# UNIVERSITY OF ONTARIO INSTITUTE OF TECHNOLOGY

# Information Visualization:
# Text Visualization

## Dr. Christopher Collins

# TEXT VISUALIZATION

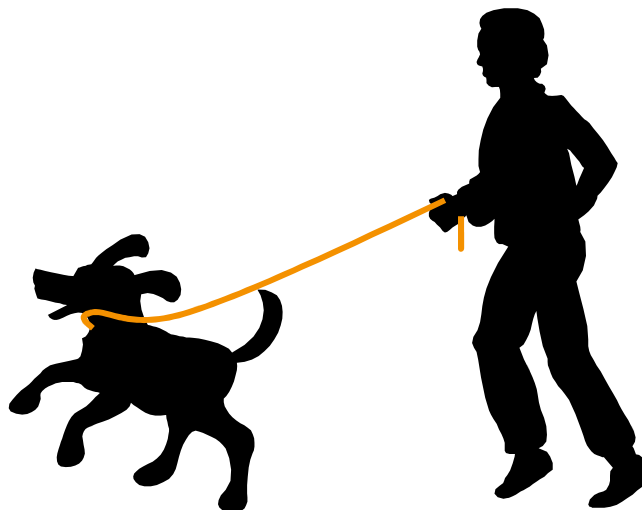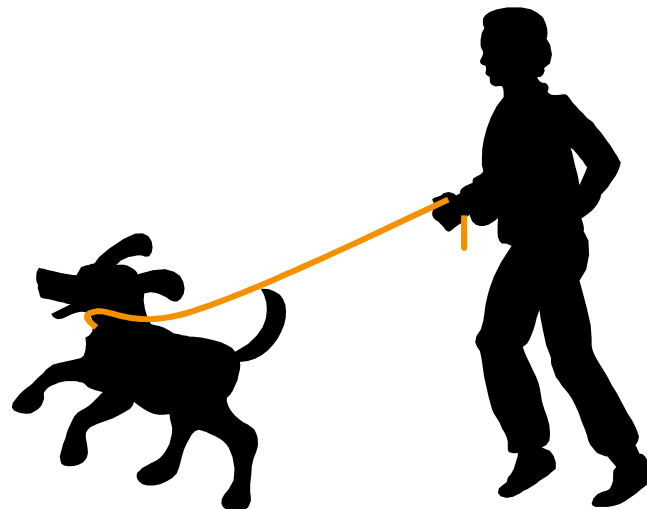# The dog.

# The excited dog .

# The man.

# The man walks.
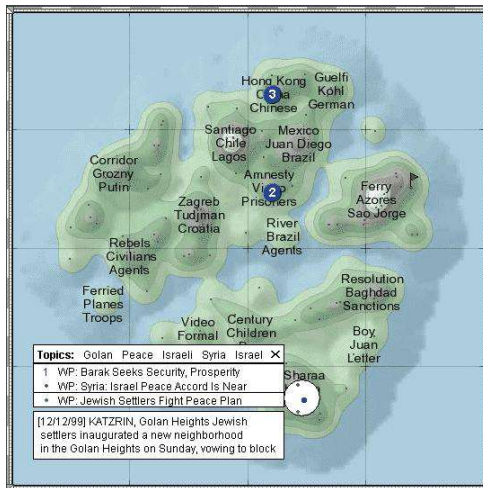
# The man walks the excited dog.

As the man walks the excited dog, he daydreams of the coming spring, and is filled with dread, as he is every year when the days drag on longer, the happy sun grinning sardonically at him as he enters his windowless workplace prison for the most hectic and stressful time of year.
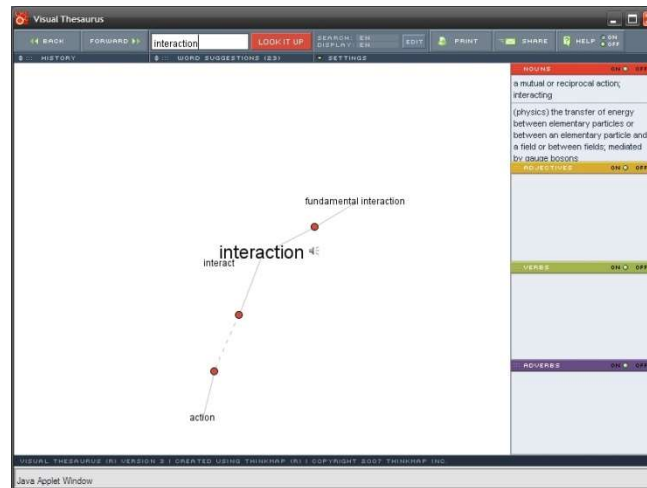
Example based on lecture notes of Marti Hearst, 2006
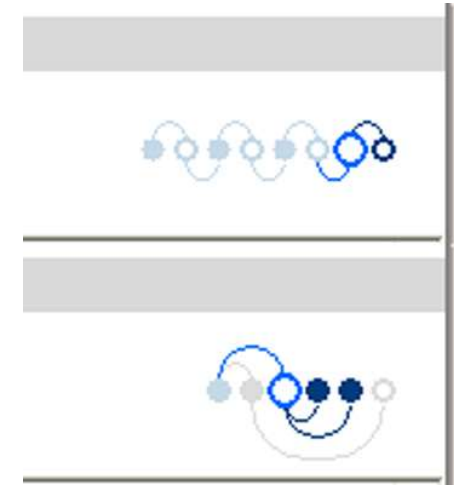
# Why Visualize Language?

