# VISUALIZING TEXT

## Petra Isenberg

# RECAP

## STRUCTURED DATA



| 0.103 | 0.176 | 0.387 | 0.300 | 0.379 |
|-------|-------|-------|-------|-------|
| 0.333 | 0.384 | 0.564 | 0.587 | 0.857 |
| 0.421 | 0.309 | 0.654 | 0.729 | 0.228 |
| 0.266 | 0.750 | 1.056 | 0.936 | 0.911 |
| 0.225 | 0.326 | 0.643 | 0.337 | 0.721 |
| 0.187 | 0.586 | 0.529 | 0.340 | 0.829 |
| 0.153 | 0.485 | 0.560 | 0.428 | 0.628 |

## UNSTRUCTURED DATA

(TODAY)

# VISUALIZING TEXT

TEXT IS DIFFERENT

COMMON

UNSTRUCTURED (MOSTLY)

HIGH-DIMENSIONAL (10,000+)

BIG!

**TEXT?**

# WHY

- To assist information retrieval
- To enable linguistic analysis
- To augment analytics on mixed data



Themescape



Visual Thesaurus



Thread Arcs

## WHY

**UNDERSTANDING**: GET THE "GIST" OF A DOCUMENT

**GROUPING**: CLUSTER FOR OVERVIEW OR CLASSIFICATION

**COMPARE**: COMPARE DOCUMENT COLLECTIONS, OR INSPECT EVOLUTION OF COLLECTION OVER TIME

**CORRELATE**: COMPARE PATTERNS IN TEXT TO THOSE IN OTHER DATA, E.G., CORRELATE WITH SOCIAL NETWORK

# WHAT IS TEXT

## DOCUMENTS

ARTICLES, BOOKS AND NOVELS
COMPUTER PROGRAMS
E-MAILS, WEB PAGES, BLOGS
TAGS, COMMENTS

## COLLECTION OF DOCUMENTS

MESSAGES (E-MAIL, BLOGS, TAGS, COMMENTS)
SOCIAL NETWORKS (PERSONAL PROFILES)
ACADEMIC COLLABORATIONS (PUBLICATIONS)
EVEN WHOLE LIBRARIES, WEBSITES, SOCIAL NETWORKS

# 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:
  - 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:
  - *America is a melting pot* (metaphor)
  - *Busy as a bee* (simile)
  - *Opportunity knocked on the door* (personification)
  - *You could have knocked me over with a feather* (hyperbole)

# LANGUAGE IS AMBIGUOUS

"I can't tell you how much I enjoyed meeting your husband."

(William Empson, *Seven Types of Ambiguity*, 1947)

"Brave men run in my family."

(Bob Hope as "Painless" Peter Potter in *The Paleface*, 1948)

# 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 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

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

## OUTLINE

- TEXT AS DATA
- VISUALIZING DOCUMENT CONTENT
- EVOLVING DOCUMENTS
- DOCUMENT COLLECTIONS

# TEXT AS DATA

**Words** are the basic unit of data.

# WORD-LEVEL ATTRIBUTES

- WORD LENGTH

- PART OF SPEECH (NOUN, VERB, ADJECTIVE, ETC.)

- FORMAT (*ITALIC,* UNDERLINE, ETC.)

- LANGUAGE (ENGLISH? LATIN? JAPANESE?)

- FREQUENCY / DIFFICULTY (IS IT COMMON?)

- SENTIMENT (POSITIVE OR NEGATIVE CONNOTATION)

- SYNONYMS / ANTONYMS / ETYMOLOGY (OTHER MEANINGS? ROOTS?)

- ENTITIES (e.g. "Calgary", "Obama", "Telus" )


- … AND MANY MORE

# AGGREGATION

REPETITION
PLAGARISM
SHARED ENTITIES
AUTHOR STYLE

**COLLECTION**

▲

• DOCUMENT

▲

• SECTION

▲

• PAGE

▲

• PARAGRAPH

▲

• SENTENCE

▲

• WORD

TENSE

SENTIMENT

SENTENCE LENGTH

READING LEVEL

# 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)

# WHAT ABOUT THESE WORDS?

automate
automates
automatic        ➡  **automat**
automation

~~a, an, the, to, ...~~

" New York
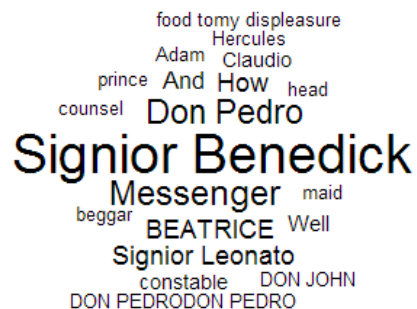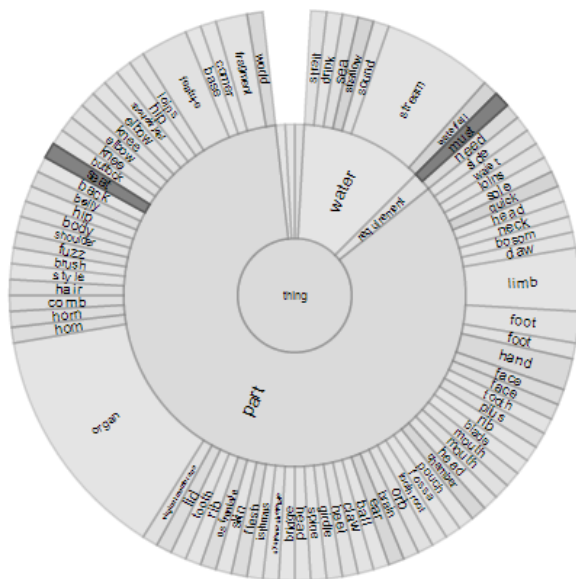" Ban Ki-moon
" Manchester United

# 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.

# NAMED ENTITY RECOGNITION

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



http://vialab.science.uoit.ca/docuburst

Much Ado About Nothing

0

0 ▮▮▮ 11.67

# TEXT PROCESSING

## TOKENIZATION: SEGMENT TEXT INTO TERMS

ENTITIES? "SAN FRANCISCO", "O'CONNOR", "U.S.A."

REMOVE STOP WORDS? "A", "AN", "THE", "TO", "BE"

N-GRAMS? CAN TAKE WORDS IN 2-WORD GROUPS (BI-GRAMS), 3-WORD (TRI-GRAMS), ETC.

## STEMMING: GROUP TOGETHER DIFFERENT FORMS

ROOTS: VISUALIZATION(S), VISUALIZE(S), VISUALLY ➡ VISUAL

LEMMATIZATION: GOES, WENT, GONE ➡ GO

FOR VISUALIZATION, SOMETIMES NEED TO REVERSE STEMMING FOR LABELS

SIMPLE SOLUTION: MAP FROM STEM TO THE MOST FREQUENT WORD

## RESULT: ORDERED STREAM OF TERMS

# TEXT PROCESSING

"The quick brown fox jumps over the lazy dog."

TOKENIZE (N=1)
[The], [quick], [brown], [fox], [jumps], [over], [the], [lazy], [dog].

TOKENIZE (N=1), REMOVE STOPWORDS, STEM
[quick], [brown], [fox], [jump], [over], [lazy], [dog]

TOKENIZE (N=2)
[the quick], [quick brown], [brown fox], [fox jumps], [jumps over], [over the]…

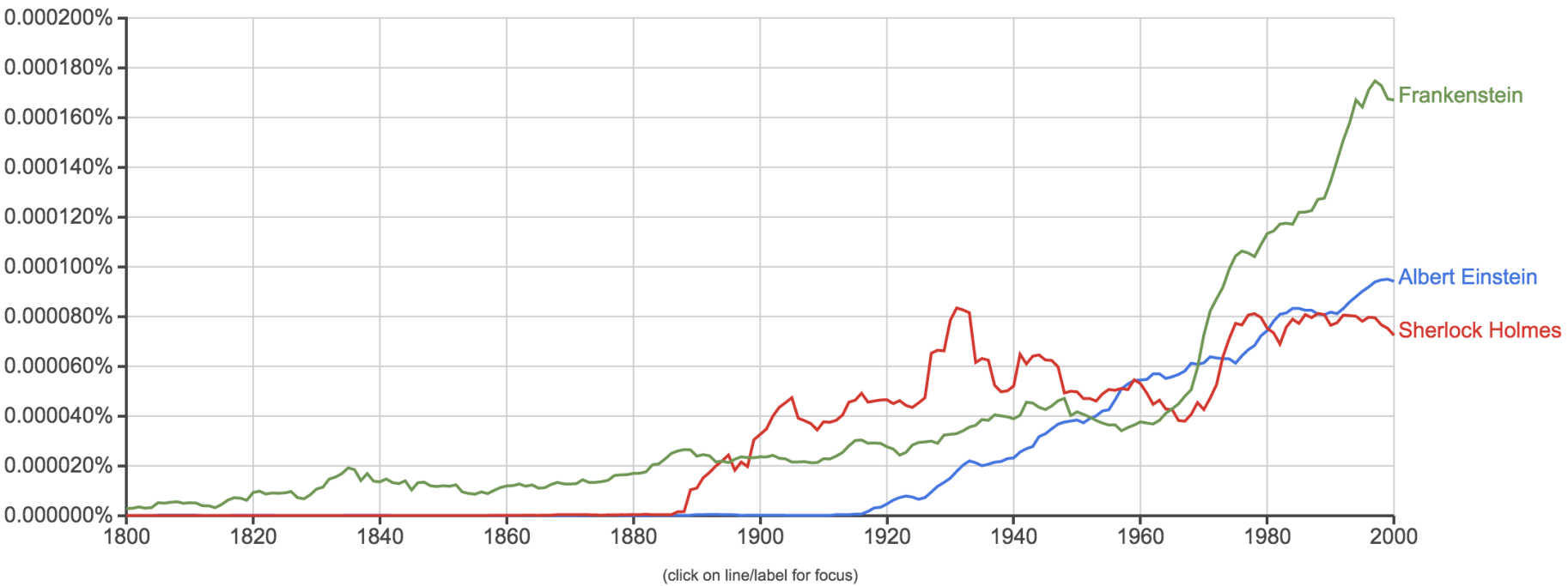TOKENIZE (N=5)
[the quick brown fox jumps], [quick brown fox jumps over], [brown fox jumps over …

