# TEXT VISUALIZATION

CS 448B | Fall 2024

MANEESH AGRAWALA

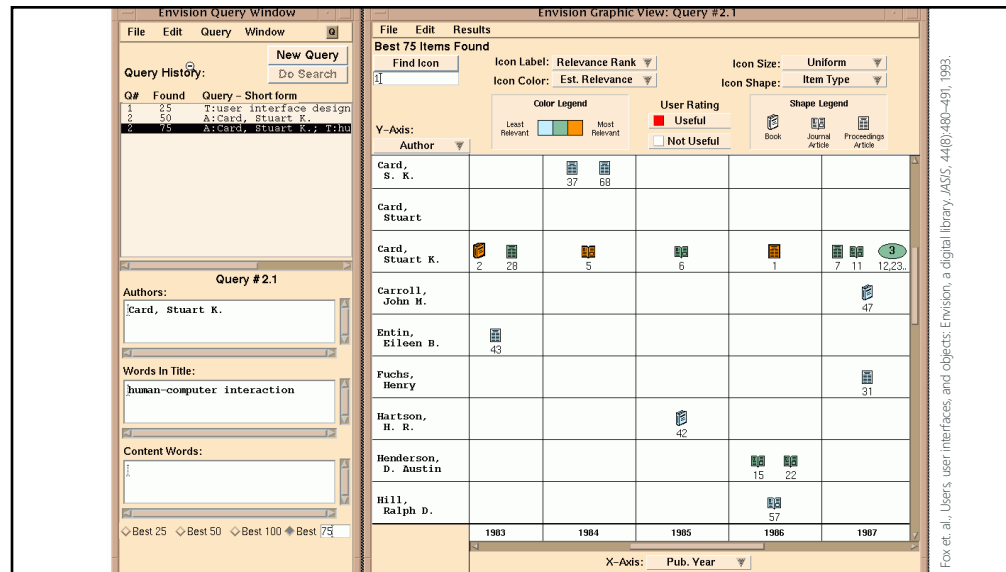1

# READING RESPONSE: QUESTIONS/THOUGHTS

While a method like ReVision (created in a research environment) would likely not be widely used until its accuracy improves significantly, I wonder if there are already AI chart generation or redesign tools being used that might not have perfect accuracy (and that are not subject to as much scrutiny as university research tools). ***Who is responsible for mistakes or misinformation in published charts? Is it the data collector, the creator of the visualization, or the editor?***

2

Tag clouds are the new mullets

**NiemanLab**

BUSINESS MODELS | MOBILE & APPS | AUDIENCE & SOCIAL | AGGREGATION & DI

**Word clouds considered harmful**

3

---

# SINGLE DOCUMENTS AND COLLECTIONS

**Documents**
- Articles, books, novels
- Computer programs
- Email, Web pages, blogs
- Tags, comments

**Collections of documents**
- Messages (e-mail, etc.)
- Social nets (profiles)
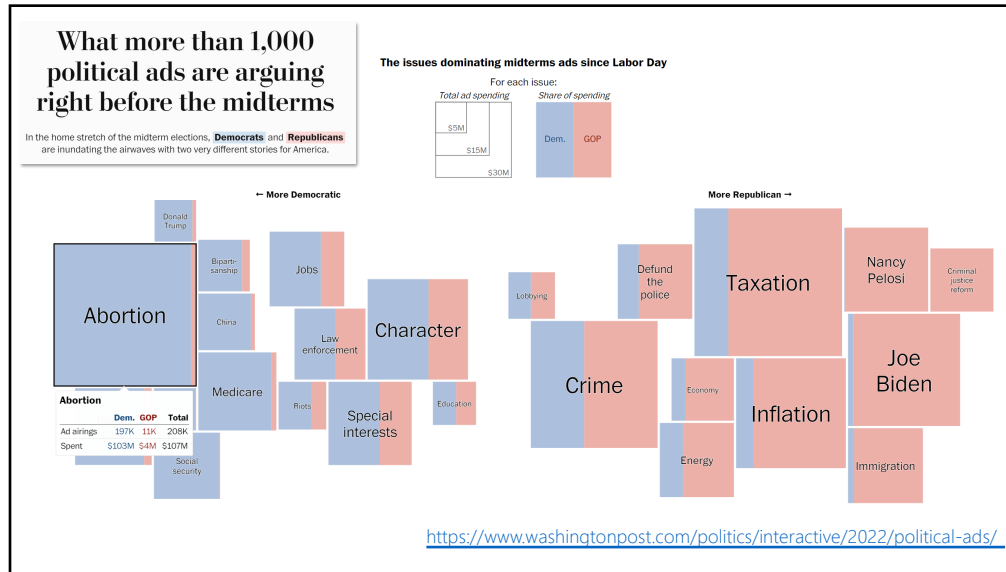- Academic collaborations (publications)



