

# Text and Document Visualization

Hendrik Strobelt - [hstrobelt@seas.harvard.edu](mailto:hstrobelt@seas.harvard.edu)  
housing day 2015



CS 171 - Visualization



**HARVARD**  
School of Engineering  
and Applied Sciences

/Users/hen> whoami



Text Visualization



Visualization for Sciences



Layout



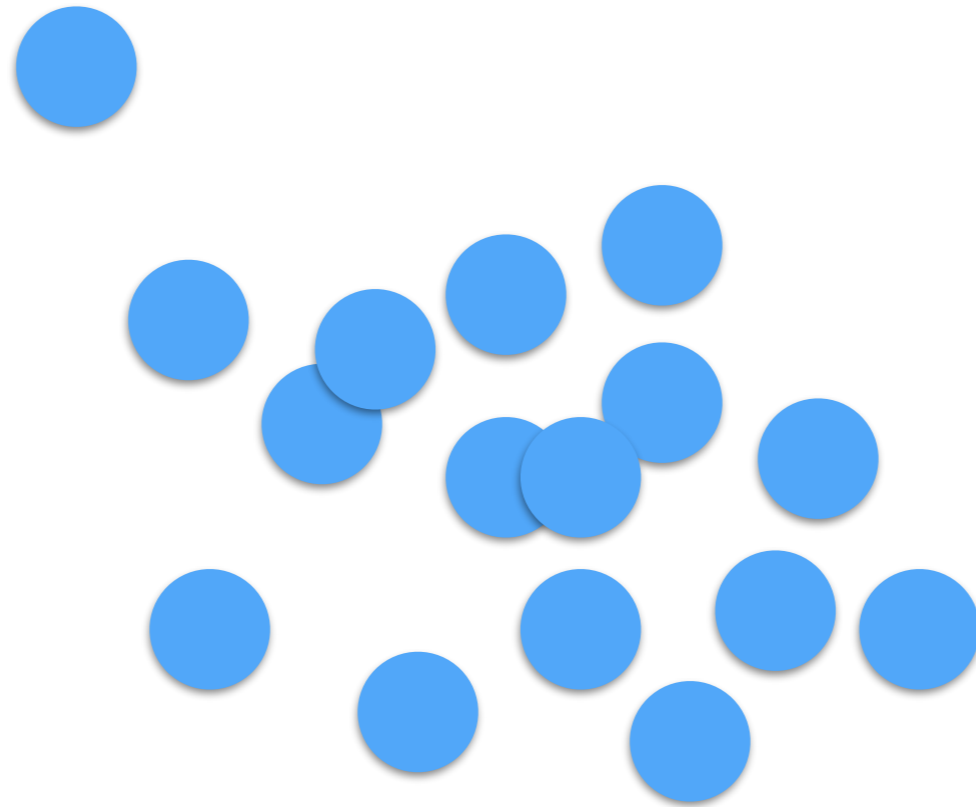
# This Week

- HW2 (due to **FRIDAY** — 11:59 pm):
  - include design studio solutions
- Section 6 special **TODAY** at 4pm MD G125

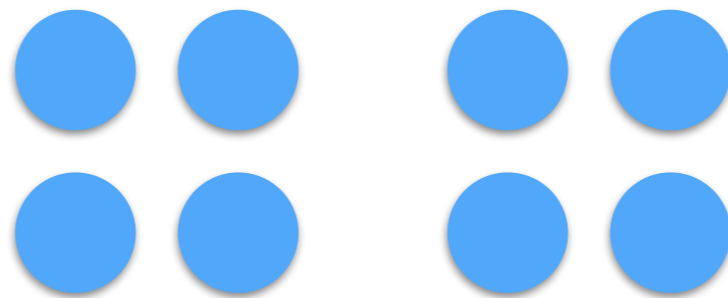
# A little experiment

**Task: How many dots?**

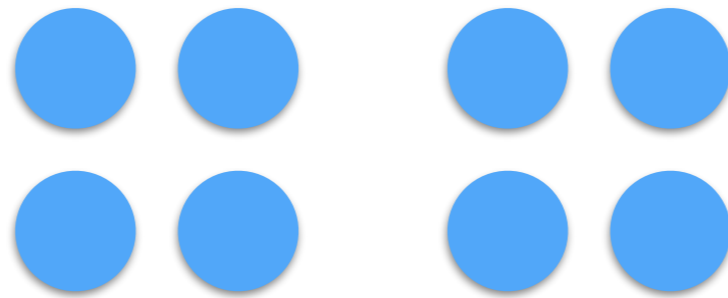
# A little experiment



# A little experiment



**Task: How many dots?**



# brief history

(western view)

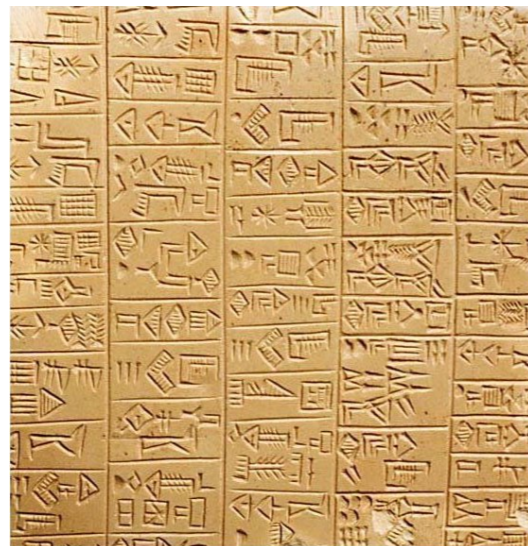
Chauvet cave  
proto-writing

~20,000 years ago



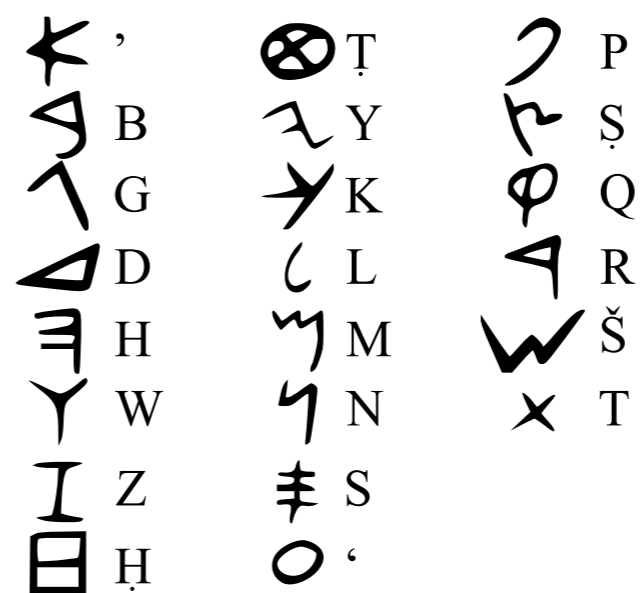
Sumerian cuneiform  
logographic

~5,000 years ago



Phoenician abjad  
predecessor of alphabet

~3,000 years ago



Latin letters

~2,500 years ago

**ABCDEF**  
**GHIJKLM**  
**NOPQRS**  
**TUVWXYZ**

abstraction

# Text

- Features of Text as representation language
  - abstract
  - general for mental concepts
  - different across population groups (countries, accents, religions,...)
  - linear perception
  - semi-structured (content: grammar, words, sentences, paragraphs,.. ; appearance: typography, calligraphy,..)
  - Legibility !!!!!

What is the challenge with Text?

Why Text Vis?

## 1.1 Text Visualization

A serious introduction to text visualization has to state that it is not a complete one. Why? When starting to work in the field, researchers are already confronted with the main problem itself, a large collection of documents covering many different aspects related to the subject text. Psychological research e.g. investigates perception and cognition of letters, the psychology of spoken and written language, or the psychology of reading. Linguistics describe inter alia models on language structure, language function, language features, etymology, and linguistic transformations. While both disciplines already fill books and would require introductions by themselves, we did so far not mention visual appearance (typography) or evolution of sign systems. As practical approach, we limit this introduction to key aspects in development of text and text visualizations taking the historic tour (Section 1.1.1), describing psychological backgrounds (Section 1.1.2), and describe landmarks in text visualization (Section 1.1.3). As further simplification we consider written text to stem from an alphabetic system.

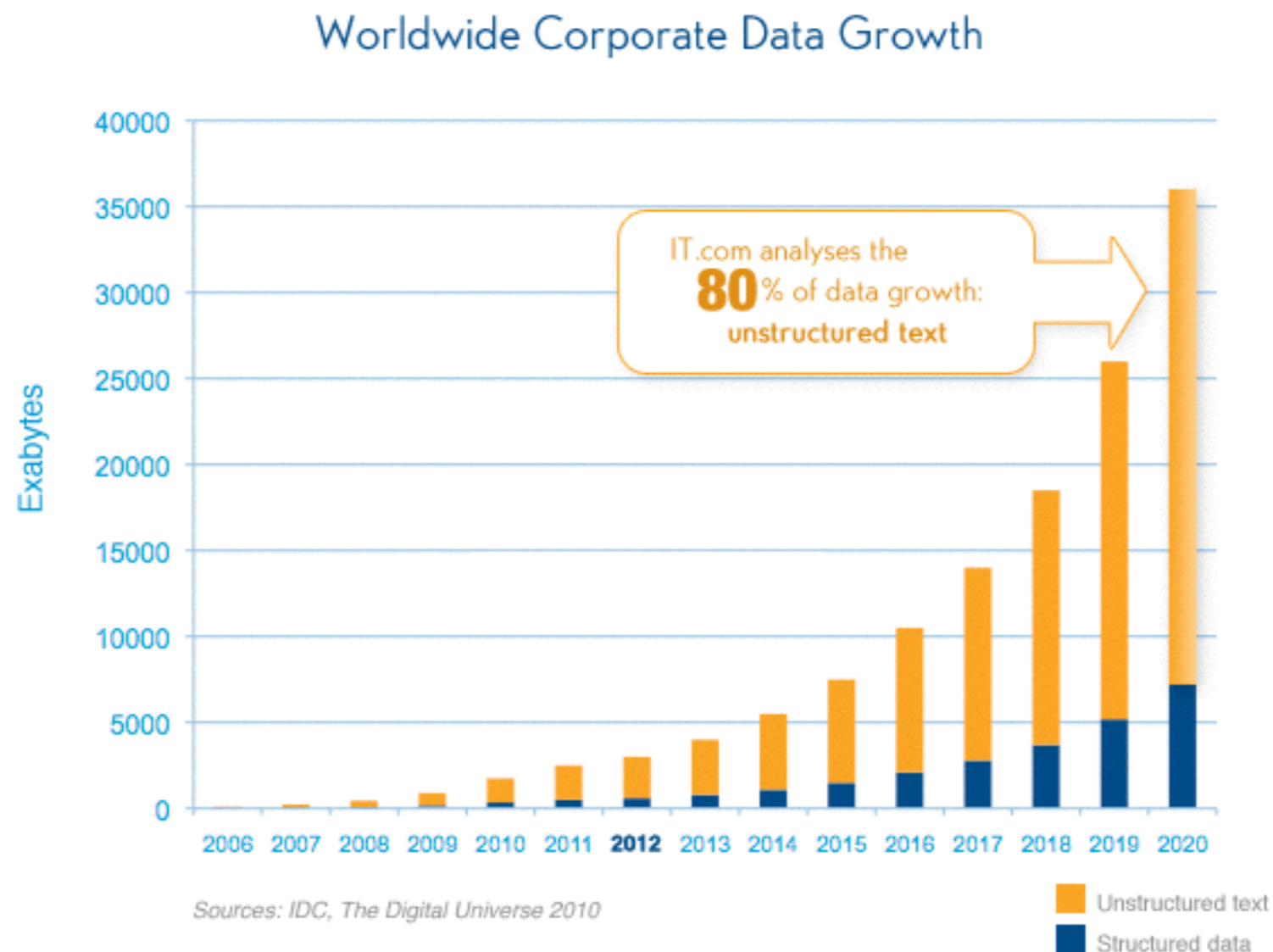
### 1.1.1 The historic trail

This section relies widely recommendable for further reading. Early humans started representing image and text representation into logographic form. Visual elements (semantics) within a language system included 24 signs for writing on papyrus vs. with hieroglyphs to an alphabet. Phoenicians have been the first known only-map their ordered set of letters. In Europe, Romans became dominant (1st century) and the medieval period developed during the 8th century. The impact on page style and decoration. The industrial revolution was invented. The successors were computers with word-processor and document distribution.

### 1.1.2 The psychological approach

We already discovered that text is nowadays as rapidly produceable and distributable as never before, but we did not throw light on how humans "consume" text. Schönplüg & Schönplüg [SS95] and Rayner & Pollatsek [RP94] provide extensive details on the psychological processes involved in reading which we summarize in this Section.

The consumption of text can be mainly split into reading as the perceptual part and understanding as the cognitive part. For reading, the human visual system performs saccadic eye movement processing lines of text. Each saccade takes on average 20 to 35 ms to bridge a range of 7 to 9



both references are

The paintings from Chauvet "covered" (Sadier et al. [SDB] story. Divergence between form evolved from pictographic elements of meaning. Their significance elements. Their significance, like the ease of development from the earliest developed 3,000 years ago. Their abjad is excessively, the Greek named

the times of Charlemagne (8th century) while printing was already allowed fast reproduction. The printing press or Schnörkel remained as a tool. Printing machines were used for content creation. Personal production of document production

# Text/Document Visualization

(focused on alphabetical languages)

- Text as Vis
- Vis for Text Documents
- Vis for large Text/Document Corpora
  - for exploring data with visualizations
  - to investigate specific properties
- Text in Vis
- TextVis Specials

# Text as Vis

- Typography:
  - typefaces (serif, sans-serif, **bold**, *italic*)
  - point size (10pt, 12pt, 24pt, 36pt.. ) - nowadays: 1/72 inch
  - line length (alignment: left, right, justified)
  - vertical: line spacing (leading)
  - horizontal: spaces between groups of letters (tracking)
  - space between pairs of letters (kerning)
  - combining letters to a glyph ligatures

*No kerning*  
A V W a

*Kerning applied*  
A V W a

fi → fi  
fl → fl

ß

# Text as Vis

- Creating a font type is an art which requires profound design knowledge
- .. or it can be a science:

Scientists have developed a way to carve shapes from DNA canvases, including all the letters of the Roman alphabet, emoticons and an eagle's head.