- To assist information retrieval

- To enable linguistic analysis

- To augment analytics on mixed data



Themescape

Visual Thesaurus

Thread Arcs

# Visualizing Language is Difficult

- Many of the common challenges still exist
- Can you name some?

# Visualizing Language is Difficult

- Many of the common challenges still exist:
  - Screen real estate / occlusion
  - Choosing appropriate visual variable mappings
  - Colour and perception issues
  - Maintaining "graphical integrity"
  - Interaction and usability

- Specific challenges for language?

# Difficult Data

- Too much data – what to use?
  - Millions of blog posts,
  - Hundreds of thousands of news stories,
  - 183 billion emails,
  - **... per day**
- Data is noisy:
  - Newswire stories are syndicated (but differ slightly)
  - 70-72% of email is spam
  - Text contains section headings, figure captions, and direct quotes

# Once you have the data...

- Most meaning comes from our minds and common understanding.

- "How much is that doggy in the window?"
  - how much: social system of barter and trade (not the size of the dog)
  - "doggy" implies childlike, plaintive, probably cannot do the purchasing on their own
  - "in the window" implies behind a store window, not really inside a window, requires notion of window shopping

(Hearst, 2006)

# Language is Ambiguous

- Words and phrases can have many meanings, determined by context and world knowledge.

- Interesting language is often figurative:
  - "Tables encourage casual interaction."

  vs

  - "I encouraged her to take a day off."

# Language is Ambiguous

- I saw <u>Pathfinder</u> on <u>Mars</u> with a telescope.

- <u>Pathfinder</u> <u>photographed</u> <u>Mars</u>.

- The <u>Pathfinder</u> <u>photograph</u> <u>mars</u> our perception of a lifeless planet.

- The <u>Pathfinder</u> <u>photograph</u> from <u>Ford</u> has arrived.

- The <u>Pathfinder</u> <u>ford</u>ed the river without <u>marring</u> its paint job.

(Hearst, 2006)

# Data Processing Decisions

- Many levels of data processing can take place:
  - Word counting
  - Stemming
  - Parsing
  - Summarization
  - Sentiment analysis
  - Topic modelling
  - Word-Sense disambiguation
- Each step of extra processing introduces uncertainty

# Visual Considerations

Supporters of Martin, who has been jailed without trial for more than two years, are calling on Prime Minister Stephen Harper to ask Mexican president Felipe Calderon to release Martin text is not preattentive under a section of the Mexican constitution that allows the government to expel undesirables from the country. Martin's supporters believe she has no chance of a fair trial in Mexico. Neither does Waage.

# Visual Considerations

Supporters of Martin, who has been jailed without trial for more than two years, are calling on Prime Minister Stephen Harper to ask Mexican president Felipe Calderon to release Martin <span style="color:red">text is not preattentive</span> under a section of the Mexican constitution that allows the government to expel undesirables from the country. Martin's supporters believe she has no chance of a fair trial in Mexico. Neither does Waage.

# Visual Considerations

- Text readability is dependent on size, orientation, font, clutter…

# Visual Considerations

- Text readability is dependent on size, orientation, font, clutter...

- More likely to need large amounts of text in language visualization

# Visualizing language is also easy!

- SO much data available for analysis
- (Mostly) readily computer readable
- Simple techniques can give instant summaries

To be, or not to be: that is the question:
Whether 'tis nobler in the mind to suffer
The slings and arrows of outrageous fortune,
Or to take arms against a sea of troubles,
And by opposing end them? To die: to sleep;
No more; and by a sleep to say we end
The heart-ache and the thousand natural shocks
That flesh is heir to, 'tis a consummation
Devoutly to be wish'd. To die, to sleep;
To sleep: perchance to dream: ay, there's the rub;
For in that sleep of death what dreams may come
When we have shuffled off this mortal coil,
Must give us pause: there's the respect
That makes calamity of so long life;
For who would bear the whips and scorns of time,
The oppressor's wrong, the proud man's contumely,
The pangs of despised love, the law's delay,
The insolence of office and the spurns
That patient merit of the unworthy takes,
When he himself might his quietus make
With a bare bodkin? who would fardels bear,
...