4

# TASKS

- What documents contain text about X?
- Which documents are of interest to me?
- Are there documents similar to this one?
- How are different words used in the collection?
- What are the main themes?
- How are themes distributed?
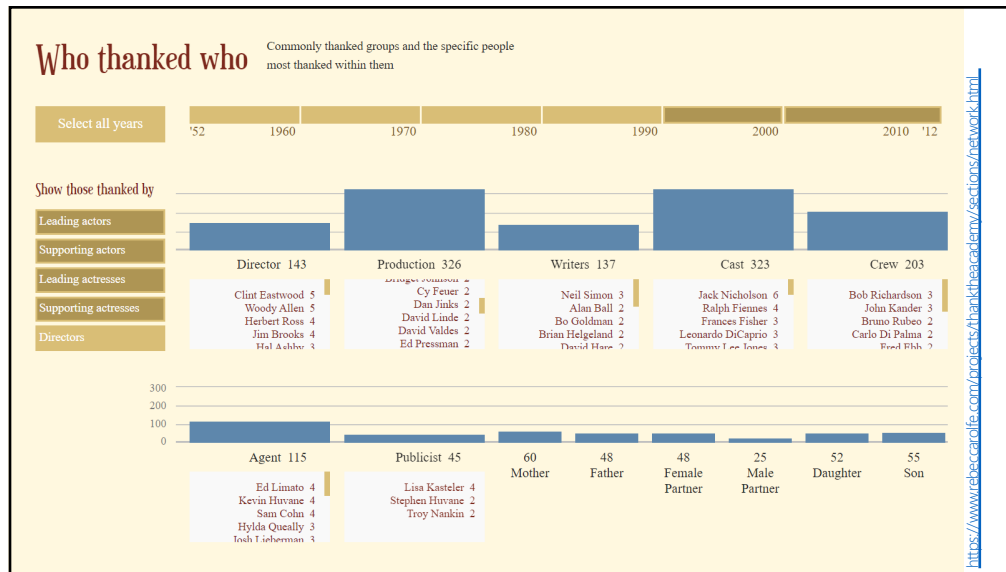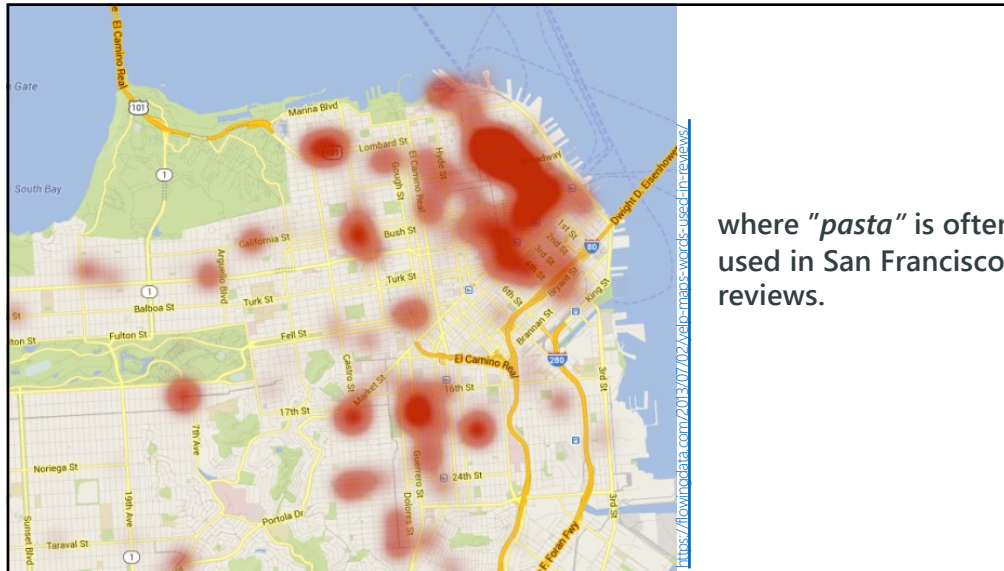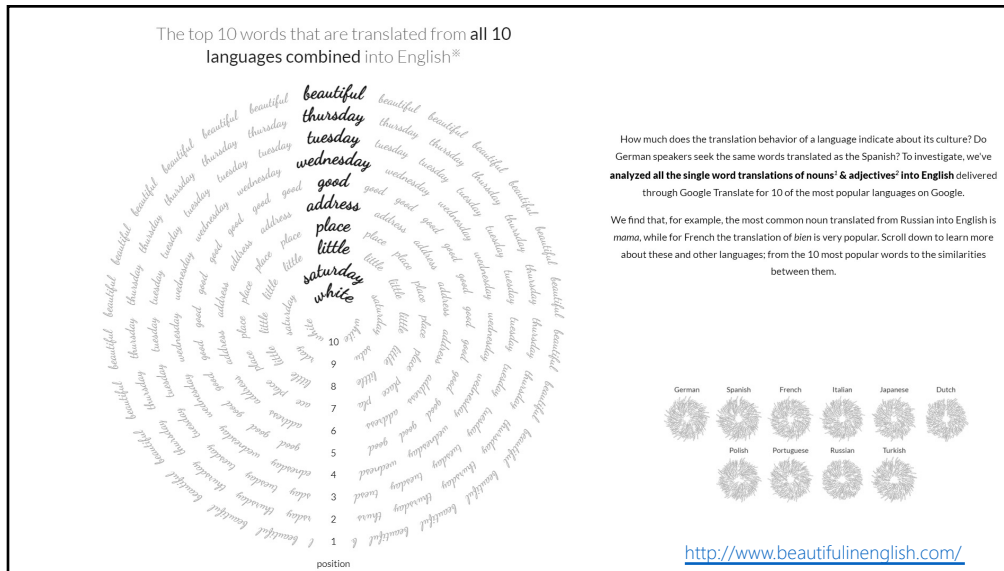- What unifies/connects documents in a collection?