Bryan Wei, a postdoctoral scholar at Harvard Medical School in Boston, Massachusetts, and his colleagues make these shapes out of single strands of DNA just 42 letters long. Each strand is unique, and folds to form a rectangular tile. When mixed, neighbouring tiles stick to each other in a brick-wall pattern, and shorter boundary tiles lock the edges in place. [...]



<http://www.nature.com/news/dna-drawing-with-an-old-twist-1.10742>

# Text as Vis

- Typesetting:
  - letterpress printing
  - Linotype machine
  - digital printing/copying (typewheel, dot-matrix, inkjet, laser)
  - digital text (resolution is key: **s m a l l** -> retina)
- Encoding text for electronic devices:
  - mapping each character to a sequence of bytes
  - Universal Character Set (UTF-[**8**,16,32]) fonts
  - exchange of typeset documents: PostScript and PDF

# Text as Vis

- rules of thumb:
  - limit the use of fonts to only a few typefaces !!
  - use “special” fonts only when appropriate
  - a good resource for fonts in web projects are google fonts

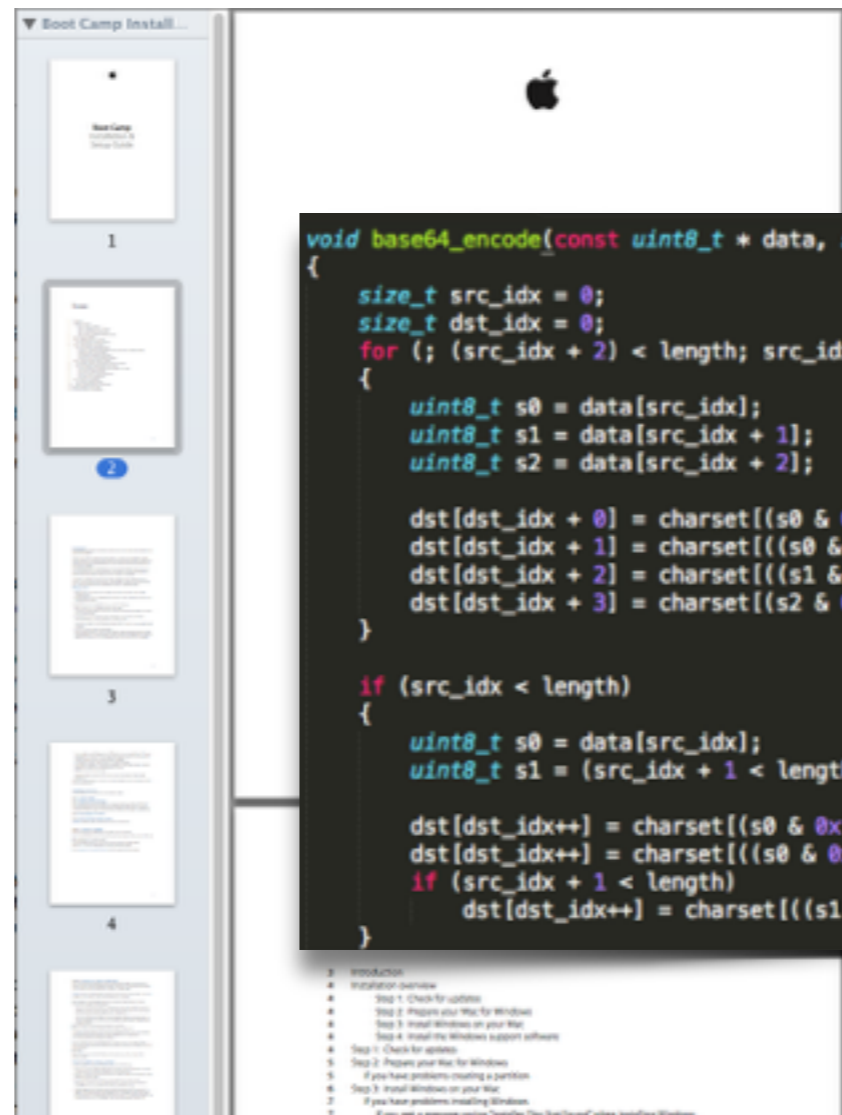


# Visualization for “Raw” Text

- in daily use..

enriched text - hypertext  
linking (graph navigation)

overview & detail

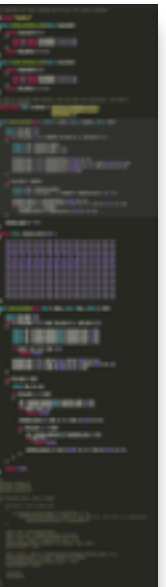


```
void base64_encode(const uint8_t * data, size_t length, char * dst)
{
    size_t src_idx = 0;
    size_t dst_idx = 0;
    for (; (src_idx + 2) < length; src_idx += 3, dst_idx += 4)
    {
        uint8_t s0 = data[src_idx];
        uint8_t s1 = data[src_idx + 1];
        uint8_t s2 = data[src_idx + 2];

        dst[dst_idx + 0] = charset[(s0 & 0xfc) >> 2];
        dst[dst_idx + 1] = charset[((s0 & 0x03) << 4) | ((s1 & 0xf0) >> 4)];
        dst[dst_idx + 2] = charset[((s1 & 0x0f) << 2) | (s2 & 0xc0) >> 6];
        dst[dst_idx + 3] = charset[(s2 & 0x3f)];
    }

    if (src_idx < length)
    {
        uint8_t s0 = data[src_idx];
        uint8_t s1 = (src_idx + 1 < length) ? data[src_idx + 1] : 0;

        dst[dst_idx++] = charset[(s0 & 0xfc) >> 2];
        dst[dst_idx++] = charset[((s0 & 0x03) << 4) | ((s1 & 0xf0) >> 4)];
        if (src_idx + 1 < length)
            dst[dst_idx++] = charset[((s1 & 0x0f) << 2)];
    }
}
```



# Visualization for “Raw” Text

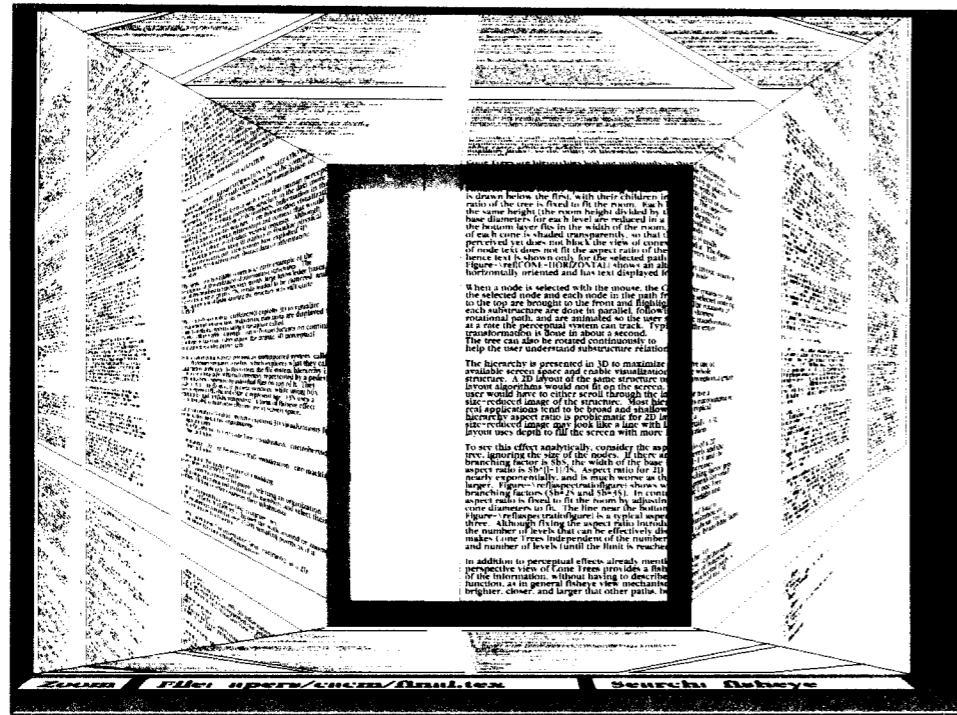


Figure 3: Document Lens with lens pulled toward the user. The resulting truncated pyramid makes text near the lens' edges readable.

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.

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 previews. 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.

So the user 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 disambiguation of keywords by their context. For instance, it

about “user” or “user interface” keyword “user” would

The technique we present to create the thumbnails is a general distortion technique for document content that high-

to a user defined interest

The global structure of a page, namely the position of images and columns, is preserved. An example is shown in Figure 1. In the keyword search application, an interest function

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

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 of the original thumbnail. Woodruff et al. [WRM02] pre-

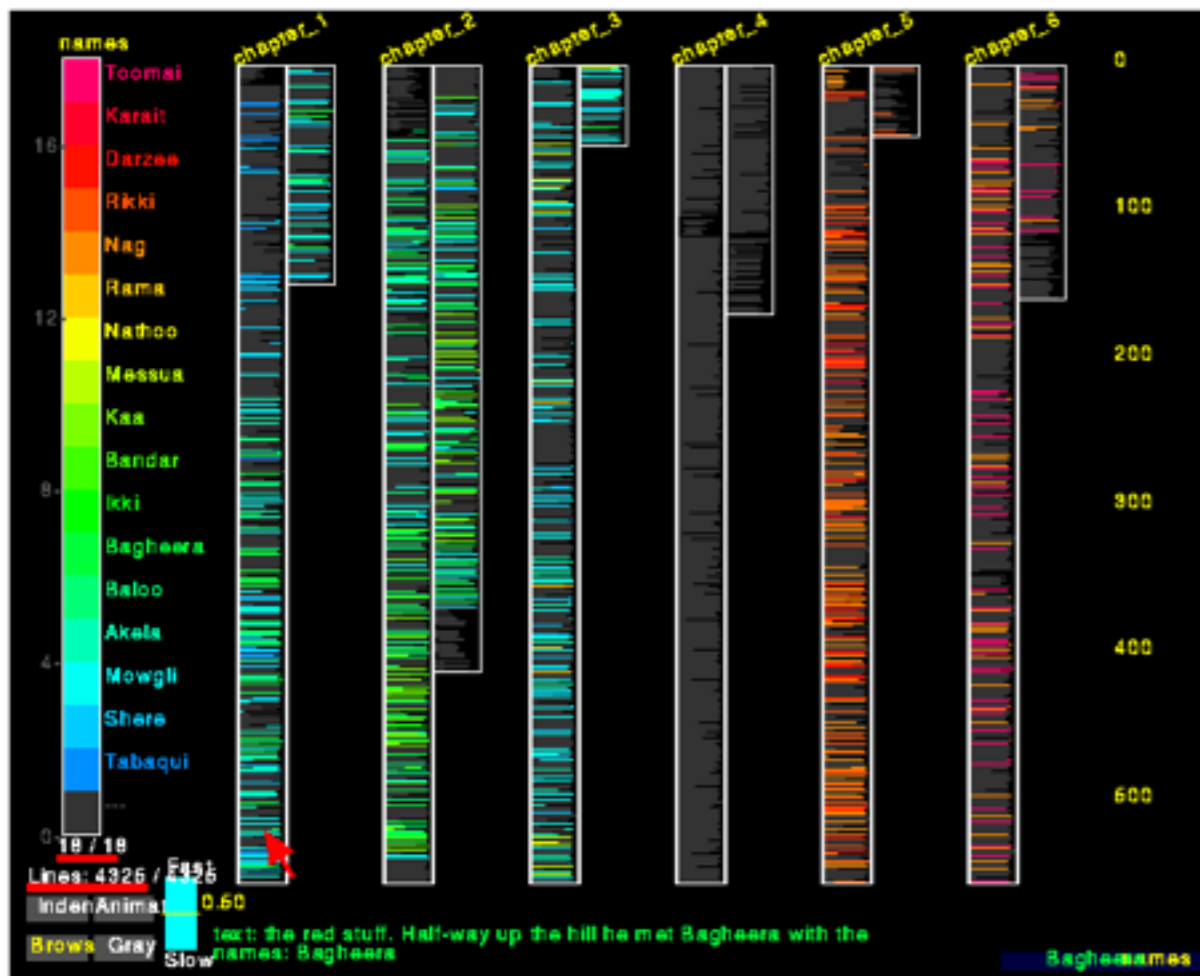
submitted to Eurographics Conference on Visualization (EuroVis) (2012)

## Document Thumbnails with Variable Text Scaling

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

*Computer Graphics Forum, volume 31 issue 3 pp.*

# Visualization for “Raw” Text

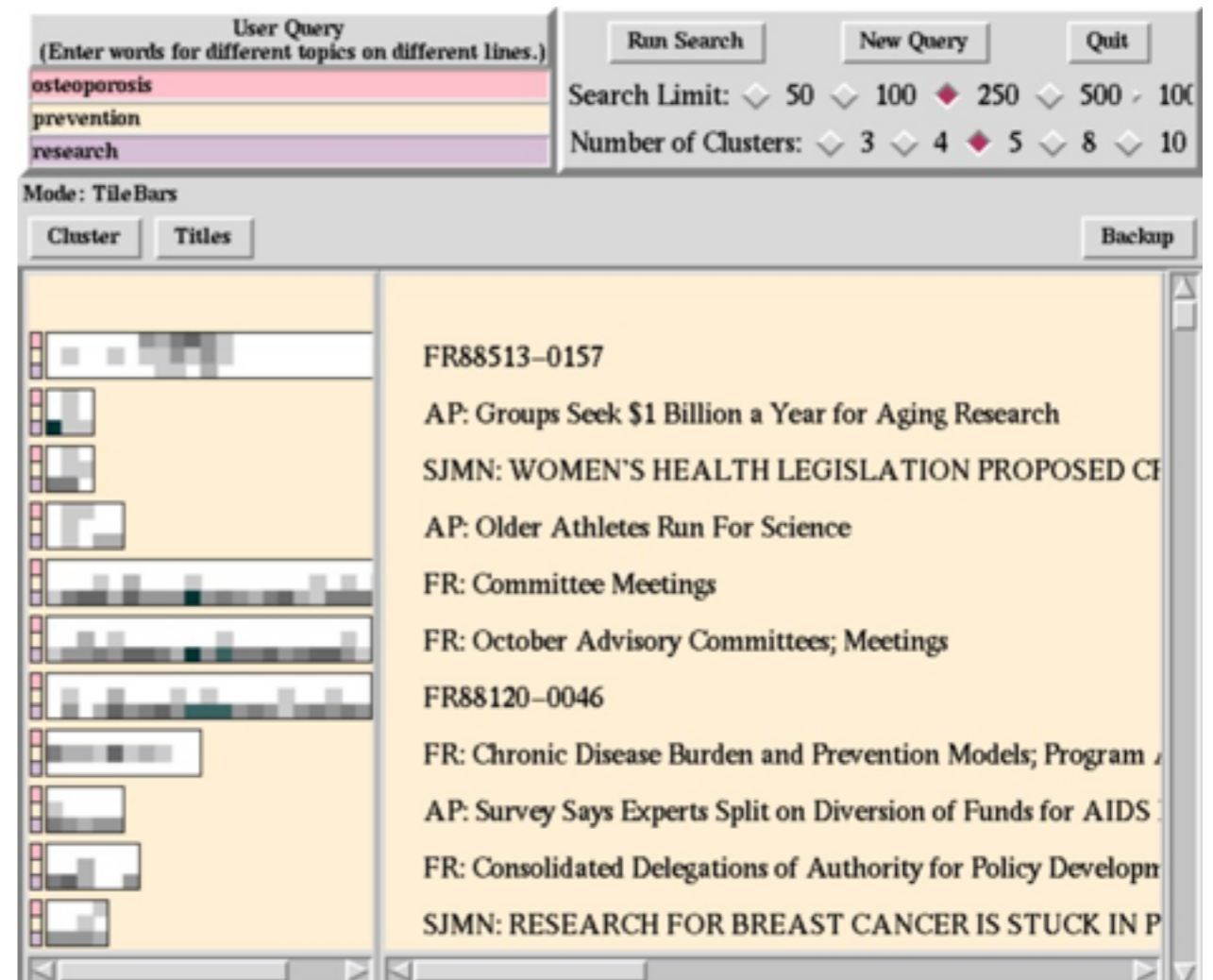


SeeSoft

Stephen G. Eick.