Text Visualization

# BACKGROUND

Mountain Peaks of Prophecy (Larkin, 1918)

# Visual Text Analytics

- Visual techniques for words, documents, sets of documents to support rapid summarization, trend analysis, exploration, search, comparative analysis, …

- Application areas include market analysis, legal studies, e-discovery, readability, literary studies, personal reflection, information retrieval and exploration, intelligence analysis

Word Clouds


TextFlow


Topic Models


Theme River


Jigsaw


Parallel Tag Clouds


Themescape

Marian Dörk et al. VisGets: Coordinated Visualizations for Web-based Information Exploration and Discovery. *IEEE Transactions on Visualization and Computer Graphics*, 14(6):1205-1212, November-December, 2008.

# Linguistic Methods

- Word Counting
- Word Scoring
- Stemming
- Stop Word Removal
- Part of Speech Tagging
- Parsing
- Word Sense Disambiguation
- Named Entity Recognition
- Semantic Categorization
- Sentiment Analysis
- Topic Modeling (some caveats)

# NLTK: Natural Language Toolkit

- NLTK.org

- Python

## NLTK 3.0 documentation

NEXT | MODULES | INDEX

### Natural Language Toolkit

NLTK is a leading platform for building Python programs to work with human language data. It provides easy-to-use interfaces to over 50 corpora and lexical resources such as WordNet, along with a suite of text processing libraries for classification, tokenization, stemming, tagging, parsing, and semantic reasoning, and an active discussion forum.

Thanks to a hands-on guide introducing programming fundamentals alongside topics in computational linguistics, NLTK is suitable for linguists, engineers, students, educators, researchers, and industry users alike. NLTK is available for Windows, Mac OS X, and Linux. Best of all, NLTK is a free, open source, community-driven project.

NLTK has been called "a wonderful tool for teaching, and working in, computational linguistics using Python," and "an amazing library to play with natural language."

Natural Language Processing with Python provides a practical introduction to programming for language processing. Written by the creators of NLTK, it guides the reader through the fundamentals of writing Python programs, working with corpora, categorizing text, analyzing linguistic structure, and more. The book is being updated for Python 3 and NLTK 3. (The original Python 2 version is still available at http://nltk.org/book_1ed.)

### TABLE OF CONTENTS

| NLTK News |
| Installing NLTK |
| Installing NLTK Data |
| Contribute to NLTK |
| FAQ |
| Wiki |
| API |
| HOWTO |

### SEARCH

[              ] Go

Enter search terms or a module, class or function name.

# Stemming

- Reduce words to their 'stems' by removing endings (morphology)
  - running -> run
  - runs -> run

- A good way to increase signal and reduce fracturing of the corpus if there aren't many words.

- Note: Keep the original words somewhere! Also keep the case if you choose to lowercase the word; you never know when you'll need this data

# Stop Word Removal

- Common words such as "and", "the", "I" are removed from view to highlight content words

- Domain specific stop words, e.g. in legal domain:

  - Court, attorney, honour, plaintiff, etc.

- Caution! These words have been shown to be useful for stylistic analysis! When working with text corpora, KEEP EVERYTHING.

# Part of Speech Tagging

- Assign grammatical roles to words
- Conventional tagsets and representation:
  - The/AT grand/JJ jury/NN commented/VBD on/IN a/AT number/NN of/IN …
- Many words are ambiguous: fly, chair, run, store, table, and more!
  - Fly/NN
  - Fly/VB

Fly/NN

Fly/VB

# Term / Concept Ambiguity

- Most meaning comes from our minds and common understanding.

- "How much is that doggy in the window?"
  - how much: social system of barter and trade (not the size of the dog)
  - "doggy" implies childlike, plaintive, probably cannot do the purchasing on their own
  - "in the window" implies behind a store window, not really inside a window, requires notion of window shopping

(Hearst, 2006)

# Parsing

- Determining language structure
- Can reveal word-word relationships
- Useful for processing negation

S
NP — VP
I
V — NP
shot
Det N PP
an elephant P NP
in Det N
my pajamas

S
NP — VP
I
VP — PP
V — NP
shot Det N
an elephant P NP
in Det N
my pajamas

https://nltk.googlecode.com/svn/trunk/doc/book/ch08.html

# Dependency Parsing

- Labelled directed graph

- Arcs represent relationships from heads to dependents



Head of sentence

# Word Sense Disambiguation

- Susan, the meeting chair, chaired the meeting well from the big chair in the front of the room.

  – Leader of a meeting

  – Action of leading a meeting

  – An object to sit upon

# Word Sense Disambiguation

- This is VERY difficult for a computer.

- Contexts are often the same and meanings can be quite fine-grained:

  - bank the financial institution, bank the building in which the financial institution is housed

- Annual contest: SENSEVAL

- My method: assume the most common sense

# Named Entity Recognition

- What are the people, places in the text?
- Use NLTK – it's very good at this.