5



6

## Slide 7

**What more than 1,000 political ads are arguing right before the midterms**

In the home stretch of the midterm elections, **Democrats** and **Republicans** are inundating the airwaves with two very different stories for America.

The issues dominating midterms ads since Labor Day

For each issue:

*Total ad spending*  *Share of spending*

$5M
$15M
$30M

Dem.  GOP

← More Democratic

More Republican →

Donald Trump

Abortion

Biparti-sanship

China

Jobs

Defund the police

Taxation

Nancy Pelosi

Criminal justice reform

Lobbying

Law enforcement

Character

Medicare

Crime

Economy

Inflation

Joe Biden

Riots

Special interests

Education

Energy

Immigration

Social security

**Abortion**

|  | Dem. | GOP | Total |
|---|---|---|---|
| Ad airings | 197K | 11K | 208K |
| Spent | $103M | $4M | $107M |

https://www.washingtonpost.com/politics/interactive/2022/political-ads/

7

## Slide 8

**Who thanked who**

Commonly thanked groups and the specific people most thanked within them

Select all years

'52   1960   1970   1980   1990   2000   2010  '12

Show those thanked by

Leading actors
Supporting actors
Leading actresses
Supporting actresses
Directors

**Director 143**
Clint Eastwood 5
Woody Allen 5
Herbert Ross 4
Jim Brooks 4
Hal Ashby 3

**Production 326**
Bridget Johnson 2
Cy Feuer 2
Dan Jinks 2
David Linde 2
David Valdes 2
Ed Pressman 2

**Writers 137**
Neil Simon 3
Alan Ball 2
Bo Goldman 2
Brian Helgeland 2
David Hare 2

**Cast 323**
Jack Nicholson 6
Ralph Fiennes 4
Frances Fisher 3
Leonardo DiCaprio 3
Tommy Lee Jones 3

**Crew 203**
Bob Richardson 3
John Kander 3
Bruno Rubeo 2
Carlo Di Palma 2
Fred Ebb 2

300
200
100
0

**Agent 115**
Ed Limato 4
Kevin Huvane 4
Sam Cohn 4
Hylda Queally 3
Josh Lieberman 3

**Publicist 45**
Lisa Kasteler 4
Stephen Huvane 2
Troy Nankin 2

60 Mother
48 Father
48 Female Partner
25 Male Partner
52 Daughter
55 Son

https://www.rebeccarolfe.com/projects/thanktheacademy/sections/network.html

8

where *"pasta"* is often used in San Francisco reviews.

https://flowingdata.com/2013/07/02/yelp-maps-words-used-in-reviews/

9



The top 10 words that are translated from **all 10 languages combined** into English*

beautiful
thursday
tuesday
wednesday
good
address
place
little
saturday
white

10
9
8
7
6
5
4
3
2
1

position

How much does the translation behavior of a language indicate about its culture? Do German speakers seek the same words translated as the Spanish? To investigate, we've **analyzed all the single word translations of nouns[1] & adjectives[2] into English** delivered through Google Translate for 10 of the most popular languages on Google.

We find that, for example, the most common noun translated from Russian into English is *mama*, while for French the translation of *bien* is very popular. Scroll down to learn more about these and other languages; from the 10 most popular words to the similarities between them.

German    Spanish    French    Italian    Japanese    Dutch

Polish    Portuguese    Russian    Turkish

http://www.beautifulinenglish.com/

12

## TODAY

### Learning Objectives

1. Considering text as data

2. Identifying descriptive words/keyphrases

3. Visualizing text in context

4. Searching across and comparing documents

13

# TEXT AS DATA

14

# TEXT AS NOMINAL DATA

**High dimensional (many words, e.g. 10,000+ common words)**

**More than equality tests**
- Correlations: *Hong Kong, San Francisco, Bay Area*
- Order: *April, February, January, June, March, May*
- Membership: *Tennis, Running, Swimming, Hiking, Piano*
- Hierarchy, antonyms & synonyms, entities, ...

**Words have meanings and relations**

15

# TEXT PROCESSING PIPELINE

**Tokenization**
Segment text into terms
Remove stop words? *a, an, the, of, to, be*
Numbers and symbols? *#cardinal, @Stanford, OMG!!!!!!!!*
Entities extraction? *Palo Alto, O'Connor, U.S.A.*

16

# TEXT PROCESSING PIPELINE

**Tokenization**
> Segment text into terms
> Remove stop words?  *a, an, the, of, to, be*
> Numbers and symbols?  *#cardinal, @Stanford, OMG!!!!!!!!*
> Entities extraction?  *Palo Alto, O'Connor, U.S.A.*

**Stemming**
> Group together different forms of a word
> Porter stemmer? visualization(s), visualize(s), visually → visual
> Lemmatization?  goes, went, gone → go

17

# TEXT PROCESSING PIPELINE

**Tokenization**
> Segment text into terms
> Remove stop words?  *a, an, the, of, to, be*
> Numbers and symbols?  *#cardinal, @Stanford, OMG!!!!!!!!*
> Entities extraction?  *Palo Alto, O'Connor, U.S.A.*