**Graphically displaying text.**

*Journal of Computational and Graphical Statistics*, 3(2):127-142, June 1994.



**TileBars: Visualization of Term Distribution Information in Full Text**

Marti Hearst

*Information Access, Proceedings of the ACM SIGCHI Conference on Human Factors in Computing Systems (CHI), Denver, CO, 1995*

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

# Structured Text Features

- simple counts
- or a bag of words (used for similarity measures):

	princess	dragon	castle
doc1	1	1	1
doc2	0	0	1

# Typical Steps of Processing to derive Text Features

- Large collections require pre-processing of text to extract information and align text. Typical steps are:
  - cleaning (regular expressions)
  - sentence splitting
  - change to lower case
  - stopword removal (most frequent words in a language)
  - stemming - demo porter stemmer
  - POS tagging (part of speech) - demo
  - noun chunking
  - NER (name entity recognition) - demo opencalais
  - deep parsing - try to “understand” text.

# Sample Text

KIEV, Ukraine — Struggling to reach a deal to form a new majority coalition in Parliament, and under excruciating pressure because of a looming economic disaster, the Ukrainian lawmakers temporarily running the country on Tuesday delayed until Thursday the naming of an acting prime minister and a provisional government.

The delay underscored the extreme difficulty that lawmakers now face in rebuilding the collapsed government left behind when President Viktor F. Yanukovich fled Kiev on Saturday and was removed from power in a vote supported by some members of his own party.

The three main opposition parties, which share little in common politically, have been in fierce negotiations, not just among themselves, but also with civic activists and other groups representing the many constituencies involved in Ukraine's three months of civic uprising.

Arseniy P. Yatsenyuk, the leader in Parliament of the Fatherland Party and a leading contender to serve as acting prime minister, pleaded with colleagues to swiftly reach an agreement on the designation of an interim government, which is needed to formally request emergency economic assistance from the International Monetary Fund.

# Text features are complicated

- Be aware!! text understanding can be hard:
  - *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?*

# Was that irony? - Noooo

Profanity sucks. (14)

Be more or less specific. (15)

Analogies in writing are like feathers on a snake. (19)

*excerpt from Rules of Writing*

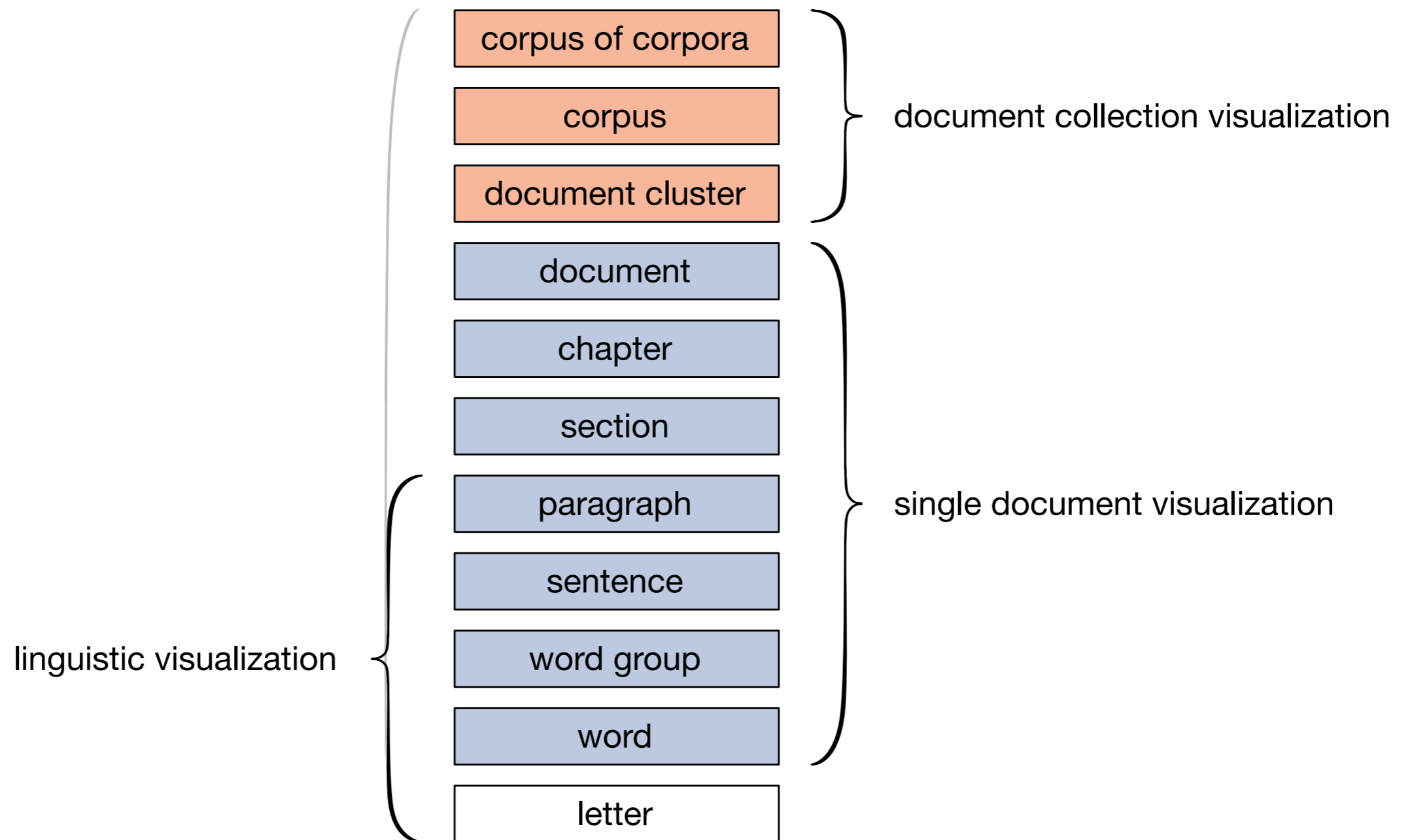
*by Frank L. Visco (June 1986 in Writers' digest)*

# Thinking about..

- or a bag of words (used for similarity measures):

	princess	dragon	castle
doc1	1	1	1
doc2	0	0	1

# Text Units Hierarchy

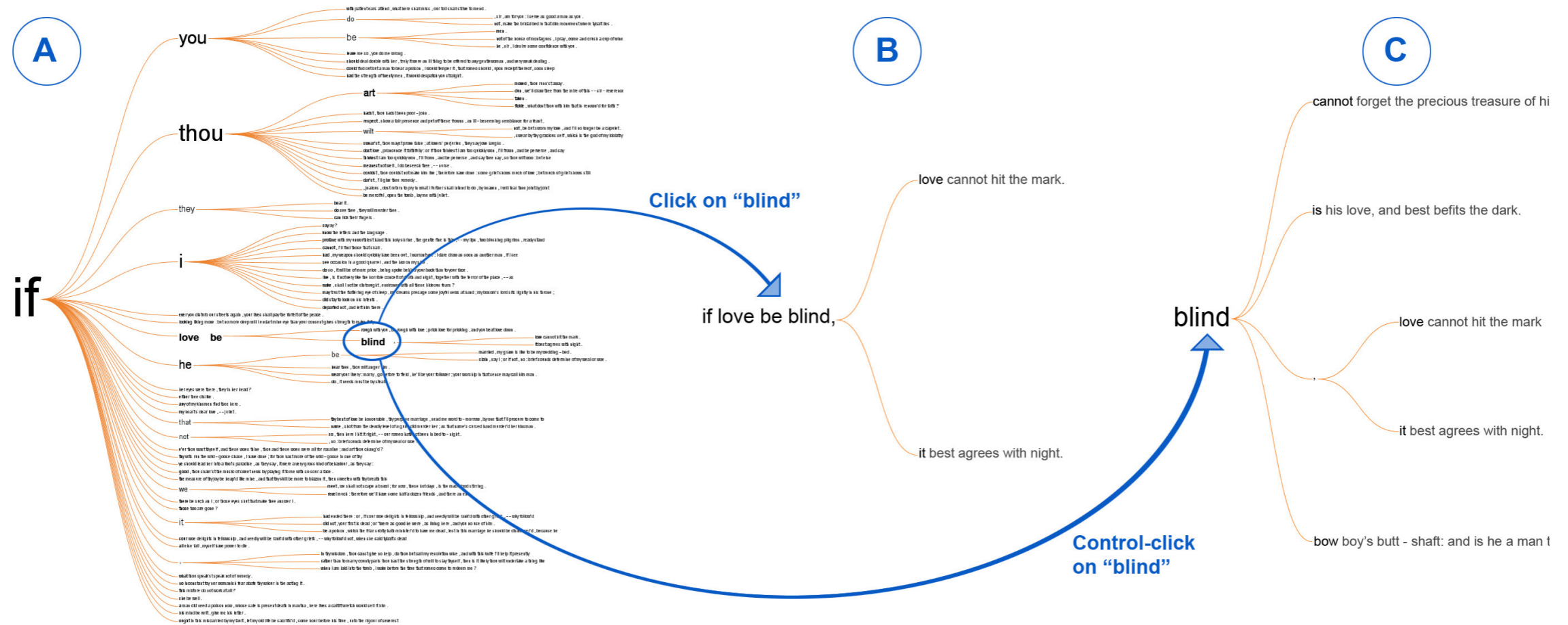


# Vis for Text Documents

- TagClouds : <http://www.flickr.com/photos/tags/>
- WordCloud (popular) — <http://www.wordle.net>