Much Ado About Nothing

# Semantic Categorization

- Placing a word into an ontology or sense thesaurus based on *meaning.*

- Common resources include:
  - WordNet
  - Roget's Thesaurus

# WordNet

- A large lexical database, or "digital dictionary"
- Covers most English nouns, verbs, adjectives, adverbs
- Organizes *synsets* by *meaning*
- Words are related to one another through many different relationship types:
- X is a kind of Y, X has part Y, an X Ys, X is Y/has property Y

# Hyponymy

- The "IS-A" relation for nouns

{vehicle}
/ \
{car, automobile}    {bicycle, bike}
/          \              \
{convertible}  {SUV}  {mountain bike}

# SEMANTIC VISUALIZATIONS

# DocuBurst



Collins, C.; Carpendale, S.; Penn, G. DocuBurst: Visualizing Document Content using Language Structure.
Proceedings of Eurographics/IEEE VGTC Symposium on Visualization, June, 2009.

Mihalcea and Tarau, 2004



Wattenberg et al., 2008

# DocuBurst

games→game
taken→take

absolute,noun,10
chair,noun,2
moment,noun,11
game,noun,30
reality,noun,3
take,verb,13
represent,verb,17
...

WordNet

game IS activity
chair IS furniture

**bird**
- gallinaceous bird
- bird of prey
- passerine
- piciform bird
- parrot
- aquatic bird

**fish**
- aquatic vertebrate

**vertebrate**
- salamander
- frog
- anapsid
- amphibian
- diapsid
- reptile
- metatherian

**chordate**

animal

**mammal**

**placental**

invertebrate

mollusk
- worm
- bivalve
- gastropod
- clam

**arthropod**

**insect**
- lepidopterous insect
- homopterous insect
- hemipterous insect
- orthopterous insect
- hymenopterous insect
- dipterous insect
- beetle
- arachnid
- crustacean
- acarine
- decapod crustacean

0  [____] 119.86

Search | Filter | Options | Text Segments | Concordance Lines

Focus:

Word/Sense Details:

**POS:** noun
**Synonyms:** dipterous insect, two-winged insects, dipteran, dipteron
**Sense:** insects having usually a single pair of functional wings (anterior pair) with the posterior pair reduced to small knobbed structures and mouth parts adapted for sucking or lapping or piercing

**Try it!  http://vialab.science.uoit.ca/docuburst**

# Lexichrome



**http://lexichrome.com**

Work in Progress with Chris Kim and Saif Mohammed

# ‹ all words associated with yellow

**? ABOUT LEXICHROME**

**RELEVANCE (DESC)**   ALPHABETICAL

| | | | |
|---|---|---|---|
| **cowardly**<br>10 out of 10 | **nugget**<br>7 out of 7 | **sun**<br>7 out of 7 | **sunny**<br>9 out of 10 |
| **saffron**<br>8 out of 9 | **treasure**<br>7 out of 8 | **lion**<br>6 out of 7 | **mustard**<br>6 out of 7 |
| **radiant**<br>6 out of 7 | **bee**<br>11 out of 13 | **butter**<br>11 out of 13 | **insecure**<br>6 out of 8 |
| **sandy**<br>6 out of 8 | **scatter**<br>6 out of 8 | **lightning**<br>8 out of 11 | **beehive**<br>10 out of 14 |
| **practically**<br>5 out of 7 | **radiate**<br>5 out of 7 | **enlighten**<br>7 out of 10 | **sunshine**<br>7 out of 10 |

# lexichrome<sup>alpha</sup>

■■ PALETTE  A WORDS  📄 TEXT

❓ ABOUT LEXICHROME

Nameless here for evermore.

And the silken sad uncertain rustling
of each purple curtain
Thrilled me - filled me with fantastic
terrors never felt before;
So that now, to still the beating of my
heart, I stood repeating
`'Tis some visitor entreating entrance
at my chamber door -
Some late visitor entreating entrance
at my chamber door; -
This it is, and nothing more,'

Once upon a midnight dreary, while
I pondered weak and weary,
Over many a quaint and curious
volume of forgotten lore,
While I nodded, nearly napping,
suddenly there came a tapping,
As of some one gently rapping,
rapping at my chamber door.
`'Tis some visitor,' I muttered,
`tapping at my chamber door -
Only this, and nothing more.'

**ANALYZE**

# Descriptive Non-Photorealistic Rendering

M. Chang and C. Collins, "Exploring Entities in Text with Descriptive Non-photorealistic Rendering," in *Proc. of the 2013 IEEE Pacific Visualization Symposium (PACIFICVIS '13)*, 2013.

# Ontology Generation

Keywords $\dashrightarrow$ Bike, Bicycle

Part-Of Relations $\dashrightarrow$ Bike ← Transmission ← Derailleur

Synonyms $\dashrightarrow$ Tire = Tyre = Bike tire

Manual Adjustments $\dashrightarrow$ ~~First, Second, Third~~

# Entity Extraction

"**Brakes** failed going at 35 mph."

Stemming

Cache

"brak"

6: {brak, brak system...}

Document: 123
Entity: 6
Location: 0

# Visual Representation



(0, 1]

Low Score

High Score

0

# Main Interface

# Exploration with Lens

# Semantic Password Analysis

- What types of words do people use in their passwords?

- Do the patterns of word use represent security vulnerabilities?