# NLTK (NATURAL LANGUAGE TOOLKIT)

Tokenize and tag some text:

```
>>> import nltk
>>> sentence = """At eight o'clock on Thursday morning
... Arthur didn't feel very good."""
>>> tokens = nltk.word_tokenize(sentence)
>>> tokens
['At', 'eight', "o'clock", 'on', 'Thursday', 'morning',
'Arthur', 'did', "n't", 'feel', 'very', 'good', '.']
>>> tagged = nltk.pos_tag(tokens)
>>> tagged[0:6]
[('At', 'IN'), ('eight', 'CD'), ("o'clock", 'JJ'), ('on', 'IN'),
('Thursday', 'NNP'), ('morning', 'NN')]
```

NLTK.org
Python

Identify named entities:

```
>>> entities = nltk.chunk.ne_chunk(tagged)
>>> entities
Tree('S', [('At', 'IN'), ('eight', 'CD'), ("o'clock", 'JJ'),
           ('on', 'IN'), ('Thursday', 'NNP'), ('morning', 'NN'),
       Tree('PERSON', [('Arthur', 'NNP')]),
           ('did', 'VBD'), ("n't", 'RB'), ('feel', 'VB'),
           ('very', 'RB'), ('good', 'JJ'), ('.', '.')])
```

# NAMED ENTITY RECOGNITION

## IDENTIFY AND CLASSIFY NAMED ENTITIES IN TEXT:
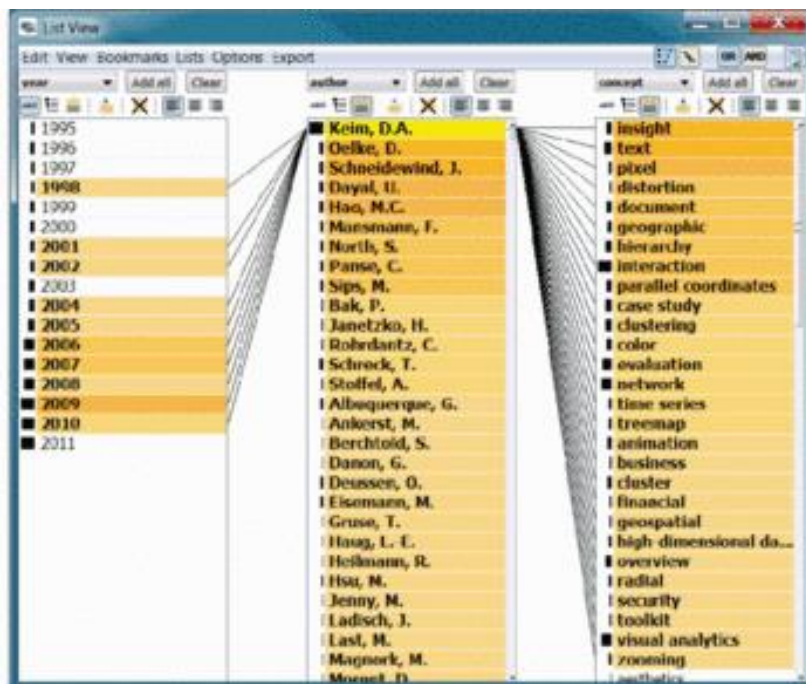
JOHN SMITH IS A **PERSON**

SOVIET UNION IS A **COUNTRY**

2500 UNIVERSITY DR IS AN

**ADDRESS**

(555) 867-5309 IS A **PHONE NUMBER**

## ENTITY RELATIONS: HOW DO THE ENTITIES RELATE?

## DO THEY CO-OCCUR IN A DOCUMENT? IN A SENTENCE?

JIGSAW

A

CENDARI
COLLABORATIVE EUROPEAN DIGITAL
ARCHIVE INFRASTRUCTURE

Home   Browse   About   Issue Report   Survey      Search      anthi ▾

B  Resources  «

C  New ▾   Save      Import ▾   Help

D  Visualizations  »

**My Projects:**
- Green Cadres
- WW1
  - Notes (1)
    - Green Cadres Notes
  - Documents (144)
  - Entities (7)
    - Event (1)
    - Organization (0)
    - Person (3)
    - Publication (0)
    - Artifact (0)
    - Place (5)
    - Tag (3)

Note 5: Green Cadres Notes

Entities (12)   Status (Open)   Assigned Users

Green Cadres Notes

Note Description [ Read Only ]        --- click here for Edit mode

In 1918, as privations and social unrest began to undermine the Austro-Hungarian war effort on the home front, a specific kind of revolt gripped the countryside in a number of regions of the empire. The so-called Green Cadres or Green Brigades were groups of armed deserters, supplemented by the local poor peasantry, who hid themselves in forested areas, staging raids on livestock and crops, attacking the local gendarmerie and military, and (in some instances) articulating social revolutionary programs. Reports on these irregular armed bands abounded in the final year of the year in many regions of both Austria and Hungary, but they were concentrated in Croatia-Slavonia (current Croatia and Serbia) and southern Moravia (current Czech republic). The Green Cadres represented a specifically rural form of unrest—largely unhitched from nationalist and party political agendas—reflecting the widespread sense of apocalyptic collapse among the rural population of Austria-Hungary.

The historical research on the Green Cadres is scant and preponderantly concentrated on the region of Croatia-Slavonia, where the Cadres where most numerous and their actions most ambitious. Communist-era Yugoslav scholarship treated the Green Cadres as proto-Bolsheviks, overemphasizing the prevalence of Leninist ideas among them. Indeed, research has revealed that soldiers returning from Russian imprisonment in 1918 played leading roles in mass desertions, mutinies, and the propagation of social-revolutionist ideas. But scholars have not identified the specific mechanisms by which former POWs became Green Cadres or how the Russian experience was reinterpreted in rural Austro-Hungarian contexts. More importantly, a comparative study of the cadres in various regions is missing because of the challenges of finding, organizing, and interpreting sources that are now fragmented in various national archival research 'siloes'.

This project seeks to open up comparative vistas on the problem of the Green Cadres. Among the possible questions it seeks to answer are: 1. How did the far-flung groups identified as Green Cadres compare to each other in terms of actions and aims; 2. Why did the Cadres appear in the places that they did?; 3. What were the social, political, and cultural factors that facilitated the formation or concentration of Cadres in specific locales? 4. What kind of deserters made up the bulk of the Cadres—deserters from the front, replacement regiments, or allotted leave after returning from Russian internment?; 5. What played a bigger role in the formation of Green Cadres: social revolutionary influences from Russian imprisonment or disillusionment with the war effort?

Most Common Person  FRAPET, Guillaume

Most Common Place  Nantes  128 docs

Most Recent  Date: 1711/1/29  1711-1-29
Oldest  Date: 1669/6/5  1669-6-5

1670   1680   1690   1700

Most Common Place  Nantes  128 docs

+
−

CENDARI NOTE-TAKING ENVIRONMENT 2015

DOCUMENT CONTENT

BUT FIRST **SOME SKETCHING**

# SKETCHING: **VISUALIZE**

IMAGINE YOU HAVE A MASTER'S
THESIS IN FROM OF YOU:

YEAR

AUTHOR

TITLE

KEYWORDS

REFERENCES

ABSTRACT TEXT

TASK:

**1)** **VISUALIZE THE MOST IMPORTANT CONTENT** FROM A SINGLE THESIS.

(~10 MINUTES)

# EXAMPLE

**Tools & Strategies for Social Data Analysis**

by

Wesley Jay Willett

A dissertation submitted in partial satisfaction of the
requirements for the degree of
Doctor of Philosophy

in

Computer Science

in the

Graduate Division

of the

University of California, Berkeley

**THESIS** WESLEY WILLETT

# TAG CLOUDS

## WORD COUNT

additional air **analysis analysts** annotation applications approach asked author average based build **chart** citizen **clustering** collaborative collection **comments** commentspace **community** complete **condition** contributions crowd crowdsourcing **data** datasets **design** different **discussion** evidence **example** experiment experts **explanations** explore features **figure** filtering **generated** group help hypotheses hypothesis **identify including** indicating information interactive interface knowledge **links** members microtasks multiple novice **number** oaen observations organize **participants** phases pp proceedings **process produced** prompt **provide quality** questions rate redundant requires **responses results** score sense share showing similar site **social source** specific state strategies study **support** systems **tags tasks tools** understanding used **users views** **visualization** web work **workers**

THESIS WESLEY WILLETT

# TAG CLOUDS

## WORD COUNT

# WHAT'S PROBLEMS DO YOU SEE WITH TAG CLOUDS?

# TAG CLOUDS

## STRENGTHS

CAN HELP WITH GISTING AND INITIAL QUERY FORMATION.