# Vis for Text Documents



## The word tree, an interactive visual concordance

M Wattenberg, FB Viégas

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

<http://www.bobdylan.com/us/songs/blowin-wind>

# Vis for Text Documents

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

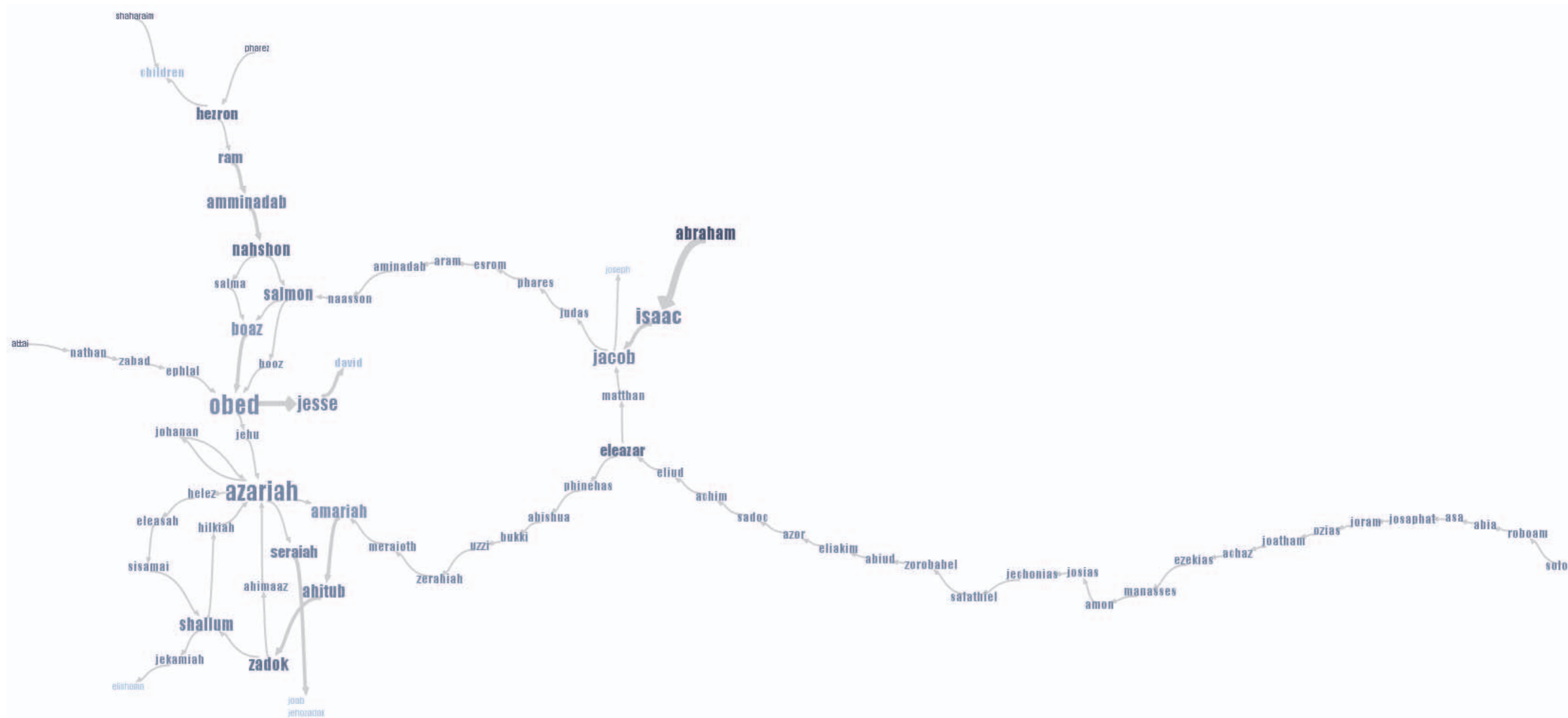
*Many Eyes creates the word graph:*

**pride** → **impertinence**

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

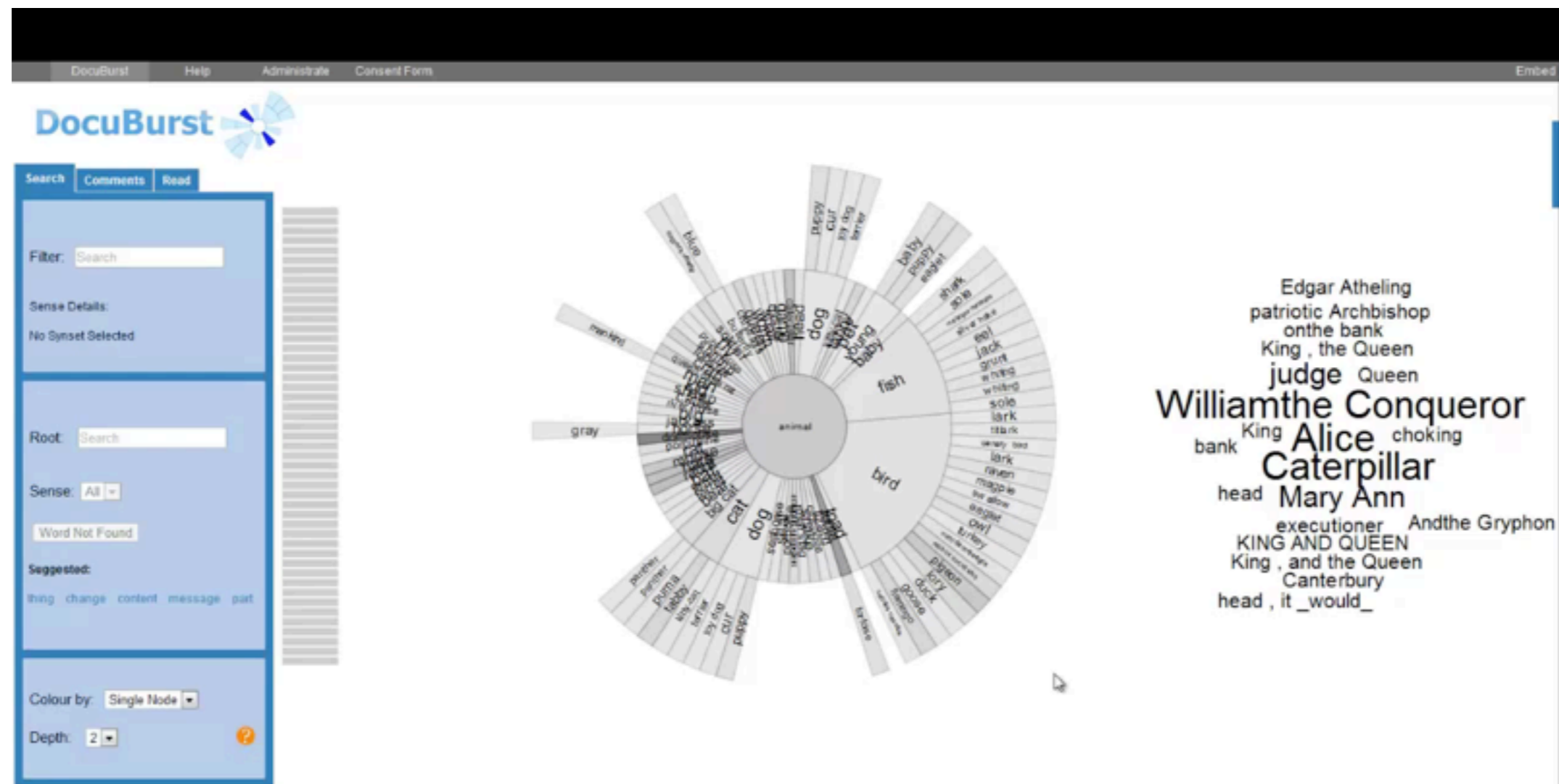
**Mapping Text with Phrase Nets.**

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

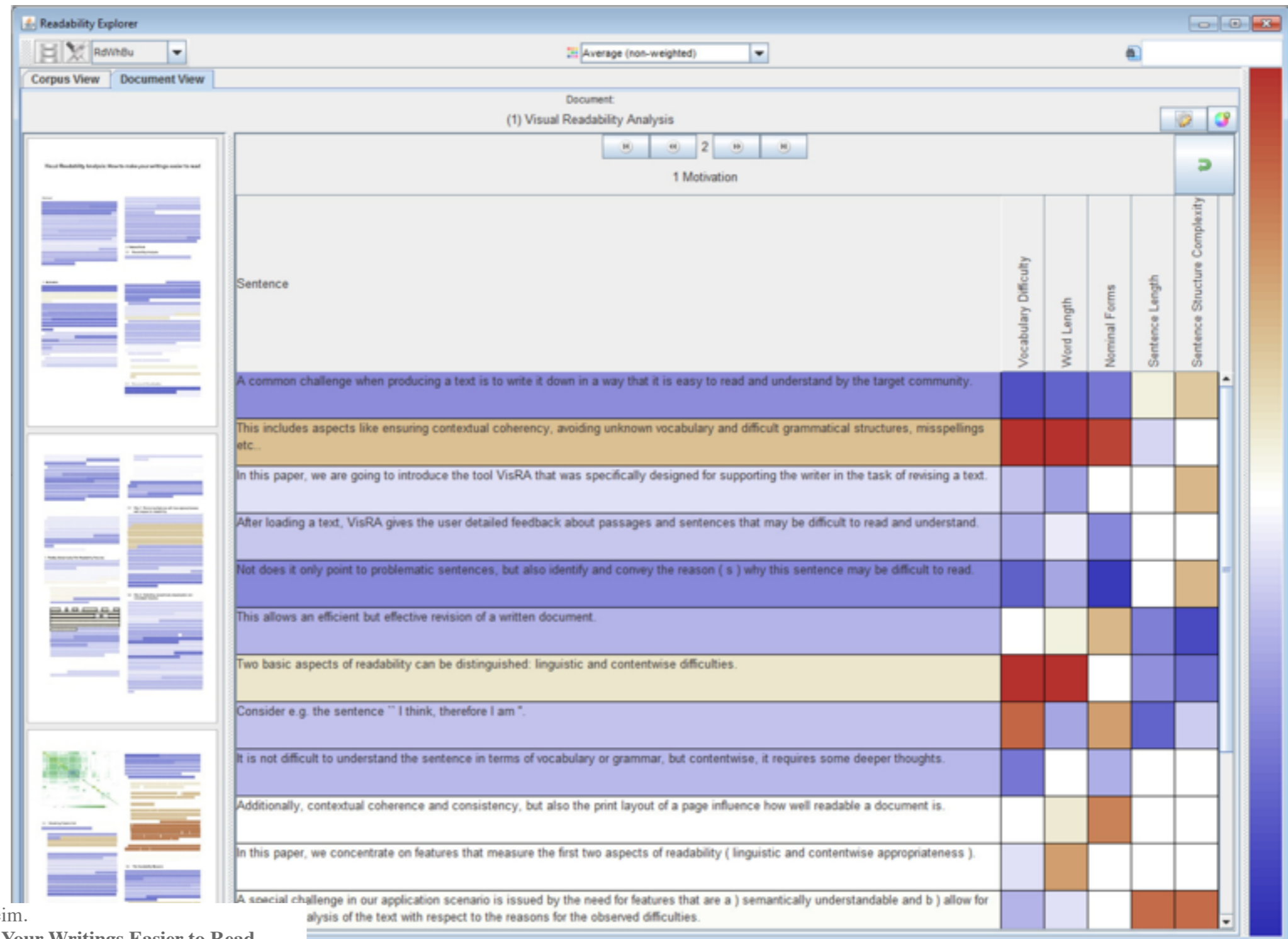


# Vis for Text Documents

- DocuBurst : <http://vialab.science.uoit.ca/docuburst/>
- based on: WordNet, see the network



# Vis for Language Analysis



D. Oelke, D. Spretke, A. Stoffel and D. A. Keim.

**Visual Readability Analysis: How to Make Your Writings Easier to Read.**

*IEEE Transactions on Visualization and Computer Graphics*, 18(5):662-674, 2012.

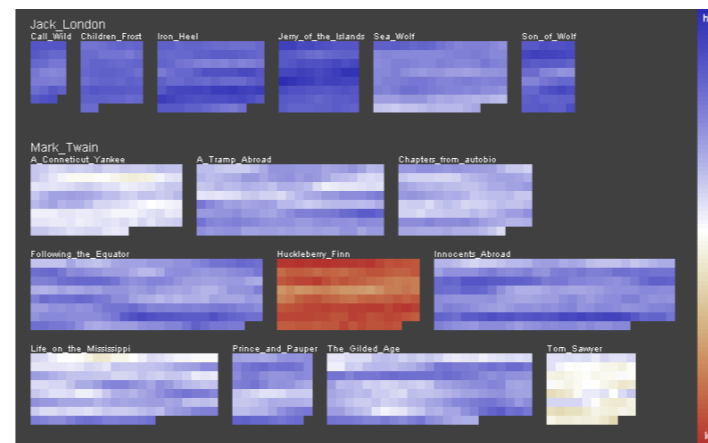
# Vis for Language Analysis

- Literature fingerprints:

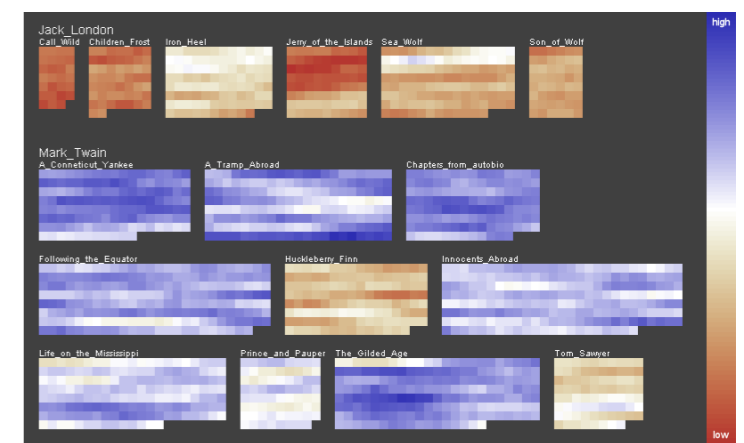
“Fingerprints of books of Mark Twain and Jack London. Different measures for authorship attribution are tested. If a measure is able to discriminate between the two authors, the visualizations of the books that are written by the same author will equal each other more than the visualizations of books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

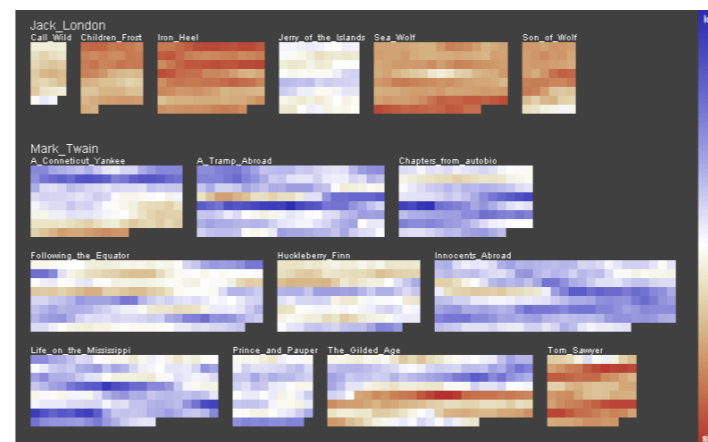
\*method to measure the vocabulary richness



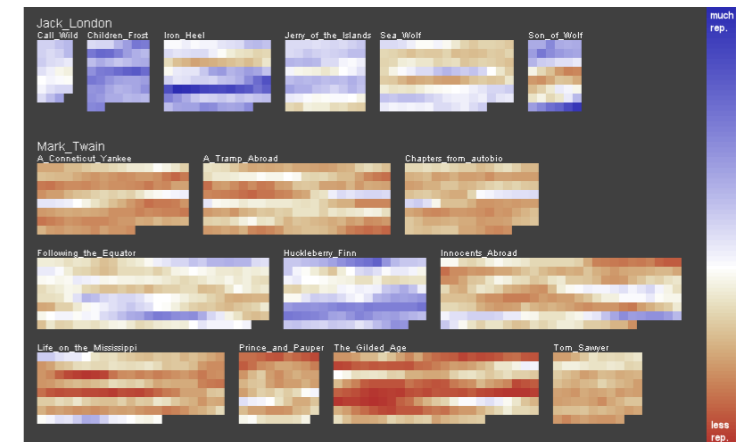
(a) Function words (First Dimension after PCA)



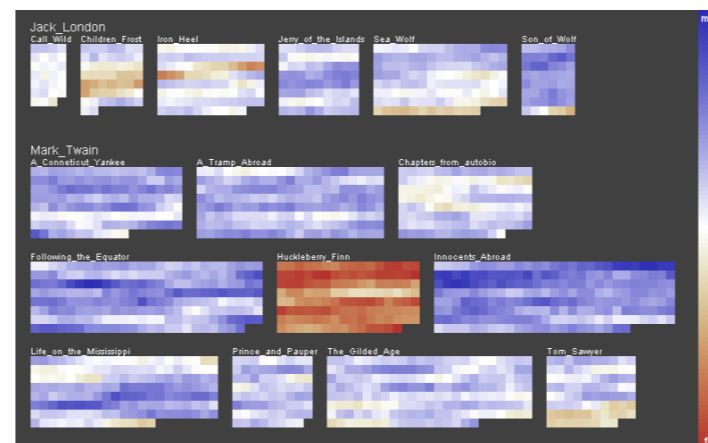
(b) Function words (Second Dimension after PCA)