R. Veras, C. Collins, and J. Thorpe, "On Semantic Patterns of Passwords and their Security Impact," *In Proceeding of the Network and Distributed System Security Symposium (NDSS'14)*, 2014.

- Extract words from 32 million passwords

- Categorize them

- Parse the results to find structure

- Create a password guessing system based on the model

| Password | Segment | Semantic tag |
| --- | --- | --- |
| hope87 | hope | wish.v.01 |
| hope87 | 87 | number |
| serenity | serenity | trait.n.01 |
| bishop5 | bishop | status.n.01 |
| bishop5 | 5 | number |
| goblue0507 | go | s.travel.v.01 |
| goblue0507 | blue | |
| goblue0507 | 507 | number |
| looted | looted | take.v.21 |
| drift21 | drift | force.n.02 |
| drift21 | 21 | number |
| candysinger | candy | s.candy.n.01 |
| candysinger | singer | musician.n.01 |
| 671soldier | 671 | number |
| 671soldier | soldier | worker.n.01 |
| bravo100 | bravo | murderer.n.01 |
| bravo100 | 100 | number |
| egobrain | ego | pride.n.01 |
| egobrain | brain | structure.n.04 |
| pitcher9 | pitcher | athlete.n.01 |
| pitcher9 | 9 | number |
| puppies | puppies | puppy.n.01 |
| church | church | religion.n.02 |
| 'ale'8 | ' | special |
| 'ale'8 | ale | alcohol.n.01 |
| 'ale'8 | '8 | num+special |

# Appropriate Levels of Detail

# Results

- Created best cracker on several measures, including percent correct guesses

- Designing strategies to help people make passwords more *semantically* secure – keep the meaning but lower the probability

# Results

- Created best cracker on measure of % correct guesses
- Place names, male names very popular
- "Cute" animals more common:
  – Monkey, dogs, cats, dolphins
- Emotional verbs like "love" are common
  – People "love" male names 4x more often than female!
- Profanity is very common

**thestar.com**

# GTA

## Is there 'love' in your online passwords?

After analyzing 32 million leaked passwords, a team of researchers from the University of Ontario Institute of Technology has discovered that "love" is the most common password verb.

| f | Tweet 55 | 8+1 1 | reddit this! | ≡ |



**By: Daniel Otis** News Reporter, Published on Fri Feb 13 2015

People are putting a little too much "love" into their online passwords.

At least that's what a team of researchers from the University of Ontario Institute of Technology (UOIT) says. They analyzed 32 million leaked passwords from the now-defunct RockYou.com website. The project was led by UOIT graduate student Rafael

---

**The New York Times** | http://nyti.ms/1xqfNJL

# The Secret Life of Passwords

We despise them – yet we imbue them with our hopes and dreams, our dearest memories, our deepest meanings. They unlock much more than our accounts.

By IAN URBINA  Video by LESLYE DAVIS

Howard Lutnick, the chief executive of Cantor Fitzgerald, one of the world's largest financial-services firms, still cries when he talks about it. Not long

Text Visualization

# LITERARY ANALYSIS

# Document Lens

# TextArc

# Literary Analysis: Semantics

# Literary Analysis: PoemViewer

- Phonetics, repetition



http://ovii.oerc.ox.ac.uk/PoemVis/

# Literary Analysis: Poemage



http://www.sci.utah.edu/~nmccurdy/Poemage/

# Literary Analysis: Patterns



http://www.cs.umd.edu/hcil/textvis/featurelens/

# Literary Analysis: Patterns

(Werschkul, 2007)

# Twitter Contrast Diagrams



Clinton's Super Tuesday Speech / Obama's Super Tuesday Speech

(Clark, 2008)

Clinton's Super Tuesday Speech

Obama's Super Tuesday Speech

# VisArgue Project



## Lexical Episode Plots

Lexical Episodes are defined as a portion within the word sequence where a certain word appears more densely than expected from its frequency in the whole text. For example, if the text contains 100 words and a certain word appears four times within the whole corpus, we would assume -with an equidistant distrubution- that this word would appear every 25 words in the text. However, if the actual distribution of this word is more dense in a certain part of the text, we define this as an episode.

Mennatallah El-Assady, http://presidential-debates.dbvis.de/

# Many Eyes

- IBM system for uploading your data and visualizing it

- You can share your visualizations (recall: empowerment aspect of Critical Visualization)

- http://www-969.ibm.com/software/analytics/manyeyes/

# Literary Analysis: Repetition

# Literary Analysis: Repetition

# Literature Fingerprinting



(c) Average sentence length

(d) Simpson's Index

(e) Hapax Legomena

(f) Hapax Dislegomena

Keim, D. A., & Oelke, D. (2007). Literature Fingerprinting: A New Method for Visual Literary Analysis. In *2007 IEEE Symposium on Visual Analytics Science and Technology* (pp. 115–122).

# Visual Readability Analysis



Figure 2: Normalization of the feature values is done relatively to the values that we observed for our ground-truth data set. The graphic shows the formulas and color scales for the 3 different cases that are possible.



colormap of overall readability score

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.