**Stemming**
> Group together different forms of a word
> Porter stemmer? visualization(s), visualize(s), visually → visual
> Lemmatization?  goes, went, gone → go

**Ordered list of terms**

18

# IDENTIFYING DESCRIPTIVE WORDS/KEYPHRASES

20

# BAG OF WORDS MODEL

**Ignore ordering relationships within the text**

**A document ≈ vector of term weights**
- Each term corresponds to a dimension (10,000+)
- Each value represents the relevance
  - For example, simple term counts

**Aggregate into a document x term matrix**
- Document vector space model

21

# DOCUMENT X TERM MATRIX

Each document is a vector of term weights
Simplest weighting is to just count occurrences

|  | Antony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth |
|---|---|---|---|---|---|---|
| Antony | 157 | 73 | 0 | 0 | 0 | 0 |
| Brutus | 4 | 157 | 0 | 1 | 0 | 0 |
| Caesar | 232 | 227 | 0 | 2 | 1 | 1 |
| Calpurnia | 0 | 10 | 0 | 0 | 0 | 0 |
| Cleopatra | 57 | 0 | 0 | 0 | 0 | 0 |
| mercy | 2 | 0 | 3 | 5 | 5 | 1 |
| worser | 2 | 0 | 1 | 1 | 1 | 0 |

22



24

25



27

# WORD/TAG CLOUDS

**Strengths**
- Compact – lots of words fit
- Draws eye to most frequent/biggest words
- Can help with gisting and initial query formation

**Weaknesses**
- Sub-optimal visual encoding (size not pos. encodes freq.)
- Inaccurate size encoding (long words are bigger)
- May not facilitate comparison (unstable layout)
- Term frequency may not be meaningful
- Does not show the structure of the text

28

# WORD WEIGHTING

**Term Frequency**
$tf_{td}$ = count(t) in document

**TF.IDF: Term Frequency by Inverse Document Freq**
$tf.idf_{td} = \log(1 + tf_{td}) \times \log(N/df_t)$

$df_t$ = # documents containing term t
N = # of documents

29

tf-idf: the effect of idf on tf

Comparison of two measures of word importance: tf (term frequency) and tf-idf (term frequency-inverse document frequency). Ranks are positions relative to other words; the highest rank is 1, which means the highest importance, the lowest rank is 271.

Source: transcript of Andy Cotgreave's interview in the "learning dataviz" episode of #SWDpodcast.
Illustration by @DmitrijsKass in R with ggplot2.

30

# LIMITATIONS OF FREQUENCY STATISTICS

**Typically focus on unigrams (single terms)**

**Often favors frequent (TF) or rare (IDF) terms**
Not clear that these provide best description of documents

**"Bag of words" ignores additional info**
Grammar / part-of-speech
Position within document
Recognizable entities

31

13

# YELP REVIEW SPOTLIGHT



[Yatani 2011]

32

# YELP REVIEW SPOTLIGHT



[Yatani 2011]

33

# TIPS FOR DESCRIPTIVE KEYPHRASES

**Understand the limitations of your language model**
Bag of words:
Easy to compute
Single words
Loss of word ordering

**Select appropriate model and visualization**
Generate longer, more meaningful phrases
Adjective-noun word pairs for reviews

34

# VISUALIZING WORDS/KEYPHRASES CONTEXT

35

# CONCORDANCE

# WORD TREE



The Word Tree, An Interactive Visual Concordance [pdf]
Martin Wattenberg & Fernanda Viégas, InfoVis 2008.

The Word Tree, An Interactive Visual Concordance [pdf]
Martin Wattenberg & Fernanda Viégas, InfoVis 2008.

38

# FILTER INFREQUENT RUNS



The Word Tree, An Interactive Visual Concordance [pdf]
Martin Wattenberg & Fernanda Viégas, InfoVis 2008.

39

46



*Alice's Adventures In Wonderland*

48

Alice's Adventures In Wonderland

49



Alice's Adventures In Wonderland

50

51



52

# GLIMPSES OF STRUCTURE

**Concordances and TextArc show local, repeated structure**
**But what about other types of patterns?**

**For example**

| | |
|---|---|
| Lexical: | \<A\> at \<B\> |
| Syntactic: | \<Noun\> \<Verb\> \<Object\> |

54

# 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
Darker color → higher ratio of out-edges to in-edges

Mapping Text with Phrase Nets
Frank Van Ham, Martin Wattenberg & Fernanda Viégas, InfoVis 2009.

55

Portrait of the Artist as a Young Man
X and Y

# EDGE COMPRESSION/NODE GROUPING

The Bible
X begat Y

58



Pride & Prejudice
X at Y

59

18th & 19th Century Novels
X's Y

62



Old Testament
X of Y

63

New Testament
X of Y

64



# ANNOUNCEMENTS

65

25

# FINAL PROJECT
## Design Reviews Dec 2 and Dec 4 (signups this week)

**Data analysis/explainer**
Analyze dataset in depth & make a visual explainer

**Deliverables**
An article with multiple different interactive visualizations
Short video (2 min) demoing and explaining the project

**Schedule**
Design Review and Feedback: 10th week of quarter, 12/2 and 12/4
Final code and video: Sun 12/8 8pm

**Grading**
Groups of up to 3 people, graded individually
Clearly report responsibilities of each member

66

# DESIGN REVIEW SIGNUPS

**Sign up for 8-10 min slot with teaching team**
Will offer slots on Mon 12/2 & Wed 12/4 during class period and perhaps others
Stay tuned for a canvas announcement about signup sheet