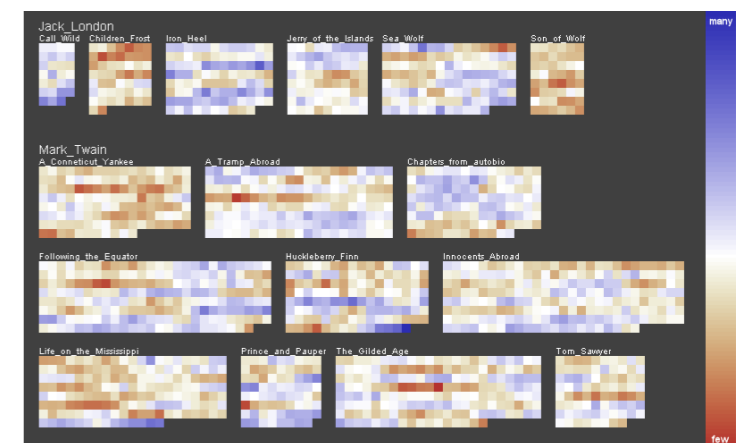
(c) Average sentence length



(d) Simpson's Index



(e) Hapax Legomena



(f) Hapax Dislegomena

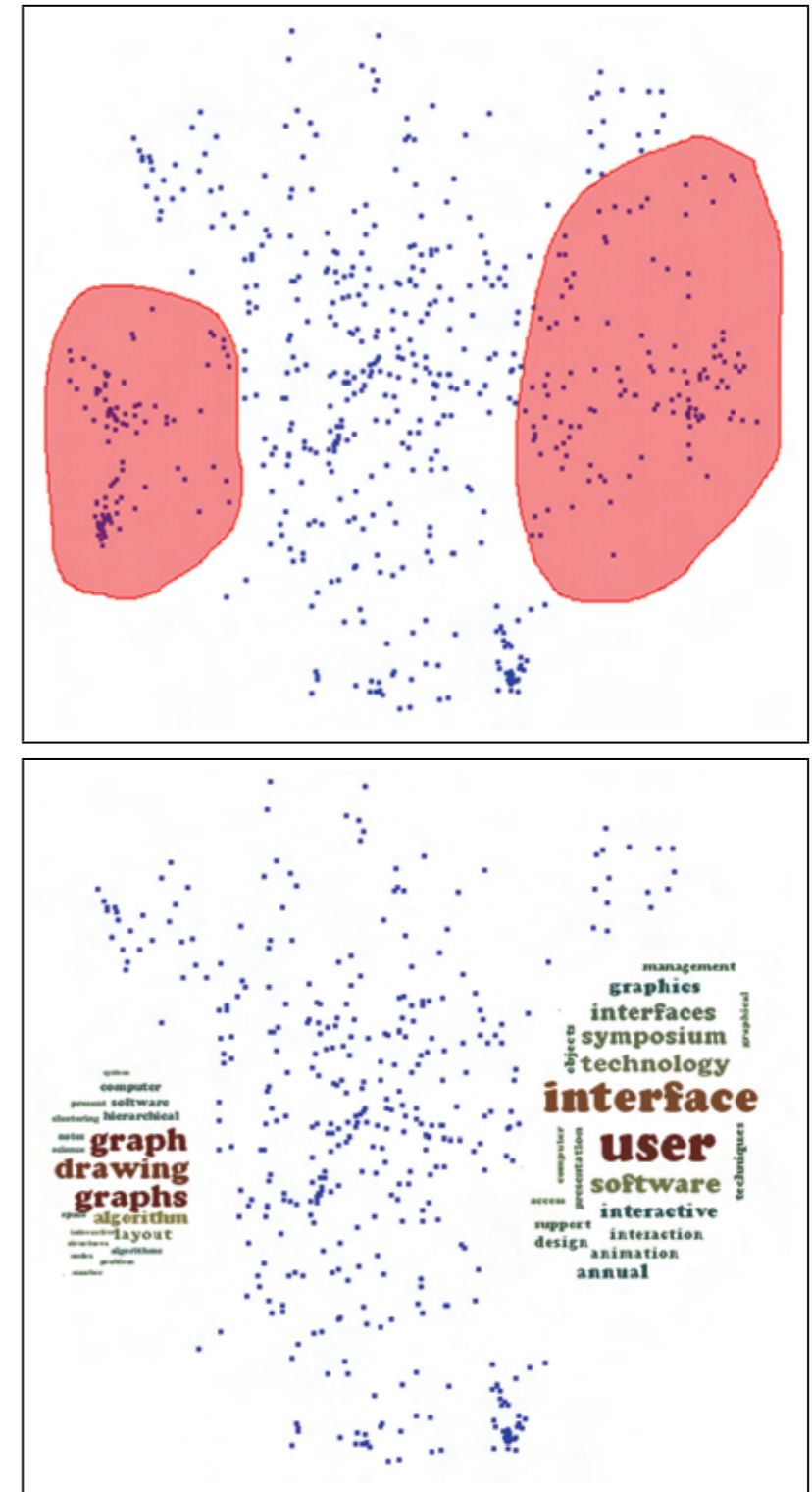
Daniel A. Keim and Daniela Oelke.

**Literature Fingerprinting: A New Method for Visual Literary Analysis.**

*Proceedings of the 2007 IEEE Symposium on Visual Analytics Science and Technology (VAST '07)*

# Visualization for Large Text Corpora

- use bag-of-words to project documents w.r.t. text similarity into a landscape
- (only) one example



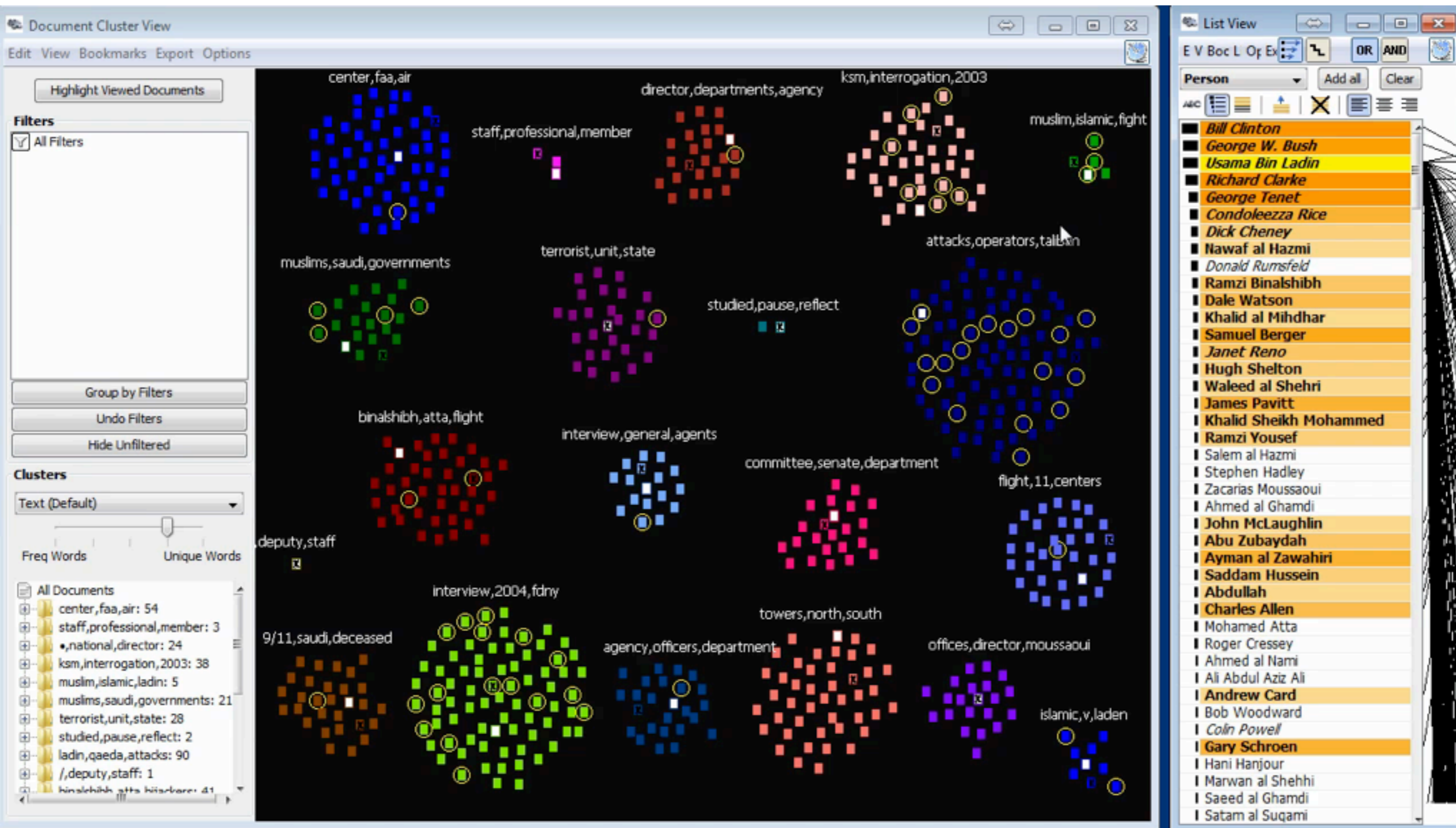
**Figure 5:** A user can interactively draw a region (polygon) containing a subset of documents of interest (top figure). Keywords are extracted from the selected document and their corresponding word cloud is built inside the user-defined region (bottom figure).

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)

# Visual Analytics for Large Text Corpora (example JigSaw)



# Vis for Large Document Collections

- documents contain more information than just text:

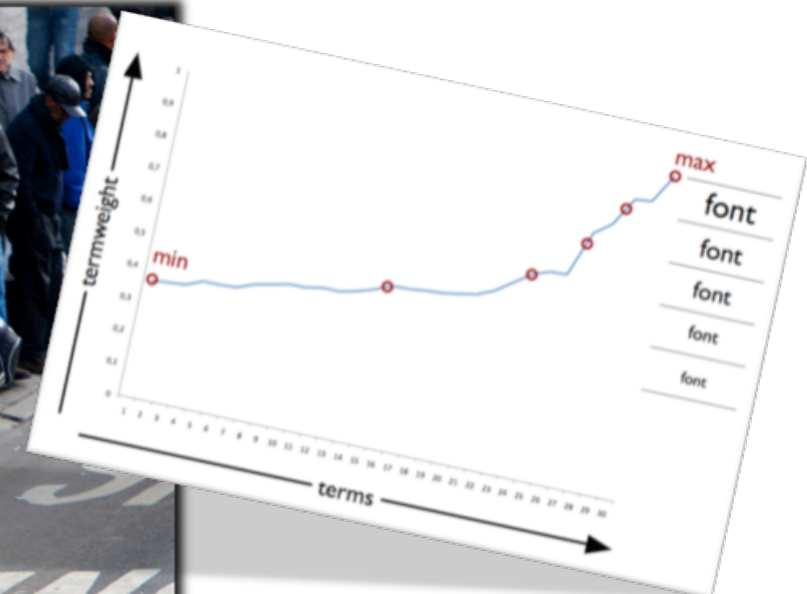
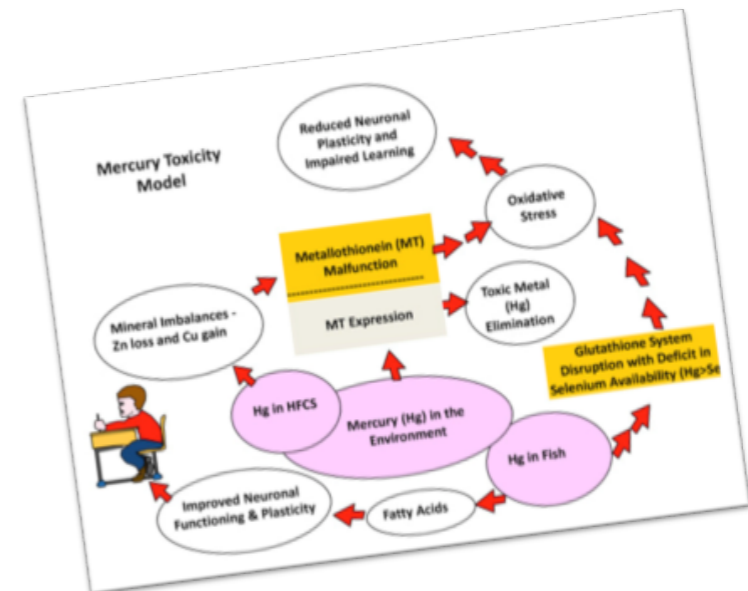
- meta information
- structure (paragraphs, text boxes,..)
- figurative content:

- parallel perception

- compact

- multi-lingual

- empathy



# Vis for Large Document Collections

- (only) three examples:
  - Bohemian bookshelf
  - DocumentCards
  - Semantics:

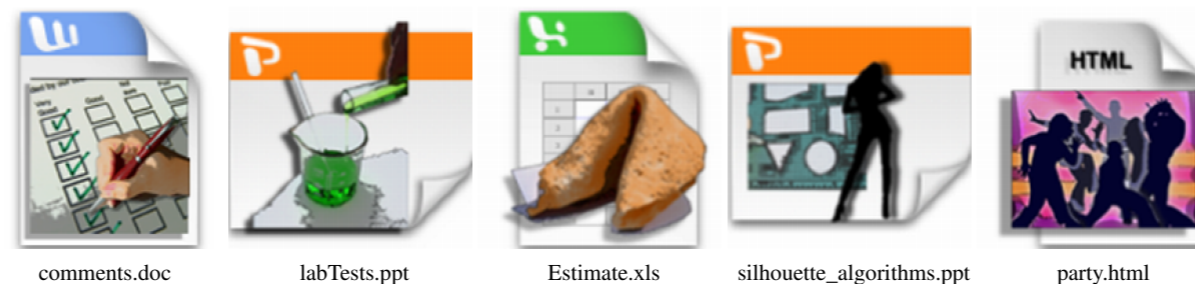


Figure 1: Semantics generated by our system for various filenames.