## WEAKNESSES

SUB-OPTIMAL VISUAL ENCODING (SIZE VS. POSITION)
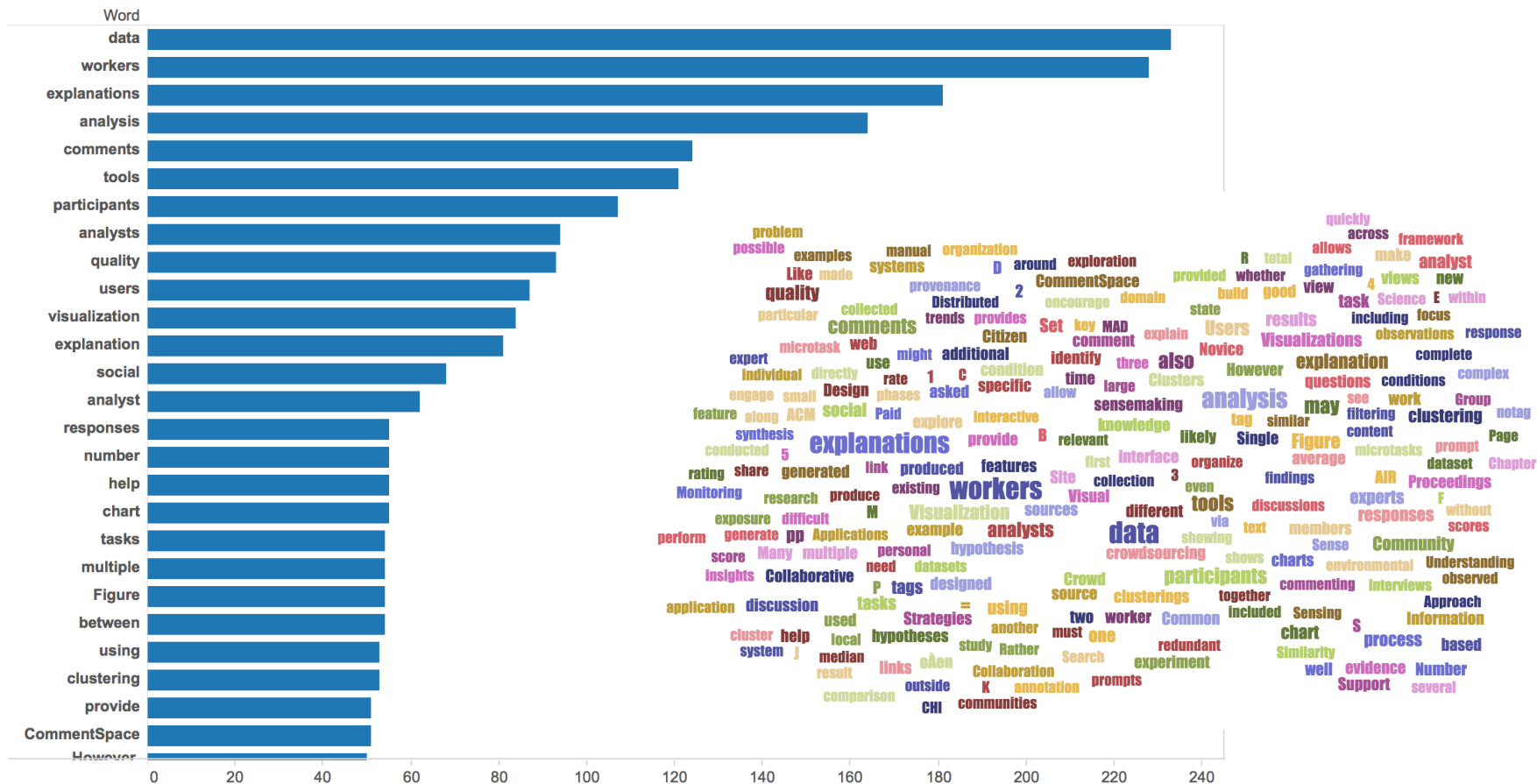INACCURATE SIZE ENCODING (LONG WORDS ARE BIGGER)
MAY NOT FACILITATE COMPARISON (UNSTABLE LAYOUT)
ORDER USUALLY MEANINGLESS (USUALLY ALPHABETICAL OR RANDOM)
TERM FREQUENCY MAY NOT BE MEANINGFUL
DOES NOT SHOW THE STRUCTURE OF THE TEXT

# WORD COUNTS

# WORDCOUNT

PREVIOUS WORD

NEXT WORD

the of and to a in that it is was i for on you he be with as by at have are this not but had his they from she which or we an there her were do you all one has will

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24

CURRENT WORD

FIND WORD: ▶   BY RANK: ▶   REQUESTED WORD: **THE**

ARCHIVE

COUNT

RANK: 1

http://wordcount.org

JONATHAN HARRIS

# CONCORDANCE

## WHAT IS THE COMMON LOCAL CONTEXT OF A TERM?

# WORD TREES

- cats are better than dogs
- cats eat kibble
- cats are better than hamsters
- cats are awesome
- cats are people too
- cats eat mice
- cats meowing
- cats in the cradle
- cats eat mice
- cats in the cradle lyrics
- cats eat kibble
- cats for adoption
- cats are family
- cats eat mice
- cats are better than kittens
- cats are evil
- cats are weird
- cats eat mice

**love the**

- **lord**
  - **thy god**
    - **with all**
      - thine heart , and with all thy soul ,
        - **and** with all thy might .
        - **that** thou mayest live .
      - thy heart , and with all thy soul , and with all thy
        - **mind**
        - **strength** , a
    - **and** ,
      - **keep** his charge , and his statutes , and his judgments , and his commandments , alway .
      - **to** walk ever in his ways ; then shalt thou add three cities more for thee , beside these three : 19
      - **that** thou mayest obey his voice , and that thou mayest cleave unto him : for he is thy life , and t
    - **to** walk in his ways , and to keep his commandments and his statutes and his judgments , that thou mayest liv
  - **your god**
    - **and to** ,
      - **serve** him with all your heart and with all your soul , 11 : 14 that i will give you the rain of your la
      - **walk** in all his ways , and to keep his commandments , and to cleave unto him , and to serve him
    - **to** walk in all his ways , and to cleave unto him ; 11 : 23 then will the lord drive out all these nations from
    - **with** all your heart and with all your soul .
    - .
  - **,**
    - **all** ye his saints : for the lord preserveth the faithful , and plentifully rewardeth the proud doer .
    - **hate** evil : he preserveth the souls of his saints ; he delivereth them out of the hand of the wicked .
    - **because** he hath heard my voice and my supplications .
- **name** of the lord , to be his servants , every one that keepeth the sabbath from polluting it , and taketh hold of my covenant
- **good** , and establish judgment in the gate : it may be that the lord god of hosts will be gracious unto the remnant of joseph
- **evil** ; who pluck off their skin from off them , and their flesh from off their bones ; 3 : 3 who also eat the
- **truth** and peace .
- **other ; or else he will hold to the one , and despise the other . ye cannot serve god and mammon .**
  - 6 : 25 therefore i say unto
  - 16 : 14 and the pharisees
- **uppermost**
  - **rooms** at feasts , and the chief seats in the synagogues , 23 : 7 and greetings in the markets , and to be called of
  - **seats** in the synagogues , and greetings in the markets .
- **father**
  - ; and as the father gave me commandment , even so i do .
  - **hath** bestowed upon us , that we should be called the sons of god : therefore the world knoweth us not , because it knew him
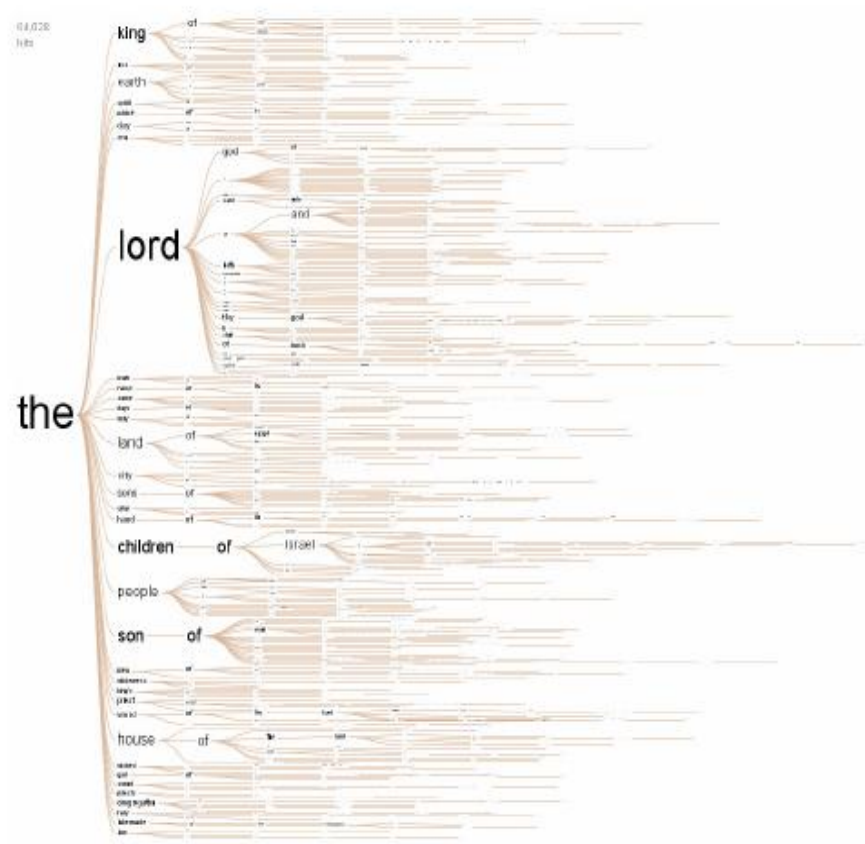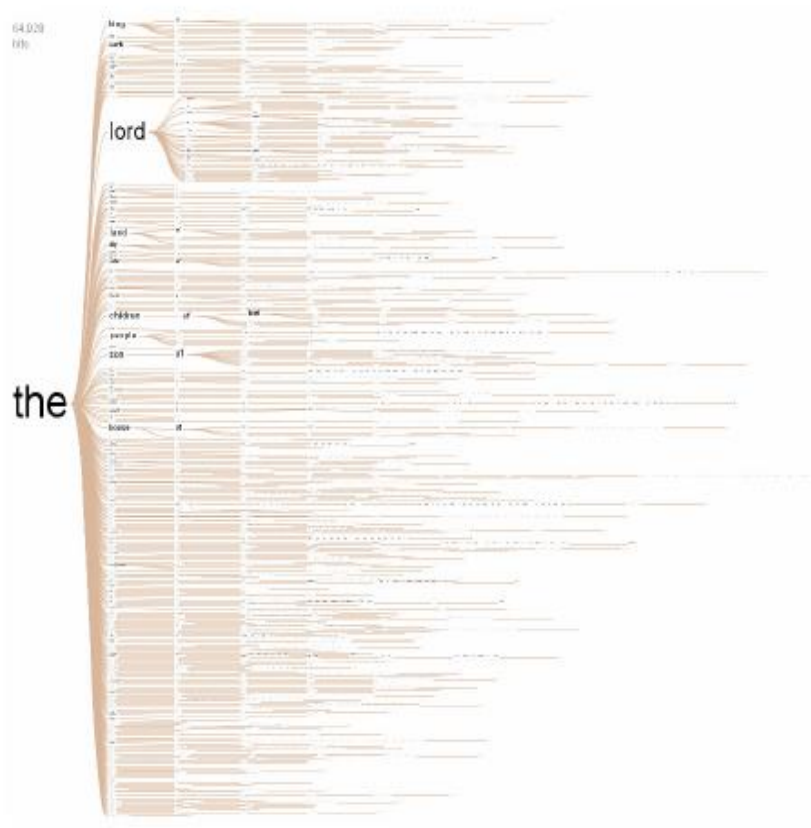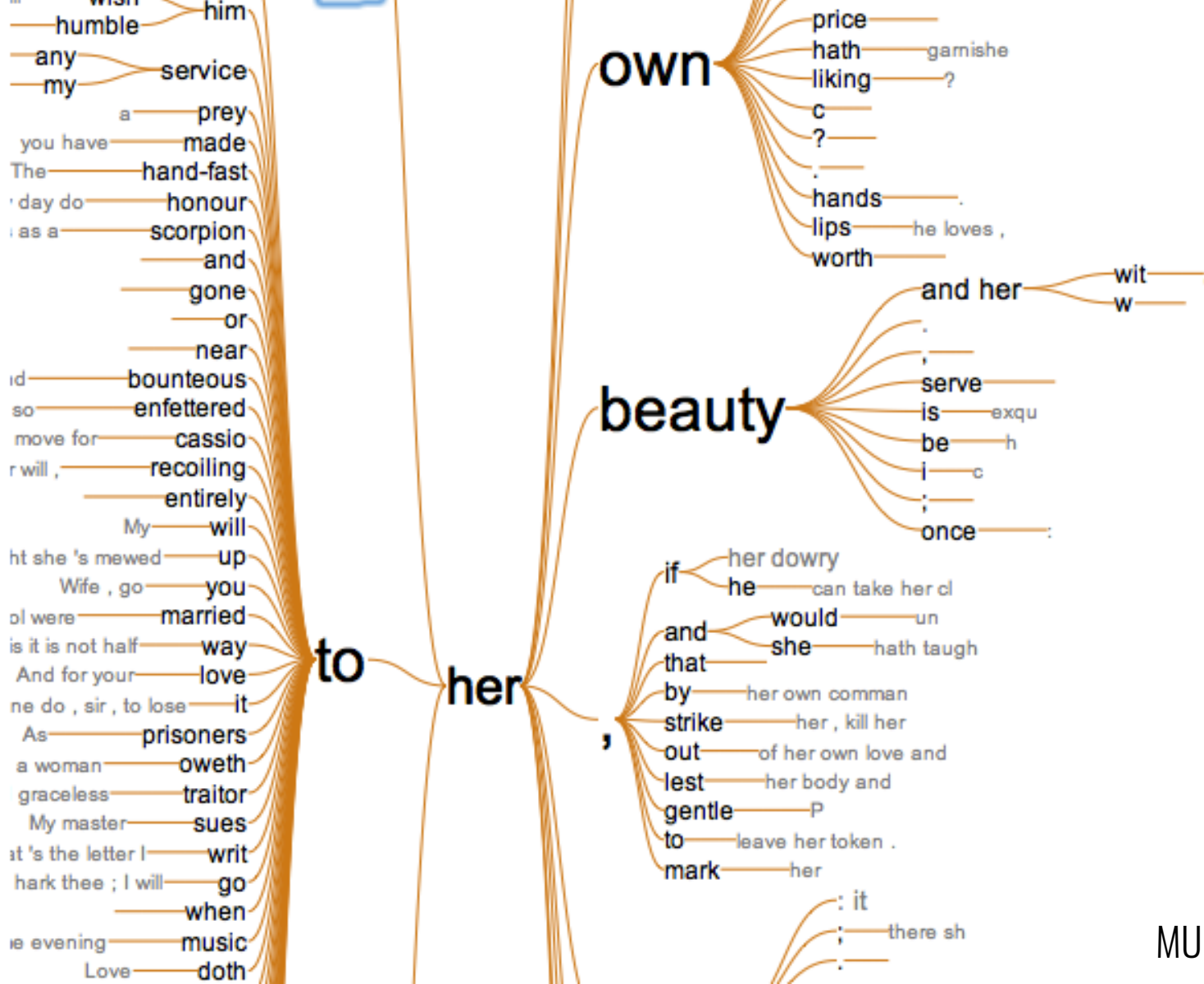- **brotherhood** .
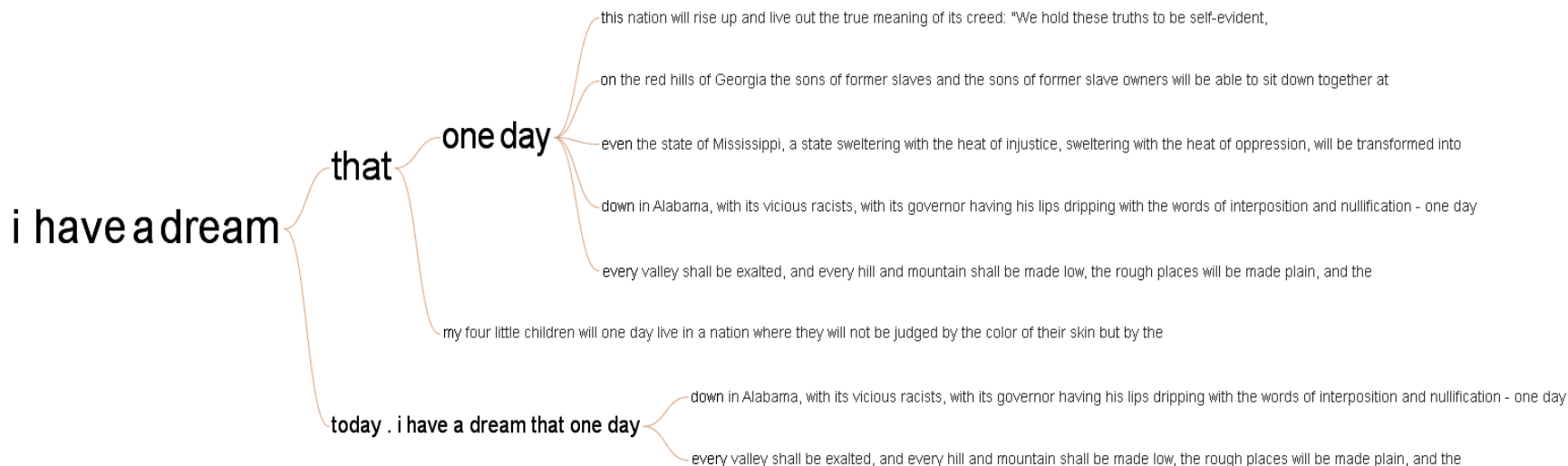- **world** , the love of the father is not in him .
- **brethren** .
- **children** of god , when we love god , and keep his commandments .