67

# SEARCHING ACROSS DOCUMENT(S)

68

**What's wrong with search results?**

69

# PROBLEMS WITH SEARCH RESULTS

Query responses don't tell you:
- How strong the match
- How frequent each term
- How term is distributed
- Overlap between terms
- Length of document

Ranking is opaque

Inability to compare results

Marti A. Hearst

70



71

VQuery
Jones & McInnes, UIST '98

72



Query Term
Suggestion
Visualization

Quintura

73

29

**Objective: Minimize time to decide which documents to look at**

TileBars, Marti A. Hearst

74



Document Length

Term 1

Term 2

Document "chunks"
(e.g., paragraphs)

Darkness to indicate
match count

75

**Objective: Minimize time to decide which documents to look at**

TileBars, Marti A. Hearst

76

---

| THE POSITIVE | | THE NEGATIVE |
|---|---|---|

Simultaneous indication of:
- Relative document lengths
- Frequency of term sets in document
- Distribution of term sets with respect to the document and each other

- Requires training to understand
- Multiple representations and doc formats *(images, page layouts, etc.)*

77

# COMPARING DOCUMENTS

81
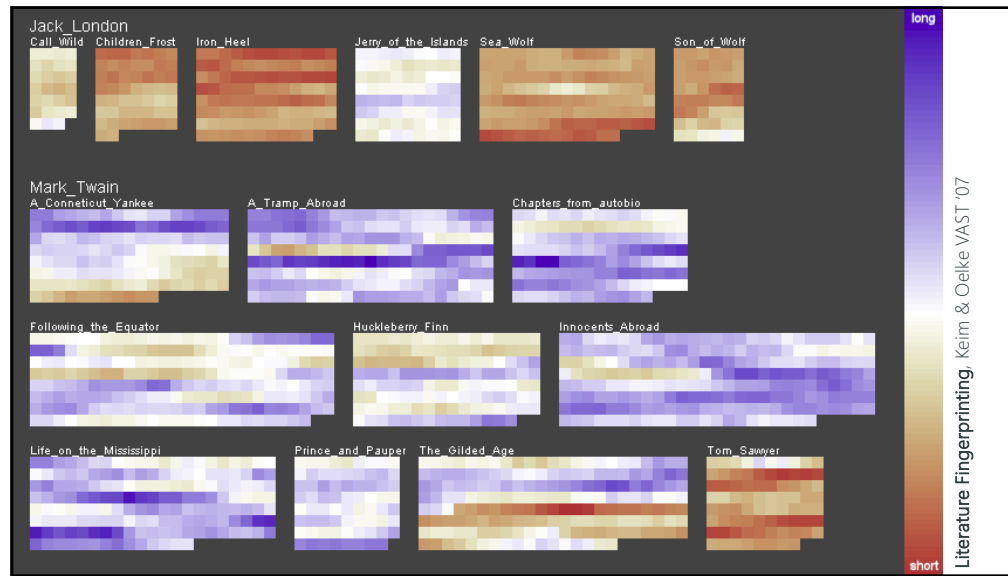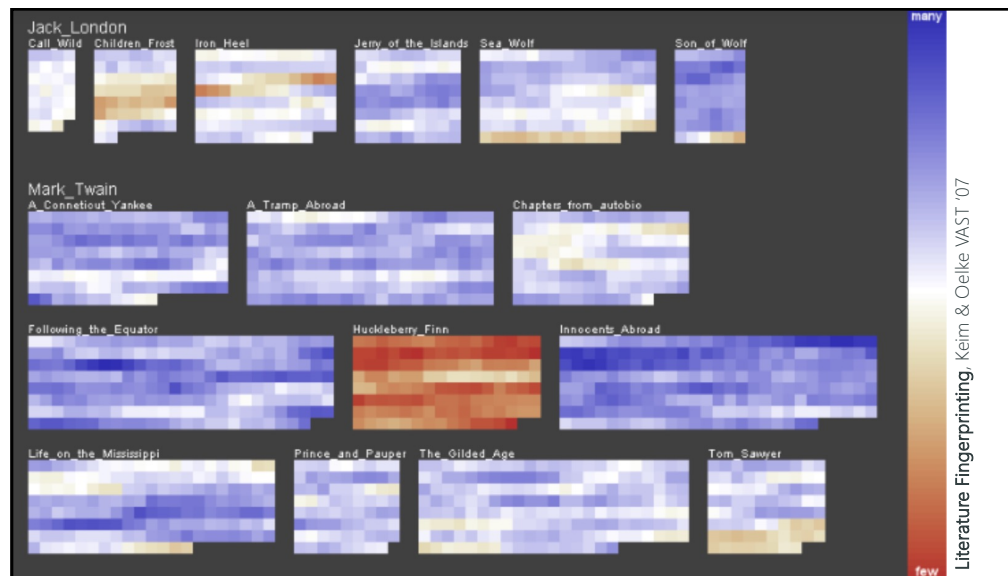
# LITERARY ANALYSIS

Features for comparison
- Word length
- Syllables per word
- Average sentence length
- Percentage by parts of speech (nouns, verbs, etc.)
- Frequencies of specific words
- Hapax Legomena (words that appear once)

82

83



84

Literature Fingerprinting, Keim & Oelke VAST '07

85



The word in context
IRAQ continues to flaunt its hostility toward America and to support terror. The Iraqi regime has plotted to develop anthrax, and nerve gas, and nuclear weapons for over a decade. This is a regime that has already used poison gas to murder thousands of its own citizens — leaving the bodies of mothers huddled over their dead children. This is a regime that agreed to international inspections — then kicked out the inspectors. This is a regime that has something to hide from the civilized world.