## Semantics: Visual Metaphors as File Icons

Vidya Setlur, Conrad Albrecht-Buehler, Amy A. Gooch,  
Sam Rossoff, Bruce Gooch

# Vis for Large Document Collections



[webpage with video](#)

Alice Thudt, Uta Hinrichs and Sheelagh Carpendale.

**The Bohemian Bookshelf: Supporting Serendipitous Book Discoveries through Information Visualization.**

*CHI '12: Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, 2012*

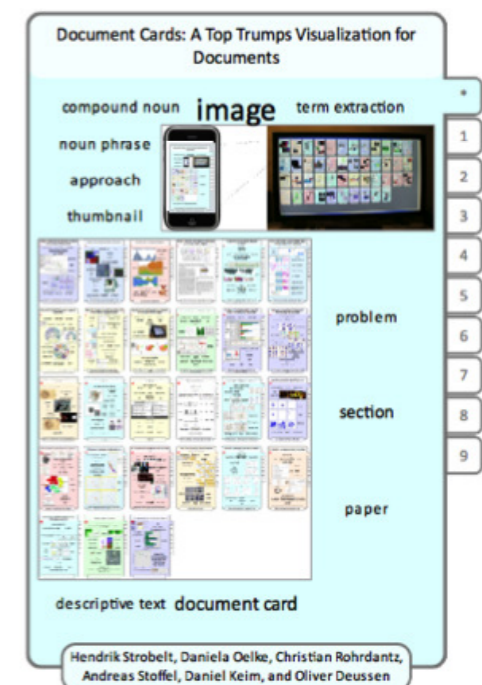
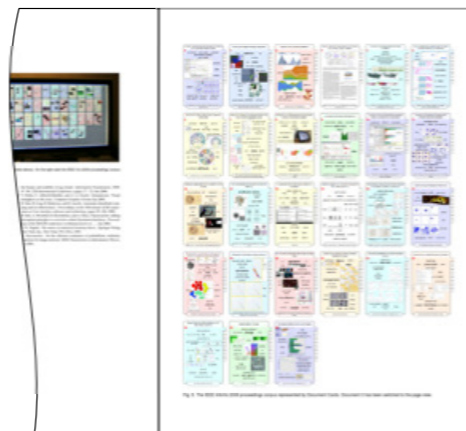
# DocumentCards

- summarize scientific documents using important terms and important figures
- design considerations:
  - Document Cards are fixed size thumbnails that are self-explanatory
  - Document Cards represent the document's content as a mixture of figure and textual representatives
  - Document Cards should be discriminative and should have a high recognizability

# DocumentCards

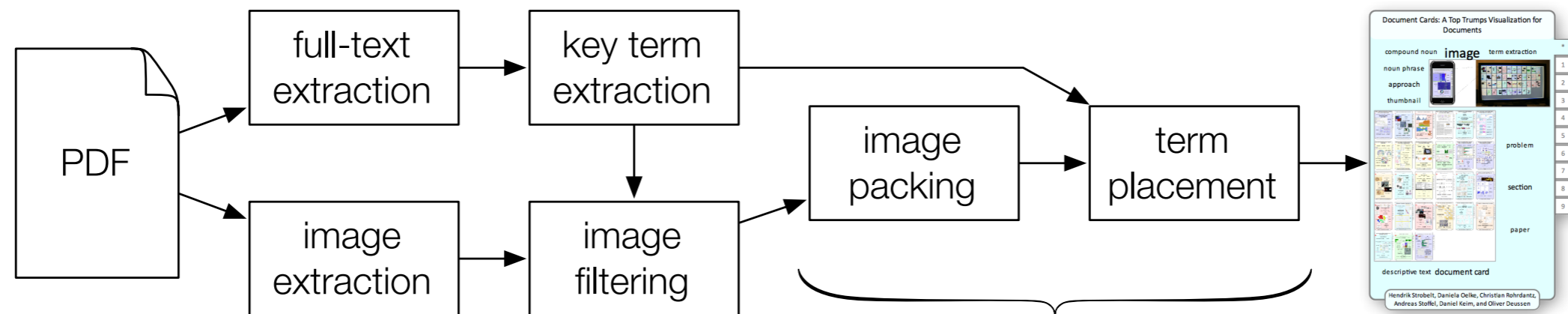


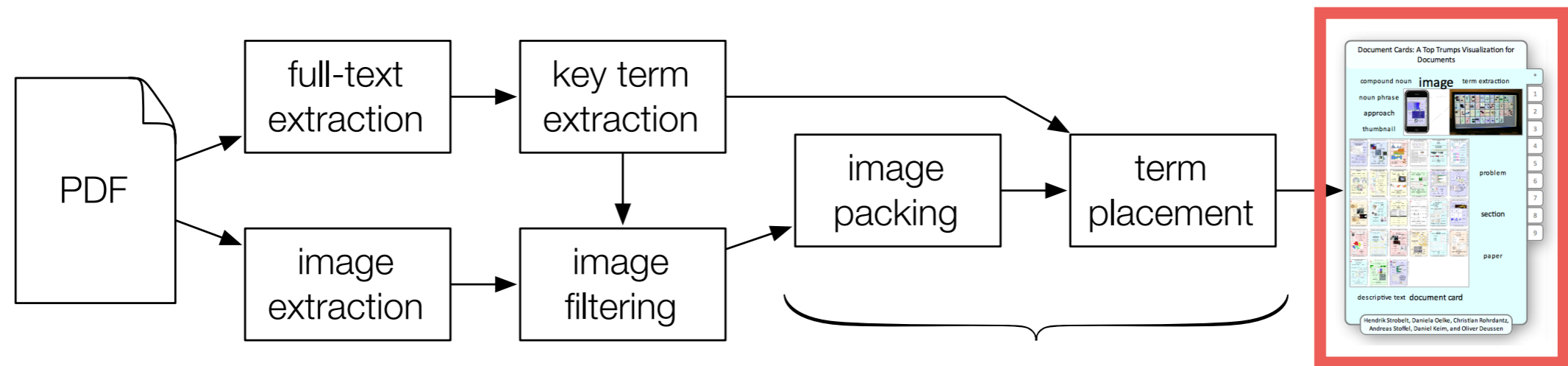
...





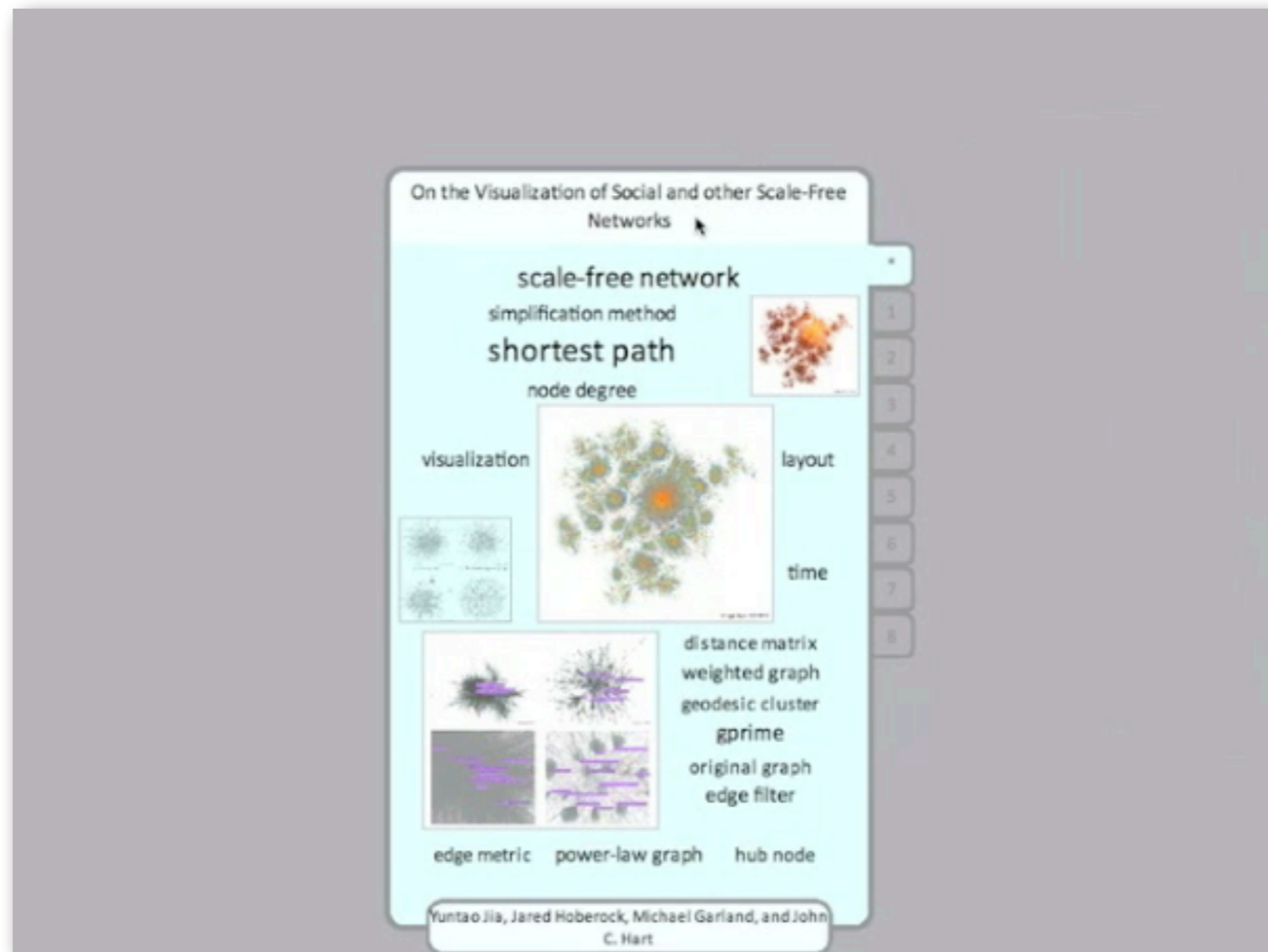
# DC - pipeline



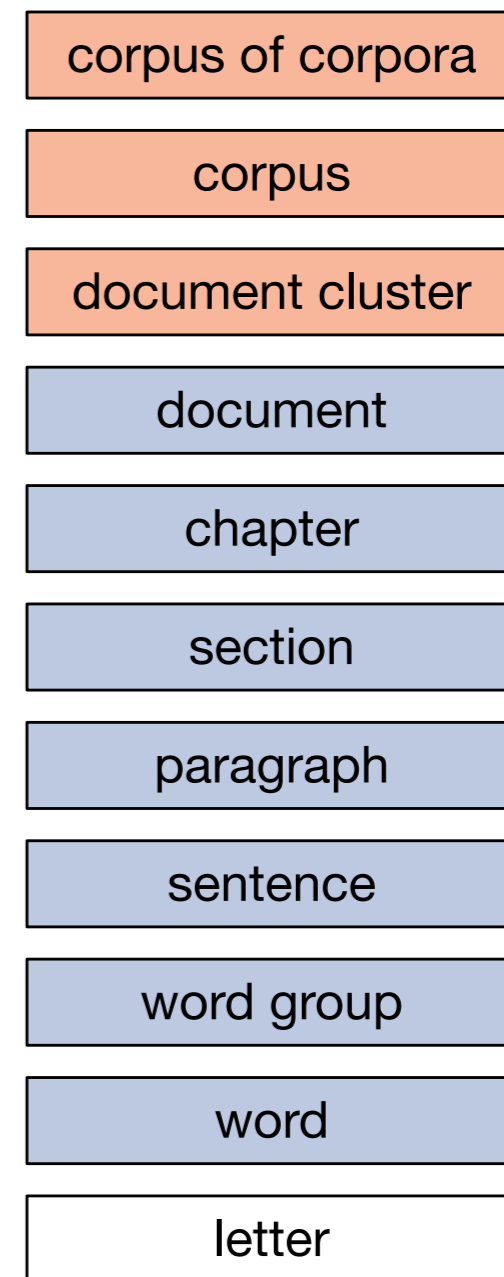
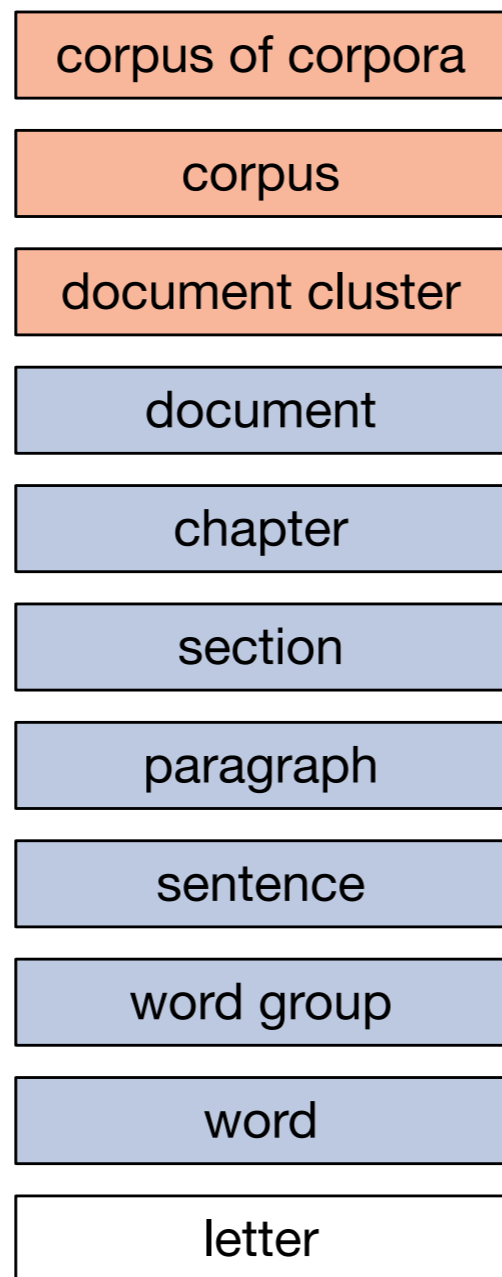
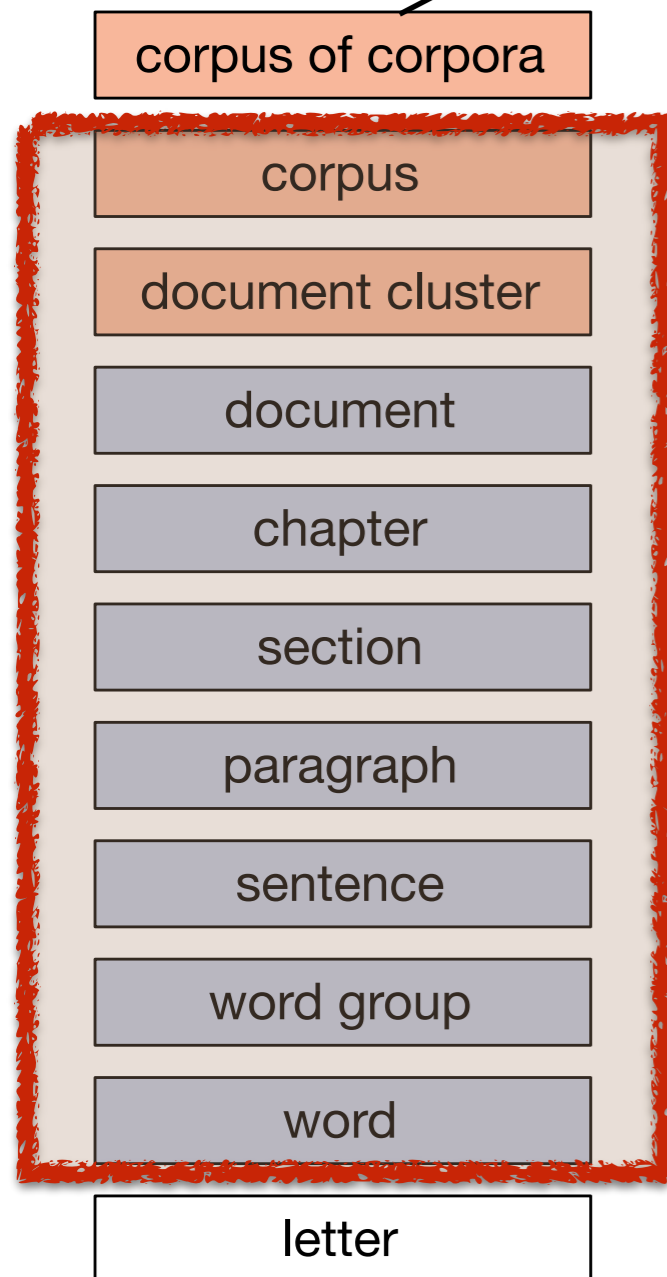