# Visual Readability Analysis



| | | 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. | | | | | |
| (b) | This includes aspects like ensuring contextual coherency, avoiding unknown vocabulary and difficult grammatical structures. | | | | | |

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.



(a)　　　　　　　　　　(b)　　　　　　　　　　(c)

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.

# Literary Analysis: Readability



**The state of our union is ... dumber:**
How the linguistic standard of the presidential address has declined

Using the Flesch-Kincaid readability test the Guardian has tracked the reading level of every State of the Union

http://www.theguardian.com/world/interactive/2013/feb/12/state-of-the-union-reading-level

# Gist Icons

- Word counts
- Automatic groupings
- Drill-down

# Gist Icons

# Computational Journalism

- Overviewproject.org

# Timeline Curator

# Social Patterns from Text

- HistoryFlow – edits on Wikipedia

Fernanda B. Viégas, Martin Wattenberg, and Kushal Dave. 2004. Studying cooperation and conflict between authors with *history flow* visualizations. In Proc. *CHI* '04.

# WordsEye.com

# SENTIMENT VISUALIZATION

# Sentiment Analysis

- Business intelligence:
  - Do people like my product/restaurant/movie/hotel?
  - Why or why not?
- Forensics and medicine:
  - State of mind analysis based on social media
- Emotional profiling / psycholinguistics
  - Understanding users -> individualization
  - Targeted advertising

# Sentiment Analysis

- Language Processing:
  - Stemming
  - POS Tagging
  - Dependency Parsing
  - Named Entity Detection
- Granularity:
  - Positive/negative/uncertain
  - 8+ emotions
  - Word, sentence, paragraph, document, corpus level

# Resources and Datasets

- NRC Word-Emotion Lexicon:
  - Saif Mohammad, 2013
    http://www.saifmohammad.com/WebPages/ResearchInterests.html

- LIWC:
  - James Pennybaker et al., 2007:
    http://www.liwc.net/

- Opinion Mining Dataset:
  - Bing Liu, 2004—current
    http://www.cs.uic.edu/~liub/FBS/sentiment-analysis.html

(Harris, 2006)

# Twitter Sentiment Viz



Healey and Ramaswamy, 2013. http://www.csc.ncsu.edu/faculty/healey/tweet_viz/

# SentimentState

- Tweets over time, categorized using an emotion lexicon

- Examine Tweets in context, filter based on time and emotions

Scantlebury and Collins, 2014. http://vialab.science.uoit.ca/sentimentstate

# SentimentState

*This movie was actually neither that funny, nor super witty.*

*This movie was actually neither that funny, nor super witty.*

# Stanford Sentiment Parser

- Recursive neural network built on top of grammatical structures

- Trained on Stanford Sentiment Treebank
  - Parse trees labelled with sentiment scores
  - Crowed-sourced and editable

Socher et al. Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank. Conference on Empirical Methods in Natural Language Processing (EMNLP 2013).

## Sentiment Trees

You can double-click on each tree figure to see its expanded version with greater details. There are 5 classes of sentiment classification: very negative, negative, neutral, positive, and very positive.



This movie does n't care about cleverness , wit or any other kind of intelligent humor

All labels are now correct

# Parsing is Needed!

- Stanford Sentiment Treebank:
    - http://nlp.stanford.edu/sentiment/treebank.html

# Challenges

- Word-counting techniques are fast, but inaccurate
  - Sarcasm, quotes, metaphorical language
- Accurate methods are slow/difficult to run over big data

# Opinion Seer



(a)

(b)

Yingcai Wu et al. 2010. OpinionSeer: Interactive Visualization of Hotel Customer Feedback. *IEEE Transactions on Visualization and Computer Graphics* 16 (6), November 2010.

Fig. 8. OpinionSeer results showing how customer opinions are correlated with trip type, gender, age range, and ratings.

# INFORMATION RETRIEVAL

Select Collections

Home | Sign In | Help

Search | View | Bookshelf

Simple | Advanced | Browse

Search for
computer

Search

Improve search results with Search Tips...

English | Go

ebrary

Search results: 25693 documents

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 Next>>

Sort results by: Score | Title | Contributor | Date

Improve search results with Search Tips...

Result 1   Score: 97%
Computer Science: Reflections on the Field, Reflections from the Field
BOOK - 216 Pages
View
Contributor: Committee on the Fundamentals of Computer Science: Challenges and Opportunities
Publisher: National Academies Press
Date: 2004
LC Call Number: QA76.C66 2004eb

Result 2   Score: 9
Essential Comp
BOOK - 322 Pages
View

Result 3   Score: 9
Group Cognitio
BOOK - 521 Pages
View
LC Call Number: LB1028.5.S696 2006eb
ISBN: 97-0-262-19539-3
Subjects: Computer-assisted instruction.,
Computer networks.