— 2002 (Paragraph 20 of 67)

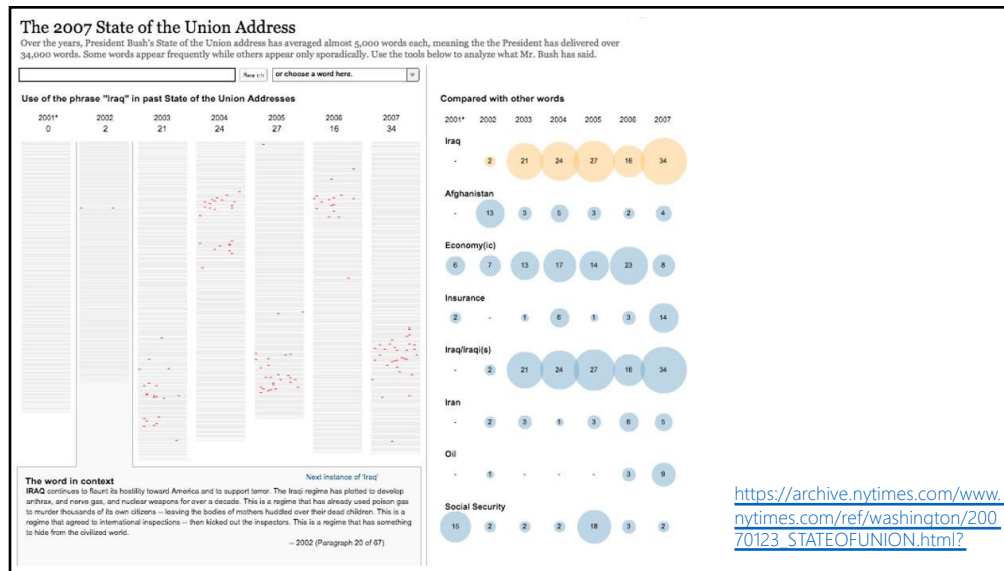https://archive.nytimes.com/www.nytimes.com/ref/washington/20070123_STATEOFUNION.html?

86

THE WORDS THAT WERE USED

READ 2007 SPEECH | FEEDBACK

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.

**LOADING**

https://archive.nytimes.com/www.
nytimes.com/ref/washington/200
70123_STATEOFUNION.html?

87

## 'Social Security'

Mr. Obama's previous speeches have rarely mentioned the term, even though it's a phrase favored by Democrats. Mr. Bush used it repeatedly in 2005 to float an ill-fated plan for privatizing the program.

SOCIAL SECURITY

ROOSEVELT | TRUMAN | EISENHOWER | J.F.K. | L.B.J. | NIXON | FRD. | CAR. | REAGAN | BUSH | CLINTON | BUSH | OBAMA

## 'power'

At one time a favorite of Republicans and Democrats, this word has fallen into post-Vietnam disuse, along with the term "strength."

POWER, POWERED, POWERFUL, POWERFULLY, POWERS

ROOSEVELT | TRUMAN | EISENHOWER | J.F.K. | L.B.J. | NIXON | FRD. | CAR. | REAGAN | BUSH | CLINTON | BUSH | OBAMA
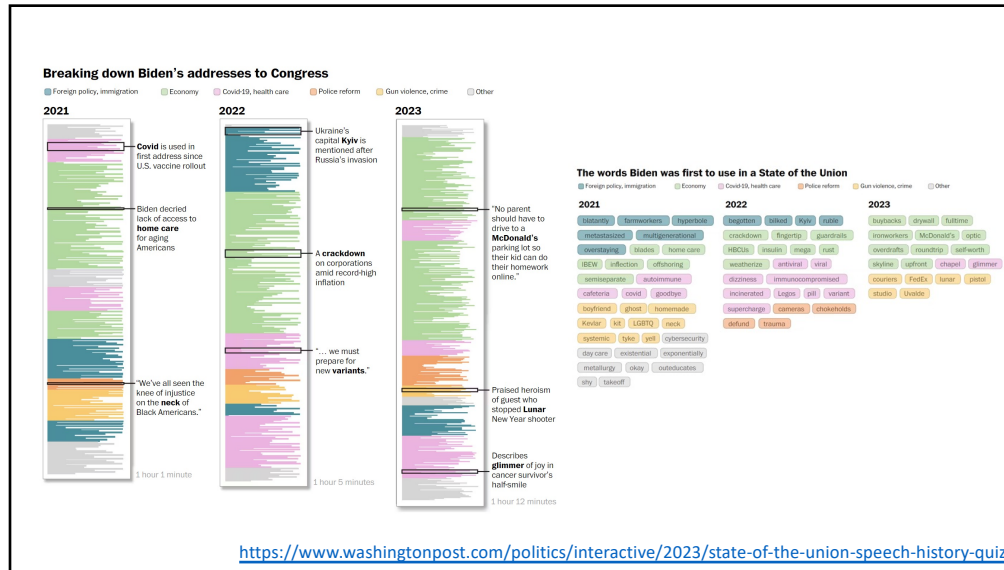
## 'innovate'

A relatively new part of the State of the Union lexicon, "innovation" has been used reliably in the last 20 years. The word has been used to describe paying farmers with surplus grain (Mr. Reagan, 1983), requiring welfare recipients to work (Mr. Reagan, 1988), and artificial retinas that help blind people to see (Mr. Clinton, 2000). Mr. Obama mentioned it 11 times on Tuesday.