# FILTER INFREQUENT RUNS

# RECURRENT THEMES IN SPEECH



i have a dream

that

one day

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

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

down in Alabama, with its vicious racists, with its governor having his lips dripping with the words of interposition and nullification - 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

today . i have a dream that one day

down in Alabama, with its vicious racists, with its governor having his lips dripping with the words of interposition and nullification - 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

# GLIMPSES OF STRUCTURE

- CONCORDANCES SHOW LOCAL, REPEATED STRUCTURE

- BUT WHAT ABOUT OTHER TYPES OF PATTERNS?

- FOR EXAMPLE

- LEXICAL:     <A> at <B>

- SYNTACTIC:     <Noun> <Verb> <Object>

# PHRASE NETS

LOOK FOR SPECIFIC LINKING PATTERNS IN THE TEXT:

'A AND B', 'A AT B', 'A OF B', ETC

COULD BE OUTPUT OF REGEXP OR PARSER

VISUALIZE EXTRACTED PATTERNS IN A NODE-LINK VIEW

OCCURRENCES = NODE SIZE

PATTERN POSITION = EDGE DIRECTION

van Ham et al

**Select a phrase**

| word1 | and | word2 |
| word1 | 's | word2 |
| word1 | of the | word2 |
| word1 | the | word2 |
| word1 | a | word2 |
| word1 | at | word2 |
| word1 | is | word2 |
| word1 | [space] | word2 |

or enter your own

`* and *`   Submit

**Filters**

Show top: 100
Hide common words ☑

**Zoom**

In 🔍 Out 🔍 Reset ⊡

X and Y

flushed

coughing laughter

lips

mother

eyes   opened   made

repeated

voice   soul   heart

father   face   body

mr   head   neck   blood

children   heard   hair

door   gave

shame   rage

power   love

smiled   water

fleming   stephen

charles   turned   rose   fell

darkness   sweet   black

sad   blue

silence   soft   terrible   white

gloom   beautiful   long

air   cried   simple   strange   cruel

raised thought   high   low   gentle   pale   cold   unfair

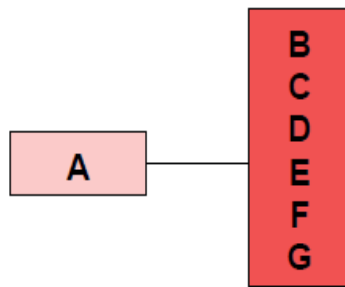red   holy   damp   dark
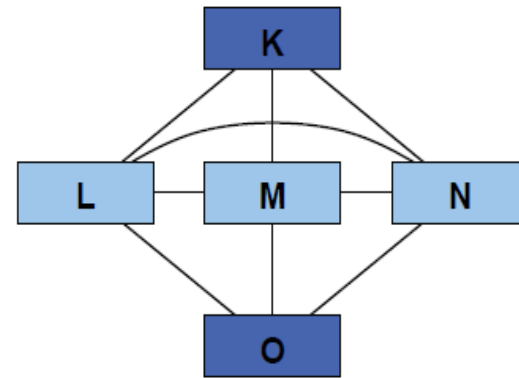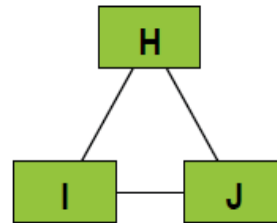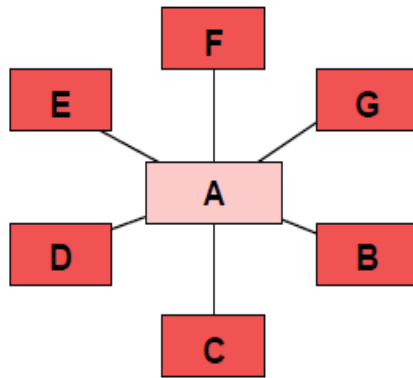
green   warm   silent

