

Why Visualize Text?



Understanding – get the “gist” of a document

Grouping – cluster for overview or classification

Compare – compare document collections, or inspect evolution of collection over time

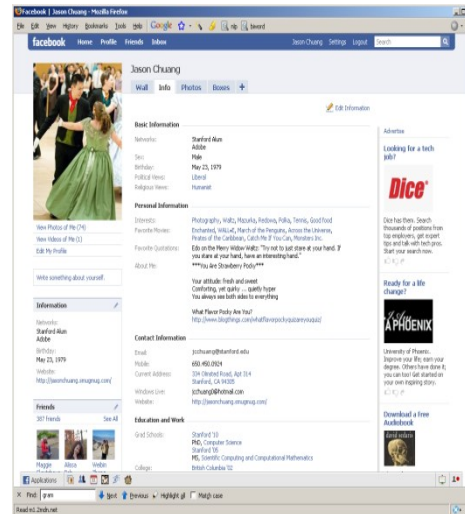
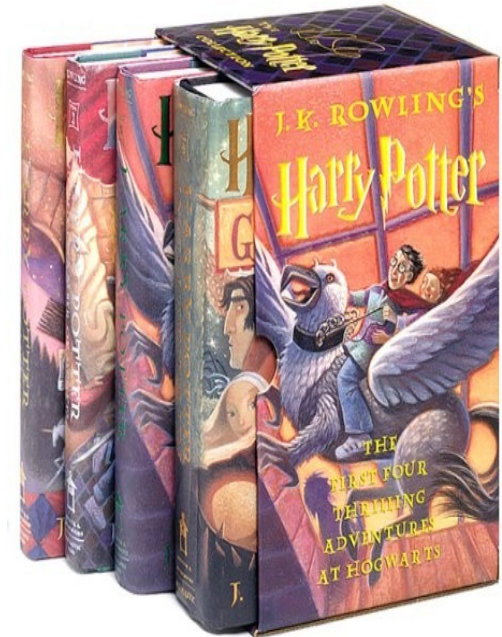
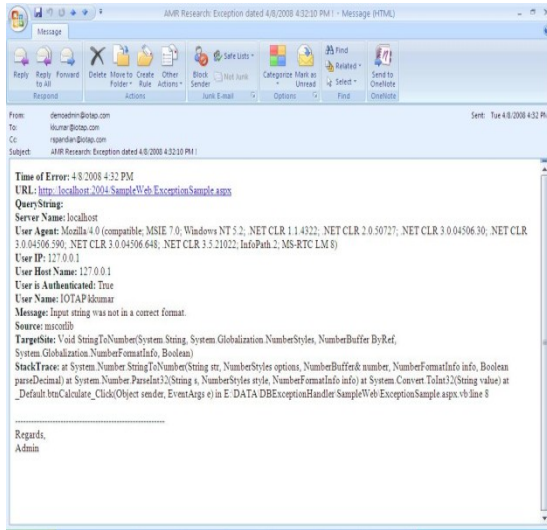
Correlate – compare patterns in text to those in other data, e.g., correlate with social network

What is Text Data



- Documents
 - Articles, books and novels
 - E-mails, web pages, blogs
 - Tags, comments
 - Computer programs, logs
- Collection of documents
 - Messages (e-mail, blogs, tags, comments)
 - Social networks (personal profiles)
 - Academic collaborations (publications)

Where Text Data?



Example: Health Care Reform



September 10, 2009

TEXT

Obama's Health Care Speech to Congress

Following is the prepared text of President Obama's speech to Congress on the need to overhaul health care in the United States, as released by the White House.

Madame Speaker, Vice President Biden, Members of Congress, and the American people:

When I spoke here last winter, this nation was facing the worst economic crisis since the Great Depression. We were losing an average of 700,000 jobs per month. Credit was frozen. And our financial system was on the verge of collapse.

As any American who is still looking for work or a way to pay their bills will tell you, we are by no means out of the woods. A full and vibrant recovery is many months away. And I will not let up until those Americans who seek jobs can find them; until those businesses that seek capital and credit can thrive; until all responsible homeowners can stay in their homes. That is our ultimate goal. But thanks to the bold and decisive action we have taken since January, I can stand here with confidence and say that we have pulled this economy back from the brink.

I want to thank the members of this body for your efforts and your support in these last several months, and especially those who have taken the difficult votes that have put us on a path to recovery. I also want to thank the American people for their patience and resolve during this trying time for our nation.

But we did not come here just to clean up crises. We came to build a future. So tonight, I return to speak to all of you about an issue that is central to that future – and that is the issue of health care.