INNOVATE, INNOVATION, INNOVATIONS, INNOVATIVE
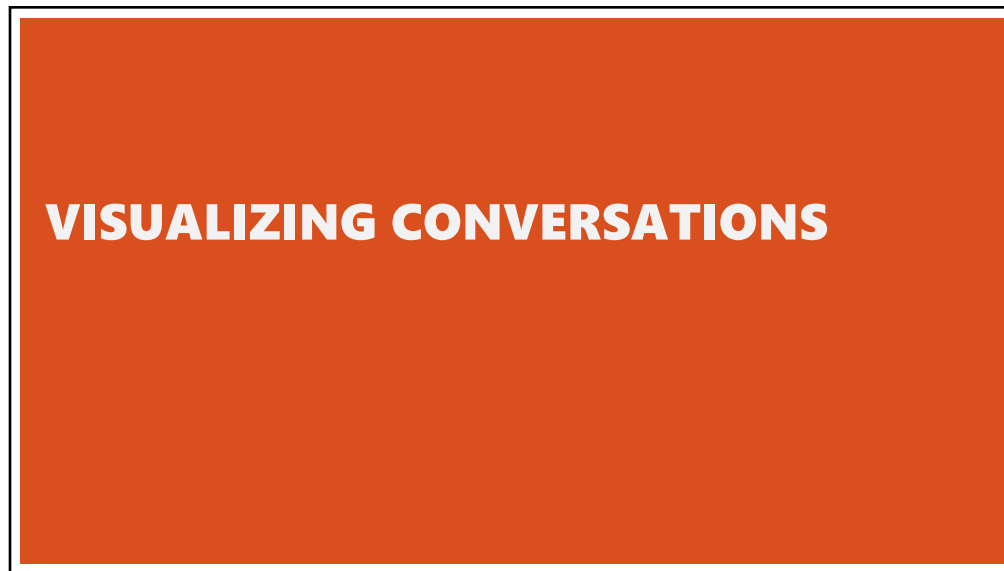
ROOSEVELT | TRUMAN | EISENHOWER | J.F.K. | L.B.J. | NIXON | FRD. | CAR. | REAGAN | BUSH | CLINTON | BUSH | OBAMA

https://archive.nytimes.com/www.nytimes.com/interactive/201
1/01/25/us/politics/state-of-the-union-words-used.html

88

Breaking down Biden's addresses to Congress

https://www.washingtonpost.com/politics/interactive/2023/state-of-the-union-speech-history-quiz/

89



# VISUALIZING CONVERSATIONS
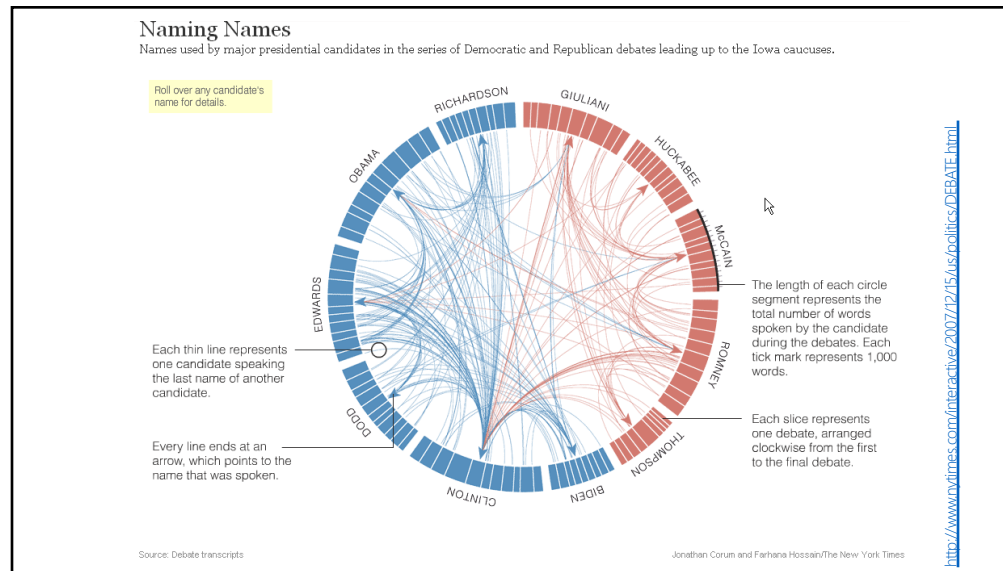
90

# MANY DIMESIONS TO CONSIDER

**Dimensions**
- Who (senders, receivers)
- What (the content of communication)
- When (temporal patterns)

**Cross-products**
- What x When → Topic "Zeitgeist"
- Who x Who → Social network
- Who x Who x What x When → Information flow

91



Naming Names
Names used by major presidential candidates in the series of Democratic and Republican debates leading up to the Iowa caucuses.