## Interaction:

- caption tooltip
- abstract tooltip
- move to orig. Pos.
- page switch
- term highlighting



DiTop



time

# Compare Corpora

- Compare topics between text collections

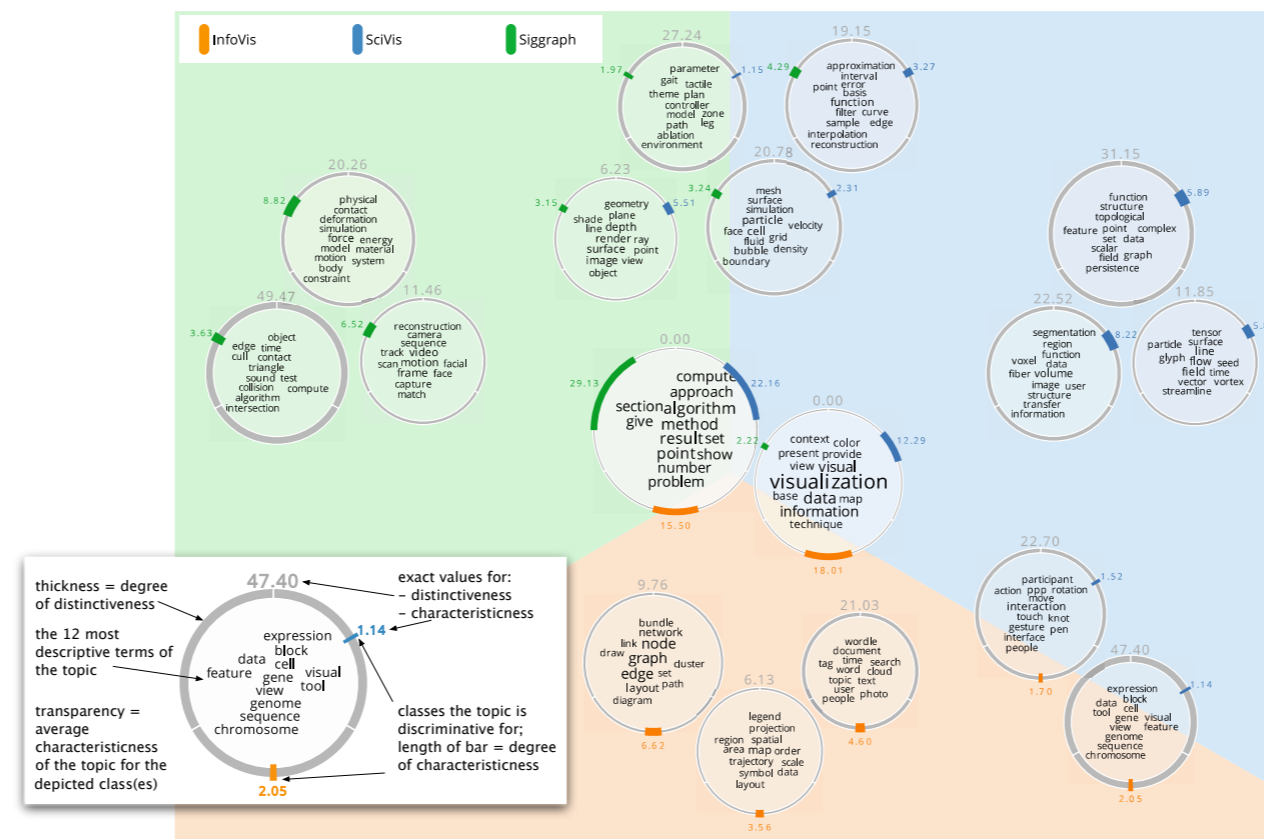
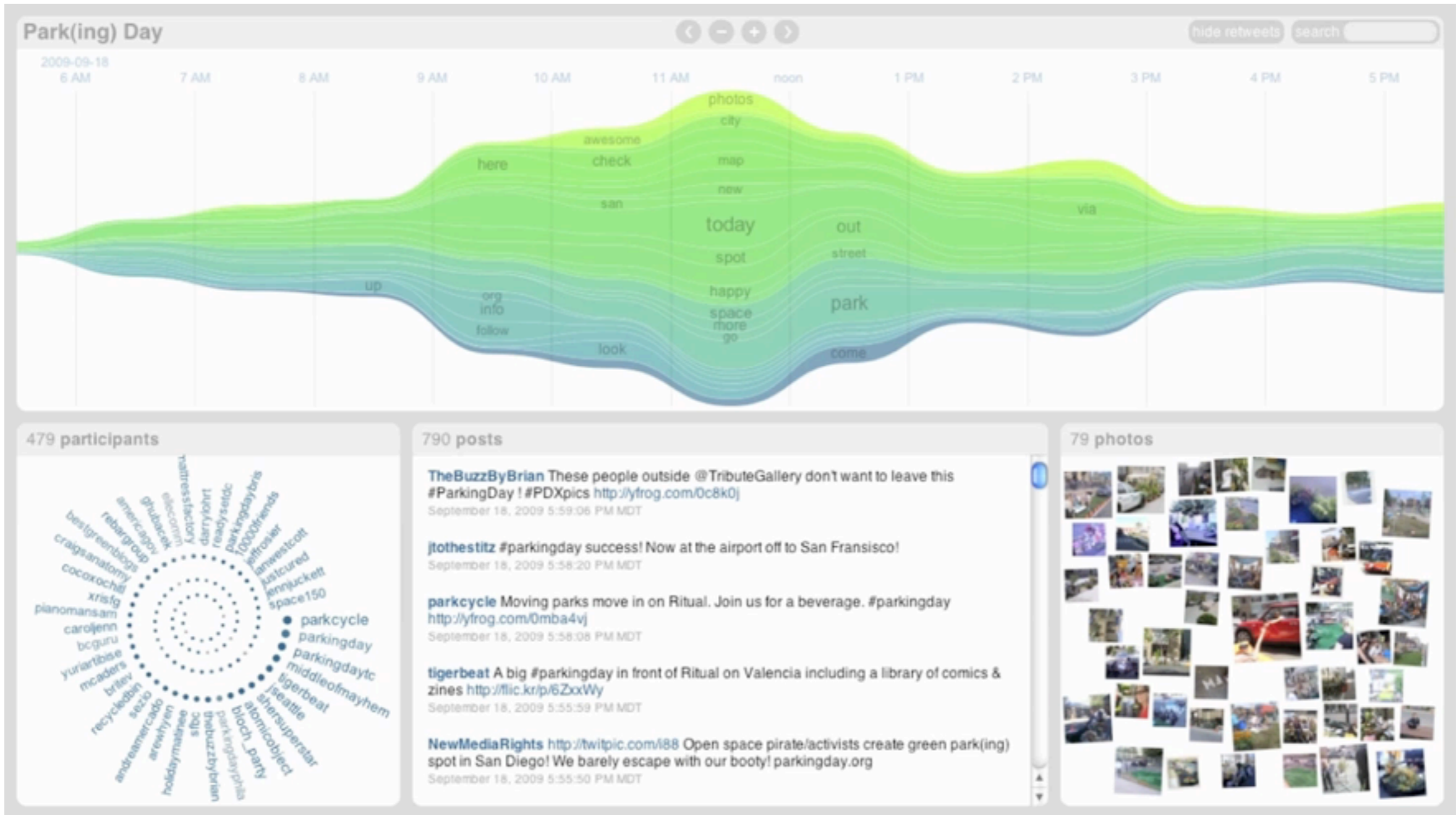


Figure 1: Comparison of 495 papers of InfoVis, SciVis, and Siggraph (discrimination threshold = 6, number of topics = 30)

# Vis for Time-Evolving Document Collections



# Vis for Time Evolving Texts

command strings.	of input in addition to or in place of typed command strings.	of input in addition to or in place of typed command strings.	of input in addition to or in place of typed command strings.	of input in addition to or in place of typed command strings.
Voice user interfaces, which accept input and provide output by generating voice prompts which are transmitted via a telephone network and heard by the user using a telephone. The user input is made by pressing keys or one keys.	Voice user interfaces, which accept input and provide output by generating voice prompts which are transmitted via a telephone network and heard by the user using a telephone. The user input is made by pressing keys or one keys.	Voice user interfaces, which accept input and provide output by generating voice prompts which are transmitted via a telephone network and heard by the user using a telephone. The user input is made by pressing keys or one keys.	Voice user interfaces, which accept input and provide output by generating voice prompts which are transmitted via a telephone network and heard by the user using a telephone. The user input is made by pressing keys or one keys.	Voice user interfaces, which accept input and provide output by generating voice prompts which are transmitted via a telephone network and heard by the user using a telephone. The user input is made by pressing keys or one keys.
Natural Language interfaces - Used for	Natural Language interfaces - Used for	Natural Language interfaces - Used for	Natural Language interfaces - Used for	Natural Language interfaces - Used for

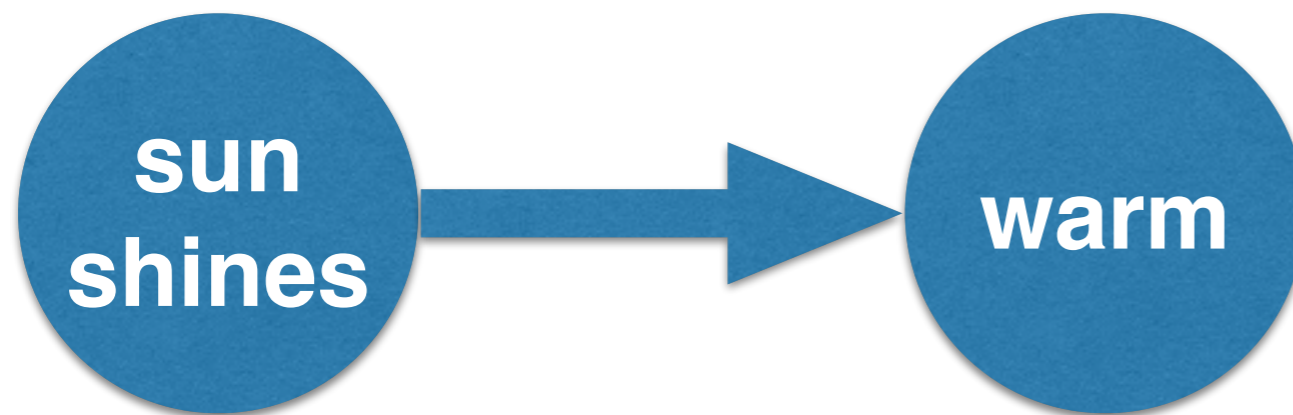
“This article examines the benefits of using text animated transitions for navigating in the revision history of textual documents. We propose an animation technique for smoothly transitioning between different text revisions, then present the Diffamation system. Diffamation supports rapid exploration of revision histories by combining text animated transitions with simple navigation and visualization tools. We finally describe a user study showing that smooth text animation allows users to track changes in the evolution of textual documents more effectively than flipping pages.”