holly   young

small

weak

humble

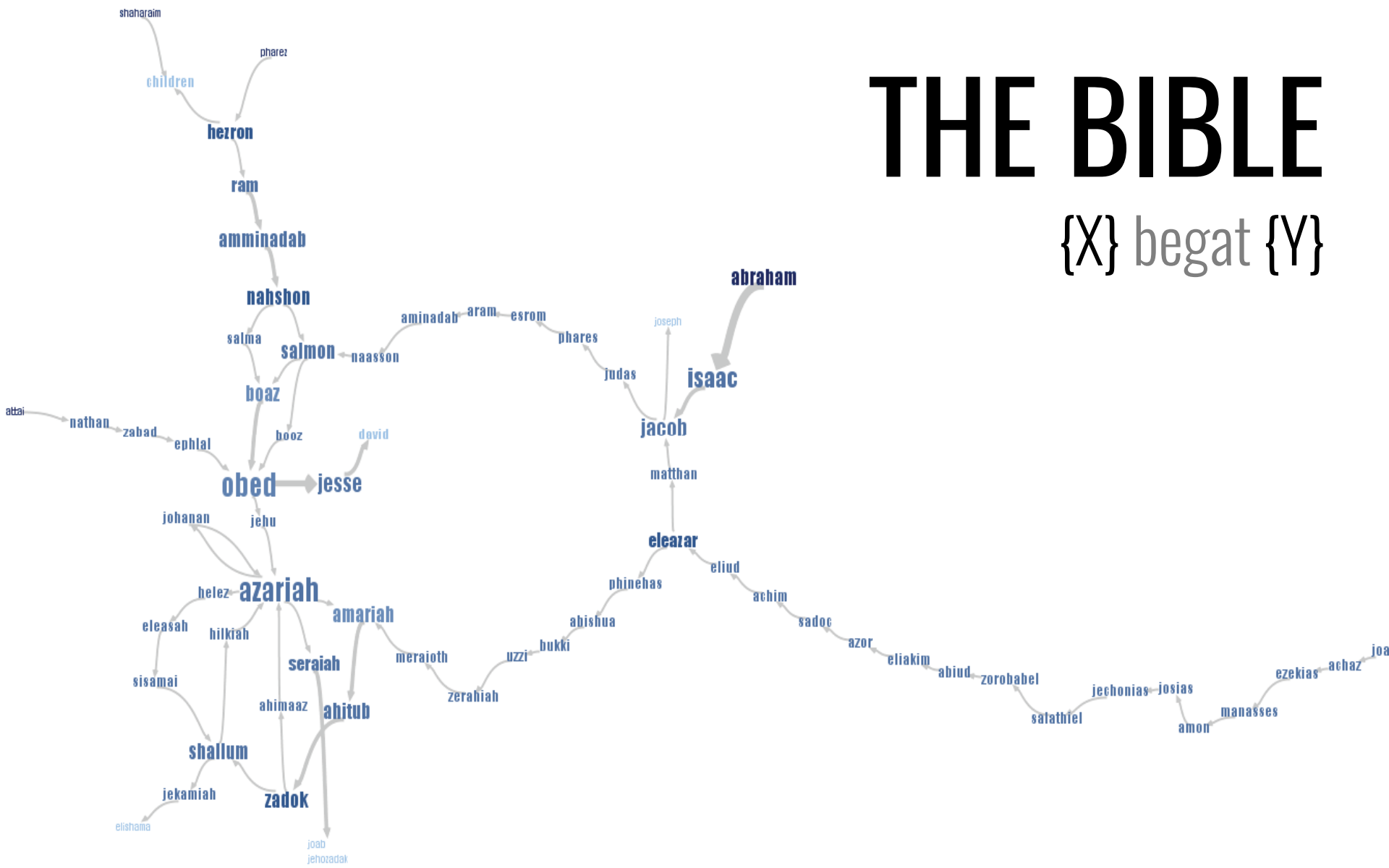PORTRAIT OF THE ARTIST AS A YOUNG MAN
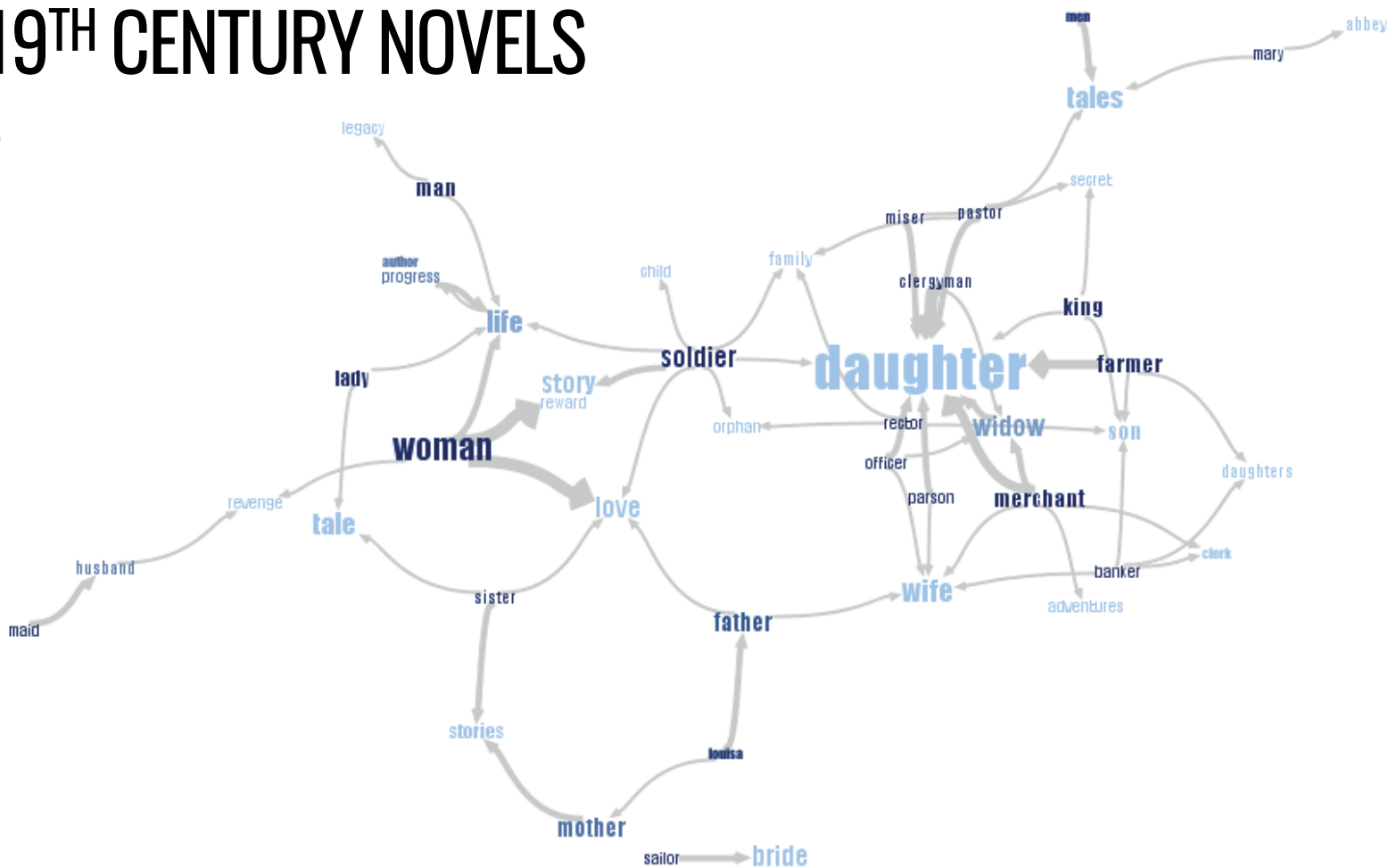JAMES JOYCE

# NODE GROUPING



(a)

(b)

(c)

# THE BIBLE

{X} begat {Y}

# 18TH & 19TH CENTURY NOVELS
{X}'s {Y}



legacy

man

author
progress

life

lady

woman

story
reward

revenge

tale

husband

maid

sister

stories

mother

sailor → bride

louisa

father

child

soldier

orphan

officer

parson

love

wife

family

miser    pastor

clergyman

rector

daughter

widow

son

merchant

banker

adventures

clerk

men        abbey
                mary

tales

secret

king

farmer

daughters

# OLD TESTAMENT
{X} of {Y}

# NEW TESTAMENT
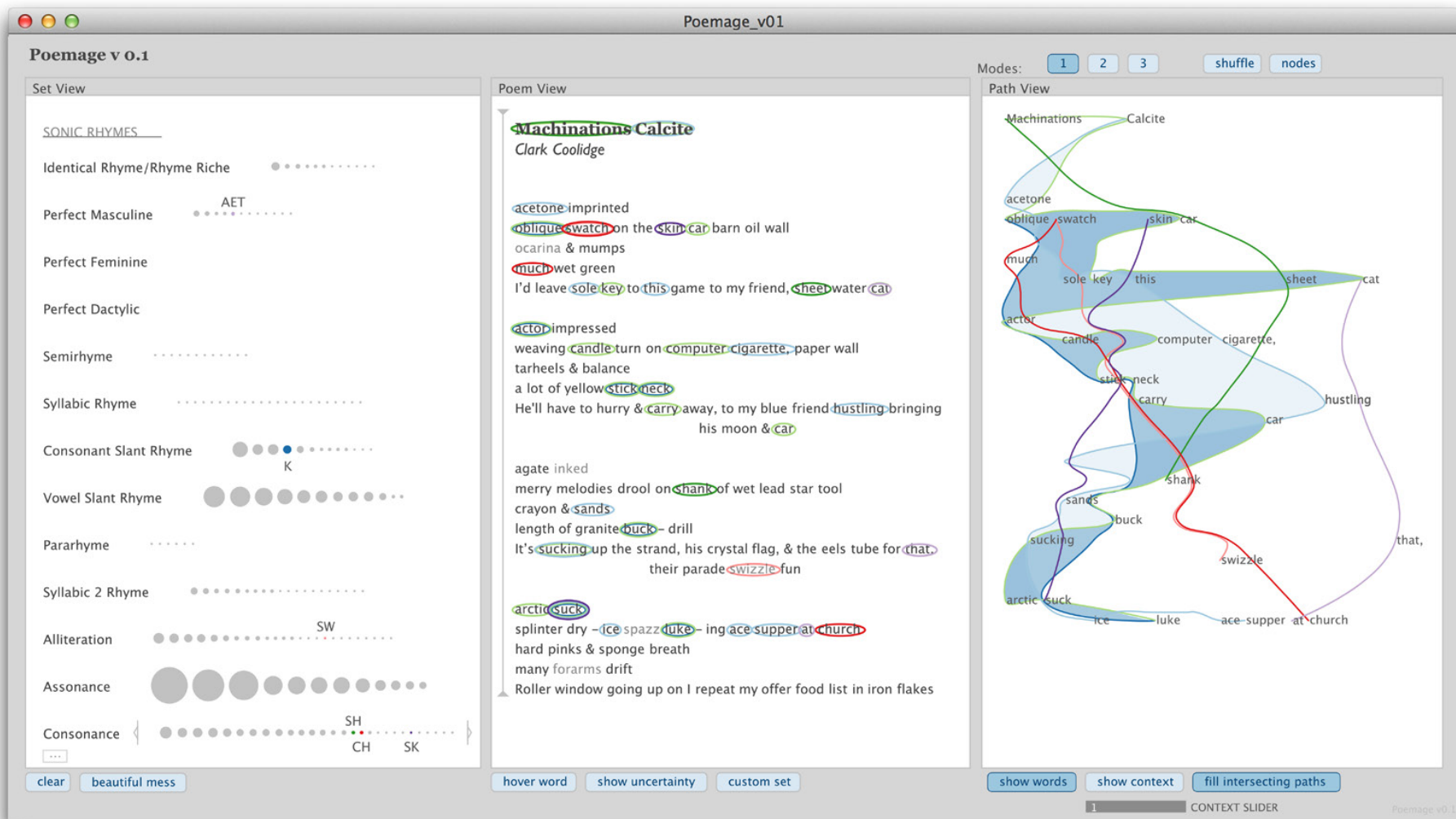{X} of {Y}

# RHYME, SPEECH, ETC.    POEMAGE McCurdy et al. 2016

# REVISIT YOUR SKETCHES?

## TASK:

**1)** VISUALIZE THE MOST IMPORTANT CONTENT FROM A SINGLE THESIS.

ARE YOUR VISUALIZATION CHOICES **EFFECTIVE**?