92

December 15, 2007

Naming Names

Names used by major presidential candidates in the series of Democratic and Republican debates leading up to the Iowa caucuses.

SIGN IN TO E-MAIL OR SAVE THIS    |    FEEDBACK

Source: Debate transcripts

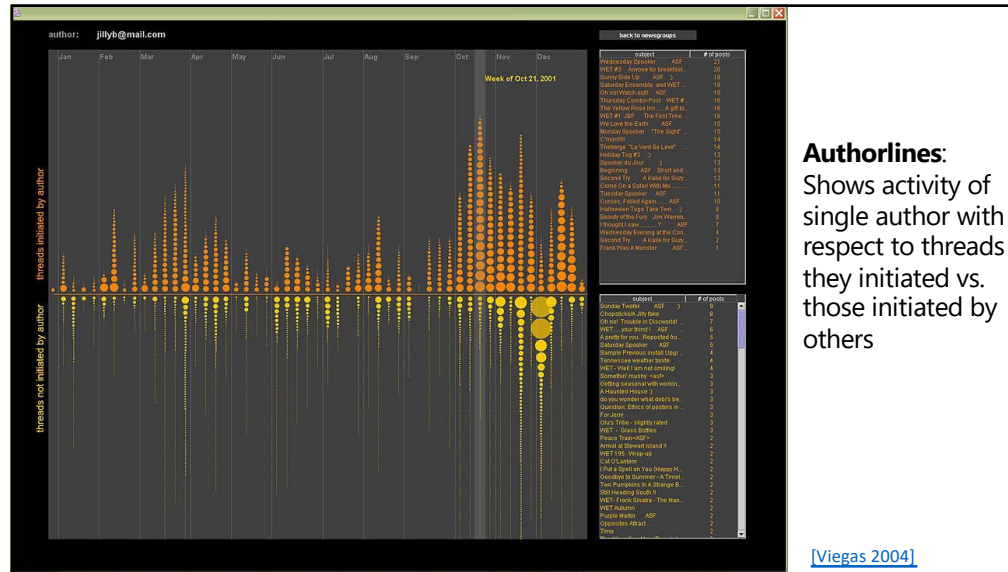Jonathan Corum and Farhana Hossain/The New York Times

http://www.nytimes.com/interactive/2007/12/15/us/politics/DEBATE.html

93

# USENET VISUALIZATION

Show correspondence patterns in text forums

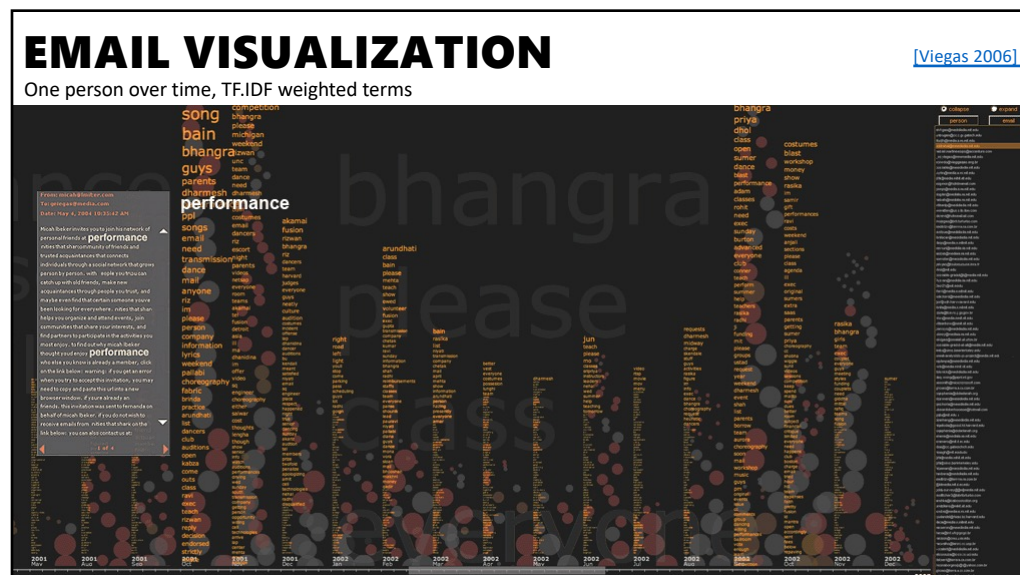Initiate vs. reply; size and duration of discussion



[Viegas 2004]

94

**Authorlines**:
Shows activity of single author with respect to threads they initiated vs. those initiated by others

[Viegas 2004]

95



**Authorlines**:
Shows activity of single author with respect to threads they initiated vs. those initiated by others

[Viegas 2004]

96

## EMAIL VISUALIZATION

[Viegas 2006]

One person over time, TF.IDF weighted terms



98

---

## SUMMARY

Text has many levels of possible visualizations

- Word, document, collection

Identifying descriptive words/keyphrases is critical

- e.g., TF, TF.IDF, regexp, …

- Domain dependent

Can go beyond standard charts and graphs but requires vigilance in design decisions (don't be seduced by possibilities)

99