Result 4   Score: 95%
Facial Analysis from Continuous Video with Applications To Human Computer Interface
BOOK - 159 Pages
View
Contributor: Colmenarez, Antonio J.
Publisher: Kluwer Academic Publishers
Date: 2004
Dewey: 004/.01/9
LC Call Number: QA76.9.H85.C653 2004ebeb
ISBN: 1-40-207802-1
Subjects: Human-computer interaction.,
Image processing -- Digital techniques.,
Computer vision.

Result 5   Score: 95%
More Than a Game : The Computer Game as Fictional Form
BOOK - 177 Pages
View
Contributor: Atkins, Barry
Publisher: Manchester University Press
Date: 2003
Dewey: 306.4/87/0285
LC Call Number: GV1469.17.S63.A85 2003ebeb
ISBN: 0-7190-6364-7
Subjects: Computer games -- Social aspects.,
Computer games.

Result 6   Score: 95%
Contributor: Lengyel, Eric

Search results: 25693 documents

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 Next>>

# Information Retrieval

- Visual query formation

- Exploration of collections

- Single/comparative document content visualization

# Visual Query Formation

- Rich specification of linguistic constraints

# VQuery

# VisGets

# Exploration of Collections

- Provide overview of:
  - entire collection
  - subset matching a query
- Clustering and categorization

# Galaxies

# Tile Bars

**Term Set 1:** law legal attorney lawsuit

**Term Set 2:** network lan

*TileBars*

| | |
|---|---|
| 1256 | Regression testing handling ha |
| 1269 | Toll fraud includes related articl |
| 1270 | In conversation Teleglobe Cana |
| 1280 | Deregulation indicates a health |
| 1298 | The last word letters to the edi |
| 1300 | What's wrong with network lice |
| 1302 | Letters letter to the editor |
| 1356 | Protecting information now vita |
| 1414 | Letters O |
| 1424 | Loose LIPS sink ships logical in |

# Lighthouse Cross-Lingual Search

# iNeATS Summarizer

# Corpus Comparison

# INTELLIGENCE ANALYSIS

# Characteristics

- Multiple data streams

- Different data types: geolocation, phone calls, travel records, paper records, video surveillance, …

- Streaming data at different rates

- High cost of failure

# Hotel Visits



Weaver, C.; Fyfe, D.; Robinson, A.; Holdsworth, D.; Peuquet, D.; MacEachren, A.M., "Visual Analysis of Historic Hotel Visitation Patterns," IEEE VAST, 2006

# Jigsaw



John Stasko and colleagues, numerous papers: http://www.cc.gatech.edu/gvu/ii/jigsaw/

# Jigsaw:
# Supporting Investigative Analysis
# through Interactive Visualization

John Stasko, Carsten Görg,
Zhicheng Liu, Kanupriya Singhal

School of Interactive Computing & GVU Center
Georgia Institute of Technology

Text Visualization

# OPEN RESEARCH AREAS

# Trust

# Unquantifiable uncertainty

| Source Data | NLP | Visualization | ??? |

# Multilanguage

- Do the same techniques work for non-Western languages?

# Term / Concept Ambiguity

- Most meaning comes from our minds and common understanding.

- "How much is that doggy in the window?"
  - how much: social system of barter and trade (not the size of the dog)
  - "doggy" implies childlike, plaintive, probably cannot do the purchasing on their own
  - "in the window" implies behind a store window, not really inside a window, requires notion of window shopping

(Hearst, 2006)

# Finding 'Sweet Spots' in the Hierarchy

+ Meta data at each level!

Different levels for each genre?

| Multi-corpora |
| Document Collection |
| Document Clusters |
| Document |
| Chapter |
| Section |
| Paragraph |
| Sentence |
| Multi-word collocates |
| Word |
| Letters |

?

# VariFocalReader



S. Koch, M. John, M. Wörner, A. Müller and T. Ertl, "VarifocalReader — In-Depth Visual Analysis of Large Text Documents," in *IEEE Transactions on Visualization and Computer Graphics*, vol. 20, no. 12, pp. 1723-1732, Dec. 31 2014.

# Multi-Media Documents





Audio + Text Analysis (e.g. court proceedings)