DOES THE VIS CAPTURE THE **LENGTH**, **FORM**, AND **POSITION** OF THE IMPORTANT CONTENT?

DO YOU SHOW OR CONNECT BACK TO THE **ORIGINAL TEXT**?

# EVOLVING DOCUMENTS

# VISUALIZING REVISION HISTORY

## HOW TO DEPICT CONTRIBUTIONS AND CHANGES OVER TIME?

# DIFF



svn diff: sshconsole.js

Diff style: [Side-by-side ▾]   ☐ Enable syntax coloring

**Files Changed:**

1. sshconsole.js: 1 change [ 1 ]

/home/toddw/src/sshconsole-read-only/content/sshconsole.js

```
...                          50 lines hidden [Expand]
51    _term = new VT100(80, 24, "term");          51    _term = new VT100(80, 24, "term");
52    //_term.debug_ = 1;                          52    //_term.debug_ = 1;
53    _term.curs_set(true, true, _term_box_element);  53    _term.curs_set(true, true, _term_box_element);
54    _term.noecho();                             54    _term.noecho();
55                                                55
56    // Replace the go_getch_ function with our own, this is called   56    // Replace the go_getch_ function with our own, this is called
57    // for every keypress that is passed through the terminal to the  57    // for every keypress that is passed through the terminal to the
58    // remote server. The character is already converted into the    58    // remote server. The character is already converted into the
59    // required VT100 character sequence(s).     59    // required VT100 character sequence(s).
60    VT100.go_getch_ = function() {              60    VT100.go_getch_ = function() {
61        var vt = VT100.the_vt_;                 61        var vt = VT100.the_vt_;
62        if (vt === undefined) {                 62        if (vt === somevalue) {
63            return;                             63            return;
64        }                                       64        }
65        var ch = vt.key_buf_.shift();           65        var ch = vt.key_buf_.shift();
66        //dump("go_getch_:: ch: '" + ch + "'\n");
67        if (ch === undefined) {                 66        if (ch === undefined) {
68            return;                             67            return;
69        }                                       68        }
70        if (vt.echo_ && ch.length == 1) {       69        if (vt.echo_ && ch.length == 1) {
71            vt.addch(ch);                       70            vt.addch(ch);
                                                  71            vt.refres();
72        }                                       72        }
73        if (_ssh_channel) {                     73        if (_ssh_channel) {
74            _ssh_channel.sendStdin(ch);         74            _ssh_channel.sendStdin(ch);
75        }                                       75        }
76    }                                           76    }
77                                                77
78    var serverTextbox = document.getElementById("sshconsole_server_textbox");  78    var serverTextbox = document.getElementById("sshconsole_server_textbox");
79    var connectionText;                         79    var connectionText;
80    if ('connectionText' in window.arguments[0]) {  80    if ('connectionText' in window.arguments[0]) {
81        connectionText = window.arguments[0].connectionText;  81        connectionText = window.arguments[0].connectionText;
82    } else {                                    82    } else {
...                          174 lines hidden [Expand]
```

WIKIPEDIA HISTORY FLOW VIÉGAS ET AL 2004

ON THE ORIGIN OF SPECIES *The Preservation of Favoured Traces*

I    II    III    IV    V    VI    VII    VIII    IX

head, not known to be electrical, but which
appears to be the real homologue of the electric
battery in the Torpedo. It is generally admitted
that there exists between these organs and
ordinary muscle a close analogy, in intimate

BEN FRY

I    II    III    IV    V    VI    VII    VIII    IX    X    XI    XII    XIII    XIV

First Edition (1859)    Second Edition (1860)    Third Edition (1861)    Fourth Edition (1866)    Fifth Edition (1869)    Sixth Edition (1872)

# Shortest Edit Path

## Edit War

# VISUALIZING DOCUMENT COLLECTIONS

# SKETCHING: **VISUALIZE**

IMAGINE YOU HAVE A MASTER'S THESIS IN FROM OF YOU:

YEAR
AUTHOR
TITLE
KEYWORDS
REFERENCES
**ABSTRACT TEXT**

TASK:

**1) VISUALIZE THE MOST IMPORTANT CONTENT** FROM A SINGLE THESIS.

**2) VISUALIZE HOW SIMILAR THIS THESIS IS** TO THESES FROM OTHER STUDENTS IN THE SAME MASTER'S PROGRAM.

(~10 MINUTES)

search all...

# Analysts: GOP may regret gridlock over Scalia replacement

# Update: Uber driver arrested in Michigan rampage that killed 6

## Boris Johnson backs EU exit: London mayor confirms support for Brexit

### 'A multifaceted catastrophe': Turkey has 'so alienated everyone it cannot service

### Blasts rock Syrian city of Homs, killing at least 32

### Palestinians struggle to define those who attack Israelis

**Canada, USA renew rivalry in CONCACAF final**

Sportsnet's James Sharman met with coach John Herdman and members of the Canadian women's soccer team, who are looking to beat the USA in Sunday's CONCACAF final.
headshot Gavin Day February 20, 2016, 8:08 PM.
headshot Gavin Day February

Feb 20  17:47  |  587 related articles  |  Sportsnet.ca

US rejected North Korea peace talks offer before last nuclear test

Malaysia, south-east Asia nations warned of terror attacks

San Bernardino victims to oppose Apple on iPhone encryption

# Samsung, LG unveil new devices in bid for smartphone recovery

## Raceline Radio Program Guide: February 21, 2016

## Canada, USA renew rivalry in CONCACAF final

### 'Deadpool' dominates again with $55 million in 2nd week

### Judge blocks attempt to halt deposition of Bill Cosby's wife

### Taylor Swift donates $250K to help Kesha's legal battle

Highlights from the USC report on entertainment diversity
WATCH: 'The Simpsons' mocks US presidential candidates in 'The Debateful Eight'

Kate Moss picks crutches with her trusty £200 Wellduk boots

Chhan, Plummer win at ACTRA Toronto Awards

Kate Hudson on her new book, 'mom buddies' Parr and Alba

Mani Paltrow confirms

2014 Sexiest Woman

Bunny days, Inside

## Chan wins Four Continents figure-skating championship

Years later, ex-Raptor Vince Carter's still soaring

SPRING TRAINING Blue Jays' focus at 2016 camp is on 2017

Scientists at Brock studying Zika to see if Canadian mosquitoes can spread the virus

How Syrian refugees arriving in Canada became 'extras' in their own stories

One dead, another injured, in avalanche near Golden

## Truex comes up a few inches short in closest Daytona 500

Canadian women earn historic 19-10 rugby victory over New Zealand

Miller puts an end to Canucks' losing streak

Kesler powers Ducks to 5-2 win over Flames

Event: Canadiens goalie steals

La Loche staff, students return to school this week

Closing arguments of Duffy trial to put spotlight on how Senate

Diabolical stabbing: Police arrest man book from prison

Kill airport attacks could rattle confidence in the Canadian

### LG Unveils the LG G5, Its First Modular Smartphone [Video]

LG G5 vs LG V10: first look

EPA asks Volkswagen to make electric cars in the US

Lennox Eyes Families, Biz Users With Budget Tab8 Tablets

What If San Bernardino Suspect Had Used an Android Instead of an

Fire Emblem Fates Marriage Guide - Best Children, Pairs, Stats, Ratings

HSBC To Bring Israeli

ZTE And China

New York

Leafs get set for a busy draft with Matthias trade

Now pitcher FA Dickey axes Blue Jays' playoff payoff

Myra

L D

Serial killer Robert Pickton pens book from prison

Yanis our Everyday

Tiny Banking

Binge drinking youngsters risk developing high blood pressure and heart disease

HPV cases drop since vaccinations started

Calgary

Halton Hills

AB

Vaugh north

Asian stocks rebound in anticipation of G20 meeting

Five things to watch for in the Canadian

Off Poland and World

newsmap.jp

# DOCUMENT CARDS
## SMALL MULTIPLES FOR DOCUMENTS



Cerebral: Visualizing Multiple Experimental Conditions on a Graph with Biological Context

Aaron Barsky, Tamara Munzner, Jennifer Gardy, and Robert Kincaid

Multi-Focused Geospatial Analysis Using Probes

Thomas Butkiewicz, Wenwen Dou, Zachary Wartell, William Ribarsky, and Remco Chang

Stacked Graphs Geometry & Aesthetics

Lee Byron and Martin Wattenberg

Vispedia : Interactive Visual Exploration of Wikipedia Data via Search-Based Integration

Bryan Chan, Leslie Wu, Justin Talbot, Mike Cammarano, and Pat Hanrahan

Geometry-Based Edge Clustering for Graph Visualization

Weiwei Cui, Hong Zhou, Student , Huamin Qu, Pak Chung Wong, and Xiaoming Li

VisGets: Coordinated Visualizations for Web-based Information Exploration and Discovery

Marian Dörk, Sheelagh Carpendale, Christopher Collins, and Carey Williamson

Who Votes For What? A Visual Query Language for Opinion Data

Geoffrey M. Draper, and Richard F. Riesenfeld

Exploration of Networks Using Overview+Detail with Constraint-based Cooperative Layout

Tim Dwyer, Kim Marriott, Falk Schreiber, Peter J. Stuckey, Michael Woodward and Michael Wybrow

Rolling the Dice: Multidimensional Visual Exploration using Scatterplot Matrix Navigation

Niklas Elmqvist, Pierre Dragicevic, and Jean-Daniel Fekete

Interactive Visual Analysis of Set-Typed Data

Wolfgang Freiler, Kresimir Matkovi`c, Computer Society, and Helwig Hauser

Graphical Histories for Visualization: Supporting Analysis, Communication, and Evaluation

