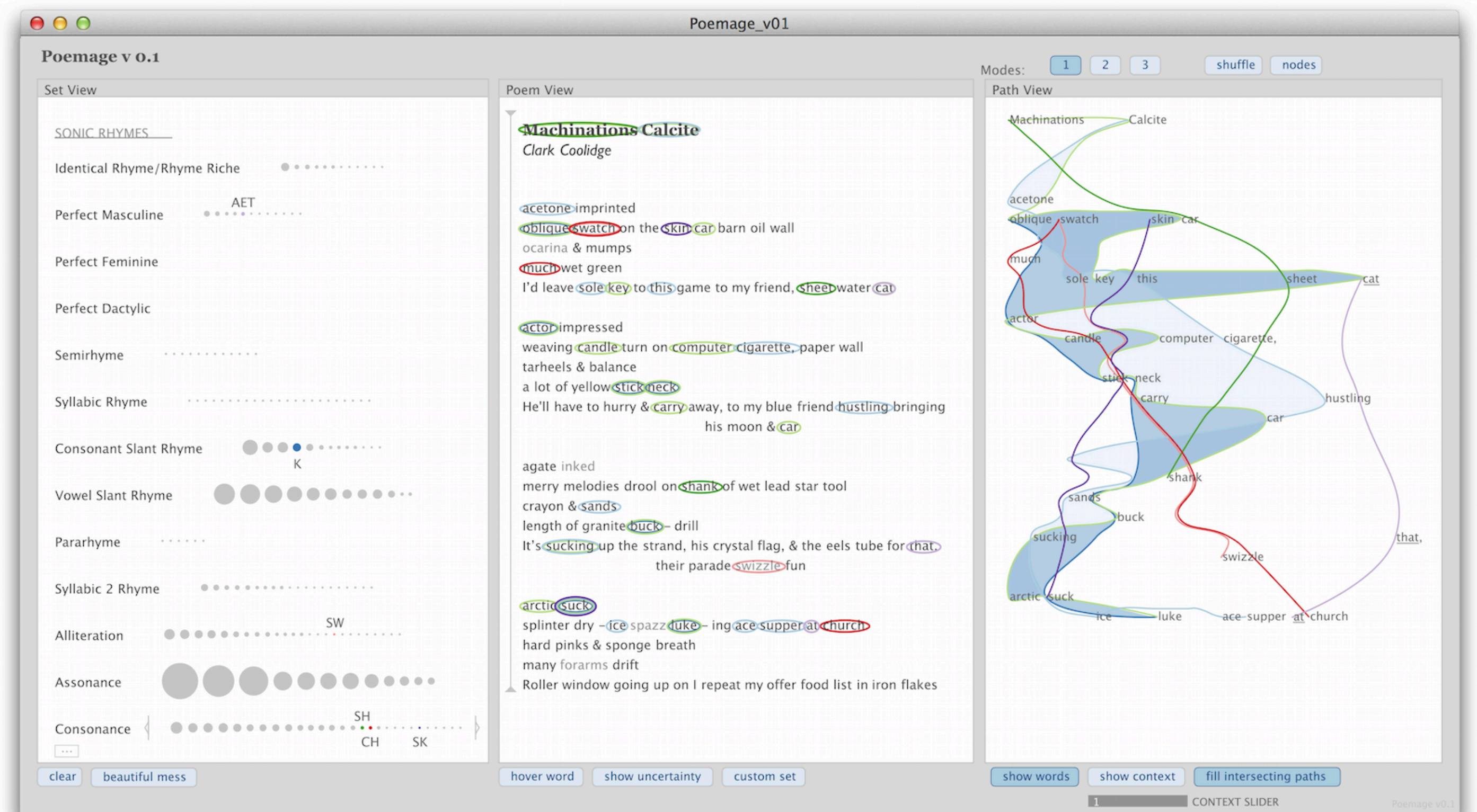
GLTR: Statistical Detection and Visualization of Generated Text, Gehrmann, Strobelt Rush: <http://gltr.io/dist/index.html>

LSTMVis: Visual Analysis for Recurrent Neural Networks, Strobelt et al.: <http://lstm.seas.harvard.edu/>

Visual Exploration of Semantic Relationships in Neural Word Embeddings. Liu et al.

Visualization for Creativity Support

Poemage: Visualizing the Sonic Topology of a Poem. McCurdy et al. <http://www.sci.utah.edu/~nmccurdy/Poemage/>



<http://textvis.lnu.se/>

Text Visualization Browser

A Visual Survey of Text Visualization Techniques

Provided by ISOVIS group

[About](#) [Add entry](#) [Contact](#)

Techniques displayed:
141

Search:

Time filter:
1976 2014

Analytic Tasks

Visualization Tasks

Data

Source

Properties

The grid displays 141 different text visualization techniques. These include word clouds of various shapes and colors, network graphs showing relationships between nodes, treemaps representing hierarchical data, and various types of charts such as bar charts, line graphs, and radar charts. Some visualizations are more complex, combining multiple techniques or using interactive elements. The techniques are arranged in a grid that is approximately 10 columns wide and 14 rows high, with the last row containing only one visualization.