# Appropriate Abstraction

# Proper Nouns

- Do not appear in ontologies like WordNet
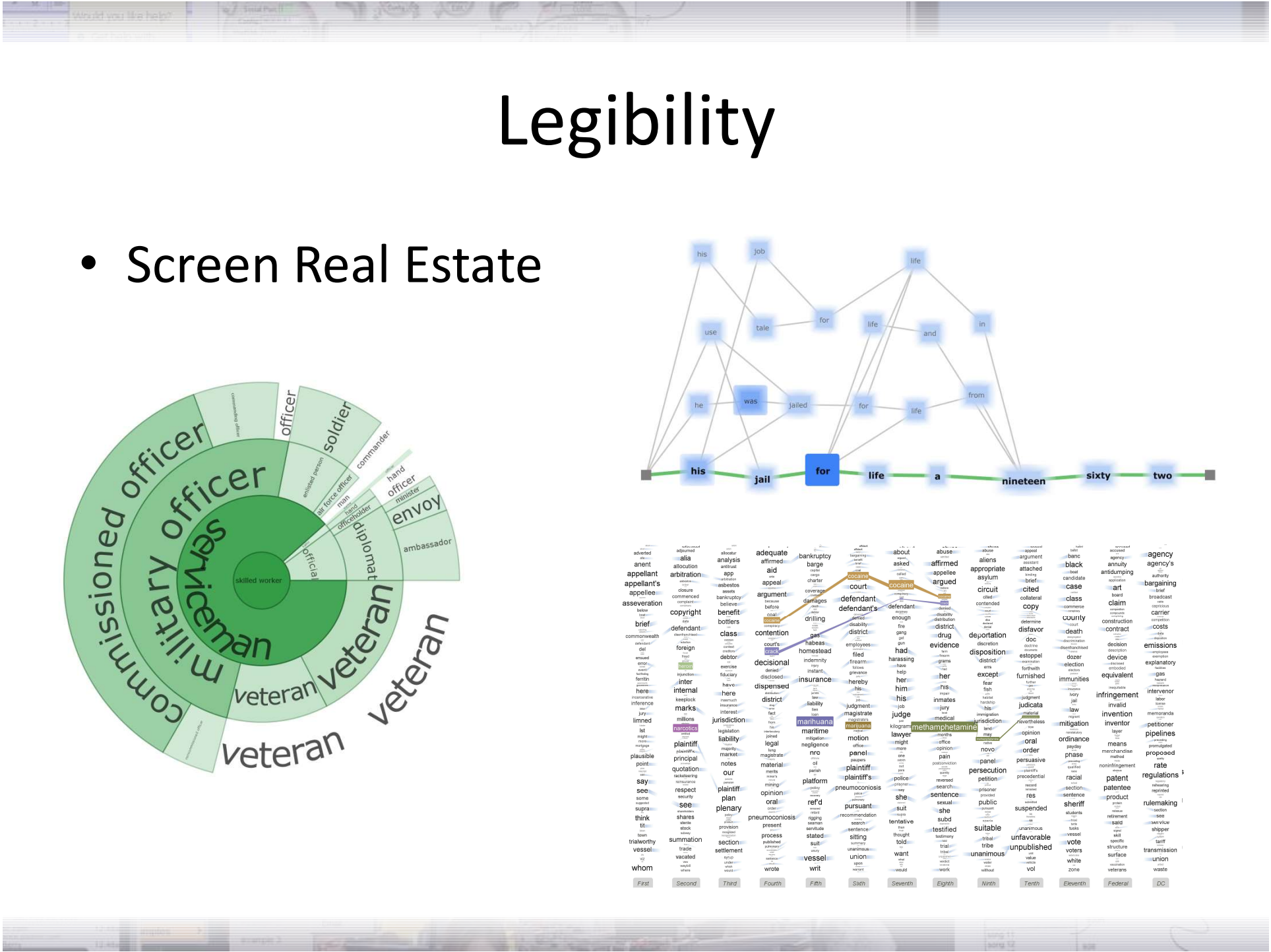- Challenging to translate

- Holy Grail: NNPs + sentiment + visualization

# Interactivity: Linking to Text

# Legibility

- Screen Real Estate

# Legibility
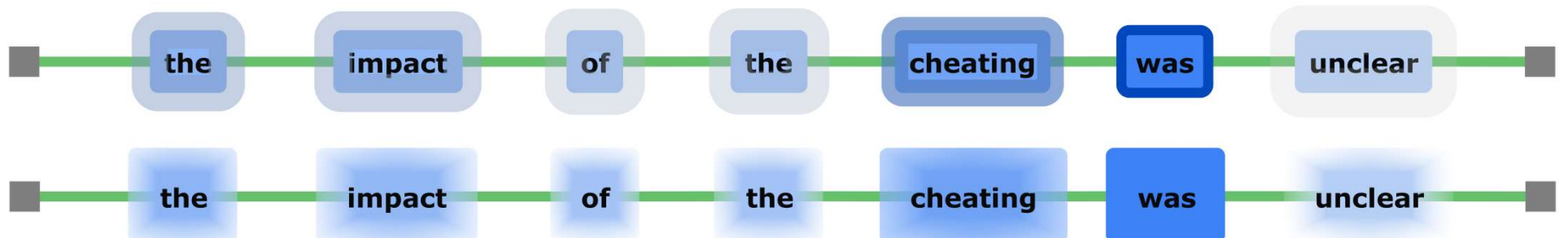
- Orientation

# Legibility

- Overlay or Background Interference

# Domains of Application

- Medicine: electronic medical records
- Business: social media analytics, corporate document collection management
- Crime Prevention and Intelligence Analysis: find threats in communications and blogs
- Legal: sift through evidence, e.g. millions of emails, to investigate fraud
- Literary and History Scholarship

# CONCLUSION

# TextVis Survey

http://textvis.lnu.se/

# Summary

- Text visualization is an exciting area of ongoing research – check out recent workshop papers at textvis.org and vis4dh.org

- Text visualization bridges visualization design, interaction design, and natural language processing