Jeffrey Heer, Jock D. Mackinlay, Chris Stolte, and Maneesh Agrawala

Improving the Readability of Clustered Social Networks using Node Duplication

Nathalie Henry, Anastasia Bezerianos, and Jean-Daniel Fekete

THEMERIVER HAVRE ET AL 1999

Nationalization of property begins
Castro bans religious TV and radio

Cuba and Soviet relations resume

US imposes embargo on Cuba

Bay of Pigs

weapons(62)
troops(53)
reform(48)
harvest(26)
church(13)

tourism(52)

oil(44)

yankee(63)

soviet(49)

imperialists(29)

kennedy(32)

cooperatives(16)

cane(7)

Nov Dec Jan Feb Mar Apr May Jun Jul Aug Sep Oct Nov Dec Jan Feb Mar Apr May

# PARALLEL TAG CLOUDS

09



| First | Second | Third | Fourth | Fifth | Sixth | Seventh | Eighth | Ninth | Tenth | Eleventh | Federal | DC |

# SUPPORTING SEARCH



**TileBars** Hearst 1999

/tmp/words22058

conscience

angel

adultery

3-

2-

1-

gen
exo
lev
deu
jos
rom
co1
col
gal
eph
phi
col
th1
th2
ti1
ti2
tit
phm
heb
jam
pe1
pe2
jo1
jo2
jo3
jud
rev

SeeSoft Eick 1994

# The 2007 State of the Union Address

Over the years, President Bush's State of the Union address has averaged almost 5,000 words each, meaning the the President has delivered over 34,000 words. Some words appear frequently while others appear only sporadically. Use the tools below to analyze what Mr. Bush has said.

Search | or choose a word here.

## Use of the phrase "Tax" in past State of the Union Addresses

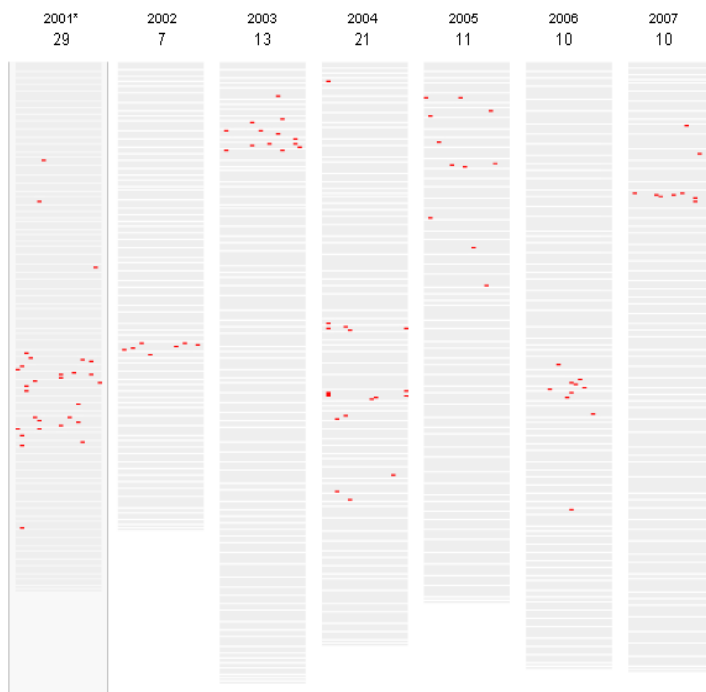| 2001* | 2002 | 2003 | 2004 | 2005 | 2006 | 2007 |
|-------|------|------|------|------|------|------|
| 29 | 7 | 13 | 21 | 11 | 10 | 10 |

### The word in context

Next Instance of 'Tax'
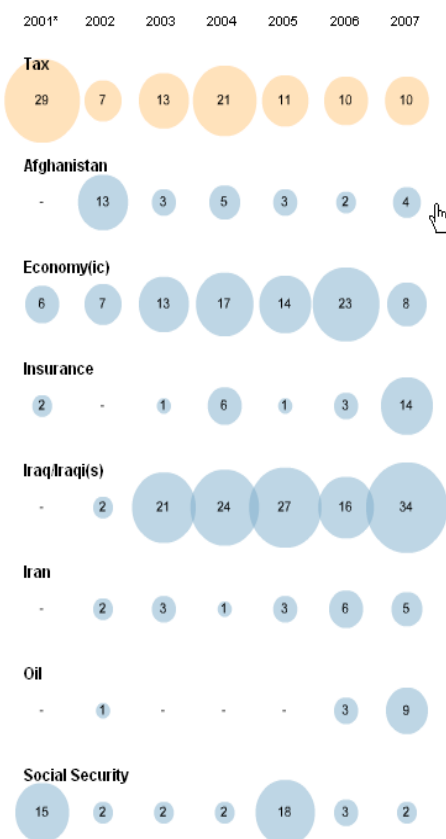
I believe in local control of schools. We should not, and we will not, run public schools from Washington, D.C. Yet when the federal government spends **TAX** dollars, we must insist on results. Children should be tested on basic reading and math skills every year between grades three and eight. Measuring is the only way to know whether all our children are learning. And I want to know, because I refuse to leave any child behind in America.

-- 2001 (Paragraph 14 of 73)

## Compared with other words

| | 2001* | 2002 | 2003 | 2004 | 2005 | 2006 | 2007 |
|------|-------|------|------|------|------|------|------|
| **Tax** | 29 | 7 | 13 | 21 | 11 | 10 | 10 |
| **Afghanistan** | - | 13 | 3 | 5 | 3 | 2 | 4 |
| **Economy(ic)** | 6 | 7 | 13 | 17 | 14 | 23 | 8 |
| **Insurance** | 2 | - | 1 | 6 | 1 | 3 | 14 |
| **Iraq/Iraqi(s)** | - | 2 | 21 | 24 | 27 | 16 | 34 |
| **Iran** | - | 2 | 3 | 1 | 3 | 6 | 5 |
| **Oil** | - | 1 | - | - | - | 3 | 9 |
| **Social Security** | 15 | 2 | 2 | 2 | 18 | 3 | 2 |

* As a newly elected president, Mr. Bush did not deliver a formal State of the Union address in 2001. His Feb. 27 speech to a joint session of Congress was analogous to the State of the Union, but without the title.

# DOCUMENT SIMILARITY & CLUSTERING

**COMPUTE SIMILARITY** BETWEEN DOCUMENTS BASED ON THE WORDS THEY SHARE

- **TF-IDF** (TERM FREQUENCY-INVERSE DOCUMENT FREQUENCY) IS COMMON

**TOPIC MODELING** APPROACHES

- ASSUME DOCUMENTS ARE A MIXTURE OF TOPICS
- TOPICS ARE (ROUGHLY) A SET OF CO-OCCURRING TERMS
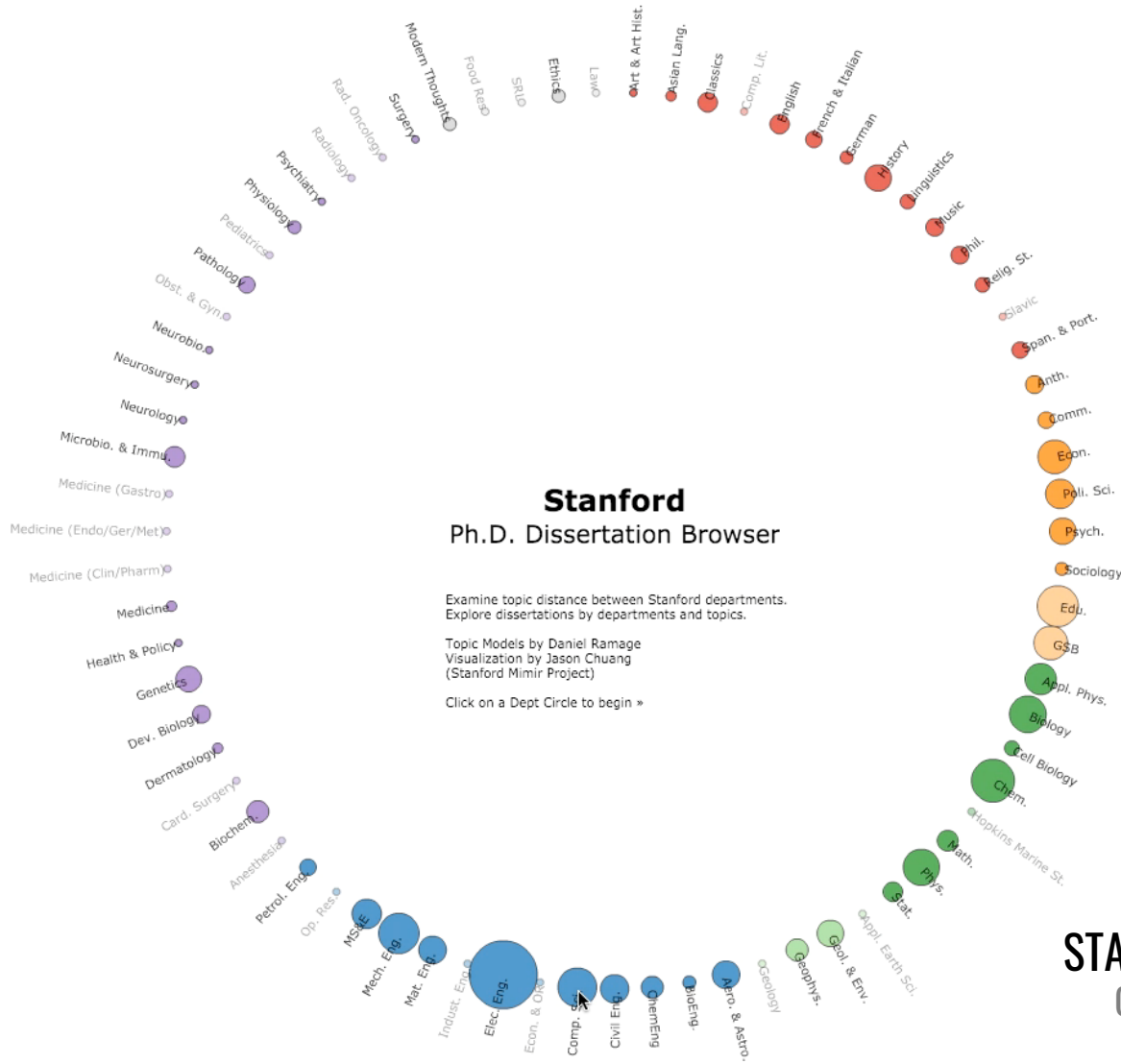- LATENT SEMANTIC ANALYSIS (LSA): REDUCE TERM MATRIX