**Video on the [webpage](#)**

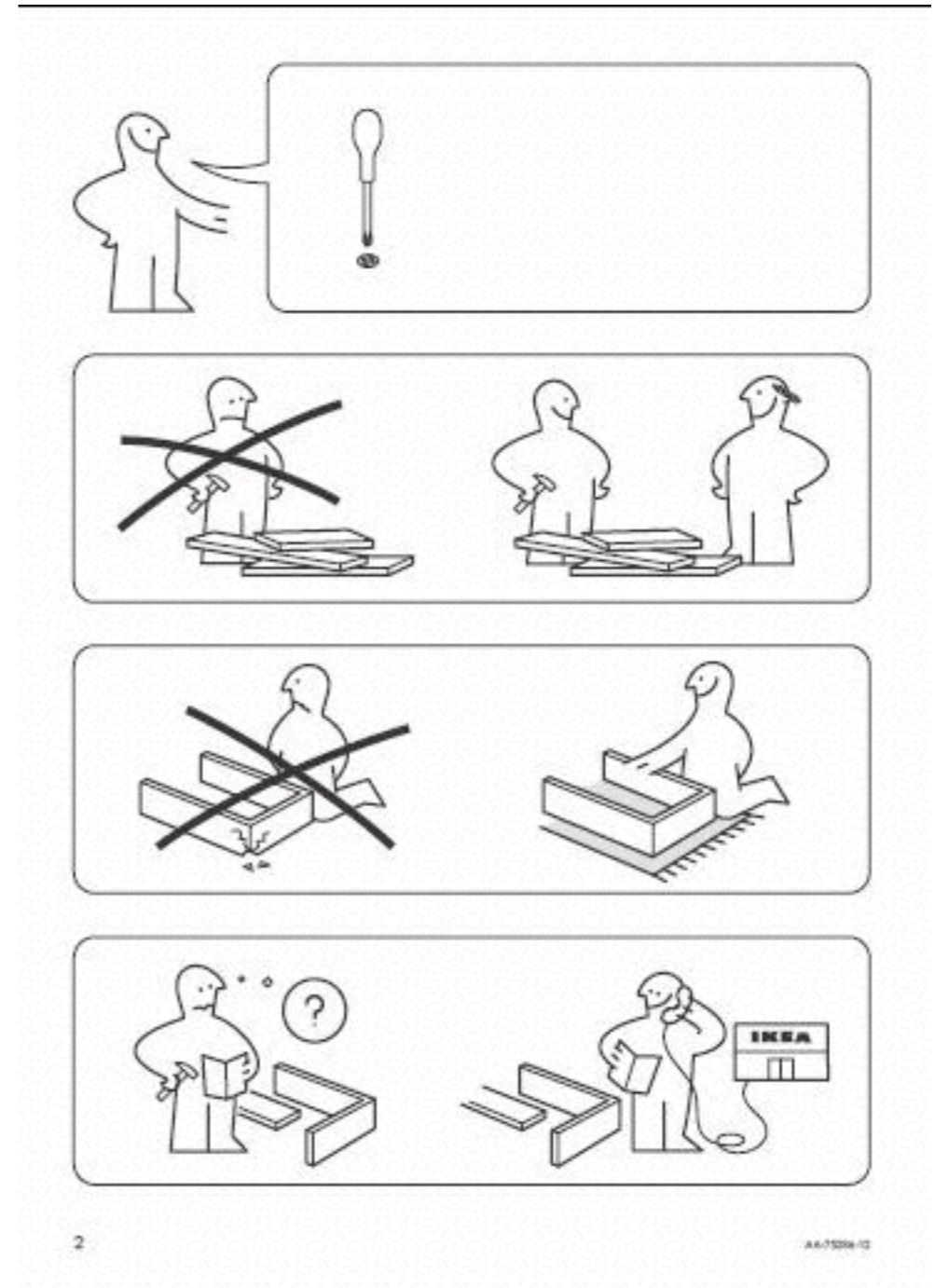
The Role of Text in Vis

# Text in Vis

- Non-Example: Ikea
- Labels:



- Map Legends



# Text in Vis Storytelling

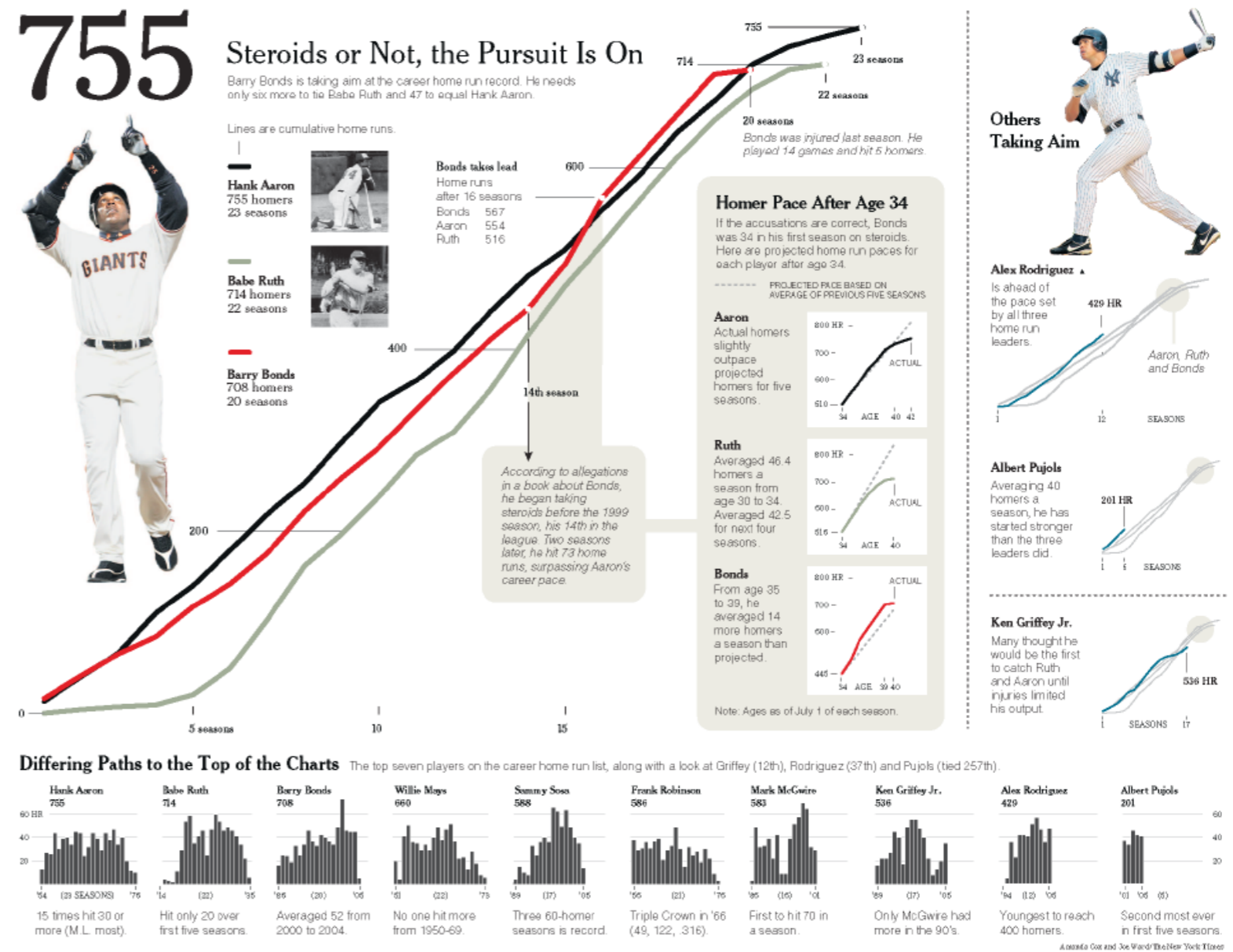


Fig. 1. Steroids Or Not, the Pursuit is On. New York Times.

**Narrative Visualization: Telling Stories with Data**

Edward Segel, Jeffrey Heer

IEEE Trans. Visualization & Comp. Graphics (Proc. InfoVis), 2010

# TextVis Specials

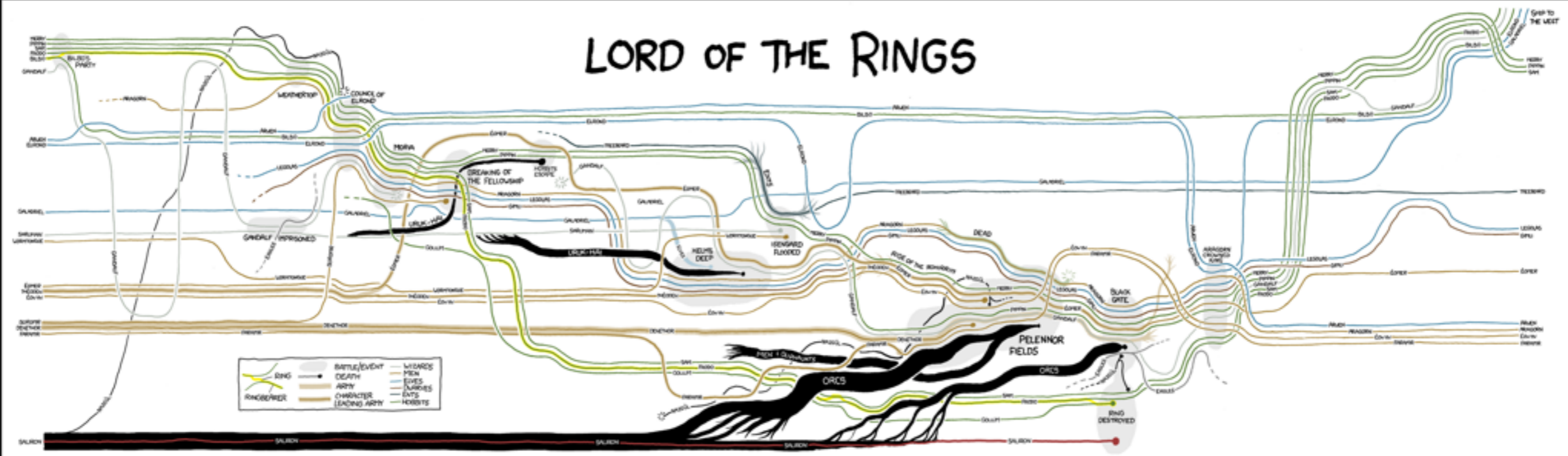
# Vis for Text Translation



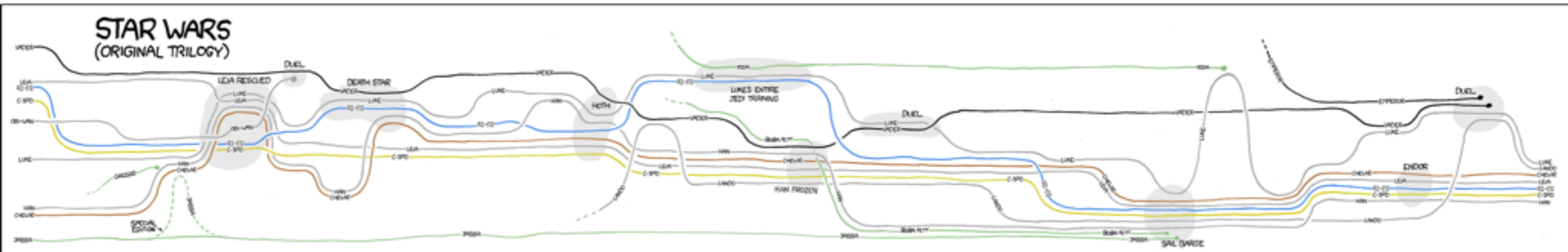
Figure 6: Translation lattice for the German sentence, “Hallo, ich bin gerade auf einer Konferenz im Nationalpark in Banff.” The statistically-identified best path (along the bottom) was incorrect and has been repaired. Photo nodes provide an alternative representation for words not in the translation vocabulary. Mouse over expands the node and reveals four photos, while other nodes move away to avoid occlusion.

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.

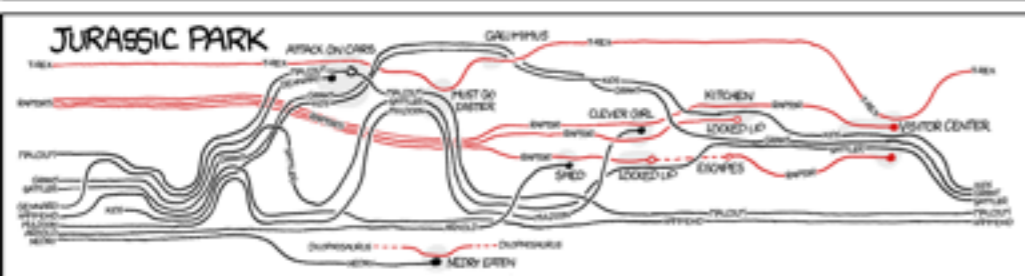
# LORD OF THE RINGS



STAR WARS  
(ORIGINAL TRILOGY)



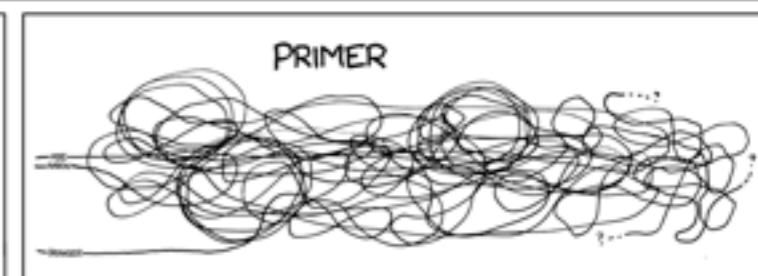
JURASSIC PARK



# 12 ANGRY MEN



## PRIMER



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

# Text to Vis conversion

“Natural language is an easy and effective medium for describing visual ideas and mental images. Thus, we foresee the emergence of language-based 3D scene generation systems to let ordinary users quickly create 3D scenes without having to learn special software, acquire artistic skills, or even touch a desktop window-oriented interface. WordsEye is such a system for automatically converting text into representative 3D scenes. WordsEye relies on a large database of 3D models and poses to depict entities and actions. Every 3D model can have associated shape displacements, spatial tags, and functional properties to be used in the depiction process.”



Figure 1: *John uses the crossbow. He rides the horse by the store. The store is under the large willow. The small allosaurus is in front of the horse. The dinosaur faces John. A gigantic teacup is in front of the store. The dinosaur is in front of the horse. The gigantic mushroom is in the teacup. The castle is to the right of the store.*

Bob Coyne and Richard Sproat. 2001.

**WordsEye: an automatic text-to-scene conversion system**

*Proceedings of the 28th annual conference on Computer graphics and interactive techniques (SIGGRAPH '01)*

# Further TextVis..

- ... on topic modeling
- ... for text exploration (human computer interaction)
- ... for search results
- ... linguistic features (e.g. vowel harmony)
- ... source code
- ... for sentiment analysis
- ... **SO MUCH MORE !!**

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

## Text Visualization Browser

A Visual Survey of Text Visualization Techniques

Provided by ISOVIS group

About Add entry

Techniques displayed:

141

Search:

Time filter:

1976

2014

Analytic Tasks



Visualization Tasks



Data

Source



Properties