- MANY, MANY APPROACHES EXIST
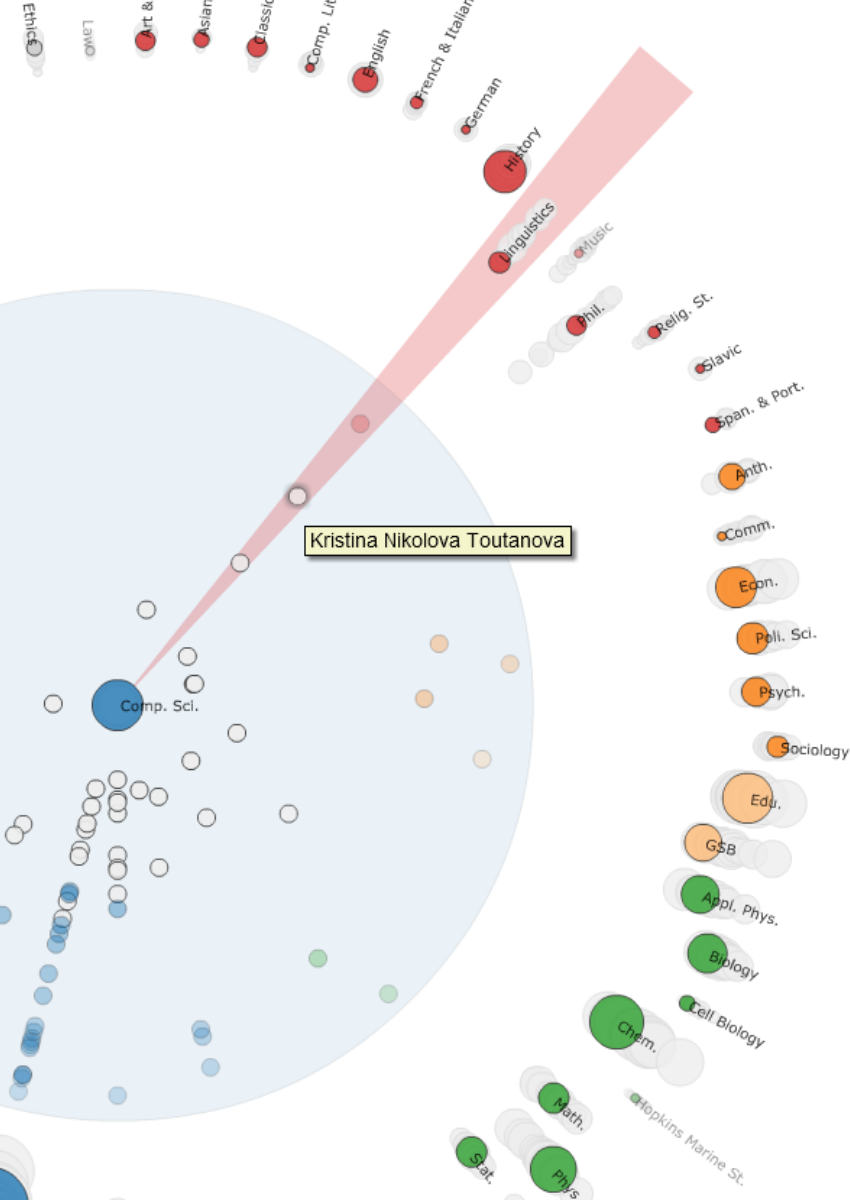
**Stanford**
Ph.D. Dissertation Browser

Examine topic distance between Stanford departments.
Explore dissertations by departments and topics.

Topic Models by Daniel Ramage
Visualization by Jason Chuang
(Stanford Mimir Project)

Click on a Dept Circle to begin »

# STANFORD DISSERTATION BROWSER
## CHUANG, RAMAGE, MANNING & HEER 2012

**Effective statistical models for syntactic and semantic disambiguation**

Student: Kristina Nikolova Toutanova
Advisor: Christopher D. Manning

Computer Science (2005)

Keywords: Syntactic, Semantic, Tree kernels, Parsing

Abstract:

This thesis focuses on building effective statistical models for disambiguation of sophisticated syntactic and semantic natural language (NL) structures. We advance the state of the art in several domains by (i) choosing representations that encode domain knowledge more effectively and (ii) developing machine learning algorithms that deal with the specific properties of NL disambiguation tasks--sparsity of training data and large, structured spaces of hidden labels. For the task of syntactic disambiguation, we propose a novel representation of parse trees that connects the words of the sentence with the hidden syntactic structure in a direct way. Experimental evaluation on parse selection for a Head Driven Phrase Structure Grammar shows the new representation achieves superior performance compared to previous models. For the task of disambiguating the semantic role structure of verbs, we build a more accurate model, which captures the knowledge that the semantic frame of a verb is a joint structure with strong dependencies between arguments. We achieve this using a Conditional Random Field without Markov independence assumptions on the sequence of semantic role labels. To address the sparsity problem in machine learning for NL, we develop a method for incorporating many additional sources of information, using Markov chains in the space of words. The Markov chain framework makes it possible to combine multiple knowledge sources, to learn how much to trust each of them, and to chain inferences together. It achieves large gains in the task of disambiguating prepositional phrase attachments.

# STANFORD DISSERTATION BROWSER
## CHUANG, RAMAGE, MANNING & HEER 2012

OFTEN, TEXT VISUALIZATIONS DO NOT REPRESENT TEXT DIRECTLY, BUT THEY REPRESENT A MODEL

WORD COUNTS, WORD SEQUENCES, CLUSTERS, ETC.

ASK:

CAN YOU INTERPRET THE VISUALIZATION?

DOES THE MODEL ACCURATELY REPRESENT THE ORIGINAL TEXT?

# LESSONS FOR TEXT VISUALIZATION

**SHOW SOURCE TEXT** (OR PROVIDE ACCESS TO IT)

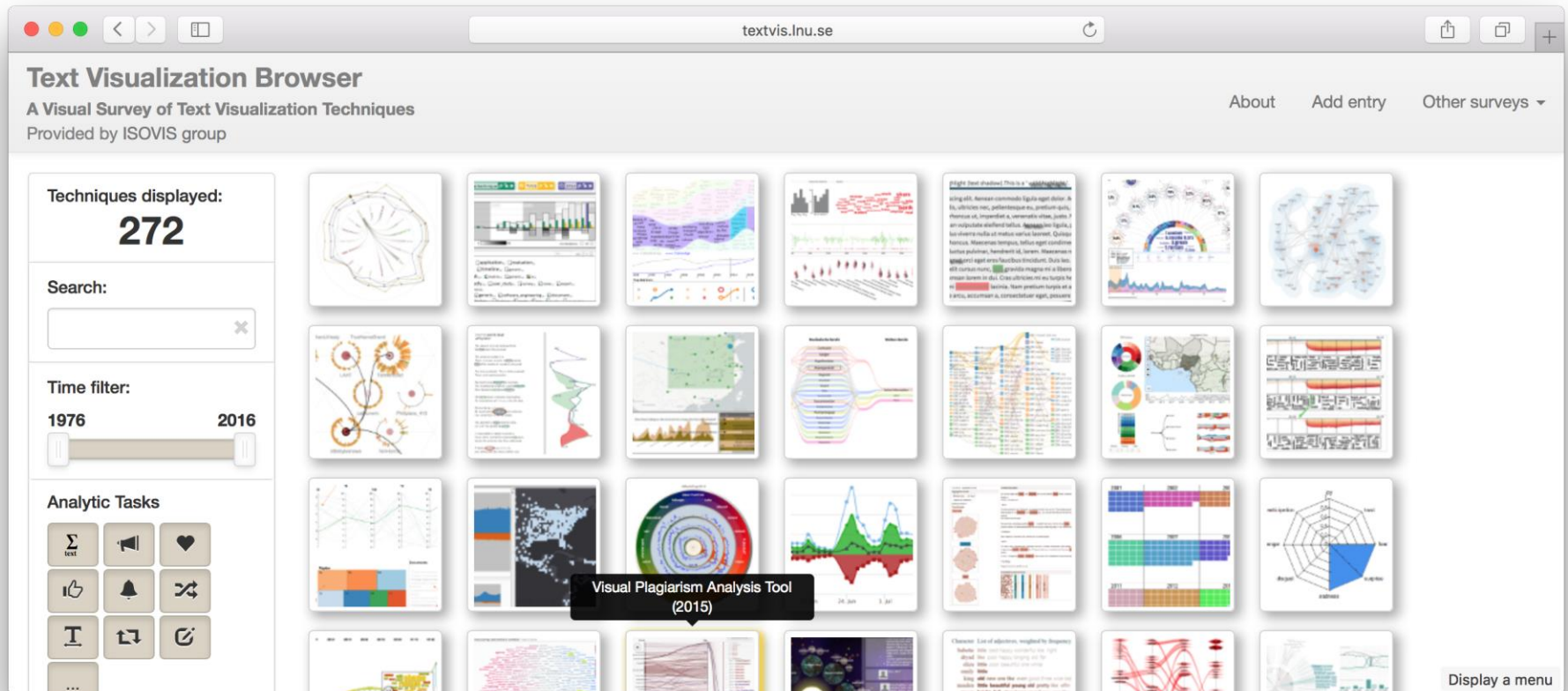WHERE POSSIBLE, USE VISUALIZATION AS INDEX INTO DOCUMENTS

GROUP DOCUMENTS IN MEANINGFUL WAYS

WILL VIEWERS UNDERSTAND THE CLUSTERS?
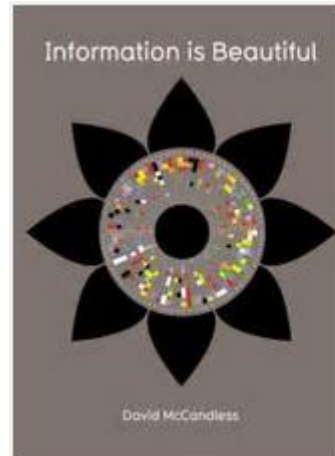
WHERE POSSIBLE **USE TEXT TO REPRESENT TEXT**

# HUNDREDS OF TOOLS & TECHNIQUES FOR TEXT AT
## http://textvis.lnu.se/

Getting to know David McCandless

HTTP://WWW.WEFEELFINE.ORG/

# QUESTIONS?

# ACKNOWLEDGEMENTS

Slides in were inspired, adapted, taken from slides by

- Christopher Collins (University of Ontario Institute of Technology)
- Wesley Willett (University of Calgary)