I am not the first President to take up this cause, but I am determined to be the last. It has now been nearly a

WordTree: Word Sequences



Search Start End

52 hits

i

- will**
 - not**
 - sign**
 - let up until those americans who seek jobs can find them -- (applause) -- until those businesses that seek capital and credit
 - back down on the basic principle that if americans can't find affordable coverage , we will provide you with a choice .
 - a plan that adds one dime to our deficits -- either now or in the future .
 - if it adds one dime to the deficit , now or in the future , period .
 - make that same mistake with health care .
 - waste time with those who have made the calculation that it's better politics to kill this plan than to improve it .
 - and i will not accept the status quo as a solution .
 - accept the status quo as a solution .
 - make sure that no government bureaucrat or insurance company bureaucrat gets between you and the care that you need .
 - protect medicine .
 - continue to seek common ground in the weeks ahead .
 - be there to listen .
- want to**
 - can stand here with confidence and say that we have pulled this economy back from the brink .
 - thank the members of this body for your efforts and your support in these last several months , and especially those who've taken the difficult votes
 - address some of the key controversies that are still out there .
 - clear up -- under our plan , no federal dollars will be used to fund abortions , and federal conscience laws will remain in place
 - speak directly to seniors for a moment , because medicine is another issue that's been subjected to demagoguery and distortion during the course of this debate
- also want** to thank the american people for their patience and resolve during this trying time for our nation .
- return** to speak to all of you about an issue that is central to that future -- and that is the issue of health care
- am**
 - not the first president to take up this cause , but i am determined to be the last .
 - determined to be the last .
- have**
 - to say that there are arguments to be made for both these approaches .
 - no **do** that these reforms would greatly benefit americans from all walks of life , as well as the economy as a whole .
 - interest in putting insurance companies out of business .
- believe**
 - it makes more sense to build on what works and fix what doesn't , rather than try to build an entirely new system from scratch .
 - (laughter) -- i believe a broad consensus exists for the aspects of the plan i just outlined : consumer protections for
 - a broad consensus exists for the aspects of the plan i just outlined : consumer protections for those with insurance , an exchange that allows individuals
- sign** this bill , it will be against the law for insurance companies to drop your coverage when you get sick or water it down when
- already** mentioned .
- just**
 - outlined : consumer protections for those with insurance , an exchange that allows individuals and small businesses to purchase affordable coverage , and a requirement that
 - want to hold them accountable .
 - can't afford it .
- realize**
 - (applause) -- i realize that many americans have grown nervous about reform .
 - that many americans have grown nervous about reform .
- would** remind you that for decades , the driving idea behind reform has been to end insurance company abuses and make coverage available for those without
- say** that rather than making wild claims about a government takeover of health care , we should work together to address any legitimate concerns you may
- faced** a trillion - dollar deficit when i walked in the door of the white house is because too many initiatives over the last decade were
- walked** in the door of the white house is because too many initiatives over the last decade were not paid for -- from the iraq
- don't** believe malpractice reform is a silver bullet , but i've talked to enough doctors to know that defensive medicine may be contributing to unnecessary costs
- know** that
 - the bush administration considered authorizing demonstration projects in individual states to test these ideas .
 - many in this country are deeply skeptical that government is looking out for them .
- think** it's a good idea , and i'm directing my secretary of health and human services to move forward on this initiative today .
- won't** stand by while the special interests use the same old tactics to keep things exactly the way they are .
- received** one of those letters a few days ago .
- understand**
 - how difficult this health care debate has been .
 - that the politically safe move would be to kick the can further down the road -- to defer reform one more year , or one
- still believe**
 - we can**
 - act even when it's hard .
 - replace acrimony with civility , and gridlock with progress .
 - do great things , and that here and now we will meet history's test .
 - i still believe that we can act when it's hard .
 - that we can act when it's hard .

WordTree: Word Sequences



Search Start End

12 hits

i will

- not
 - let up until those americans who seek jobs can find them - - (applause) - - until those businesses that seek capital and credit
 - back down on the basic principle that if americans can't find affordable coverage , we will provide you with a choice .
 - sign
 - a plan that adds one dime to our deficits - - either now or in the future .
 - it if it adds one dime to the deficit , now or in the future , period .
 - make that same mistake with health care .
 - waste time with those who have made the calculation that it's better politics to kill this plan than to improve it .
 - - and i will not accept the status quo as a solution .
 - accept the status quo as a solution .
- make sure that no government bureaucrat or insurance company bureaucrat gets between you and the care that you need .
- protect medicare .
- continue to seek common ground in the weeks ahead .
- be there to listen .

Challenges of Text Visualization



High Dimensionality

Where possible use **text to represent text...**
... which terms are the most descriptive?

Context & Semantics

Provide **relevant context** to aid understanding.

Show (or provide access to) the **source text**.

Modeling Abstraction

Determine your **analysis task**.

Understand abstraction of your **language models**.

Match analysis task with appropriate tools and models.

Topics



- Text as Data
- Visualizing Document Content
- Evolving Documents
- Visualizing Conversation
- Document Collections

Text as Data



Words are (not) nominal?

High dimensional (10,000+)

More than equality tests

Words have meanings and relations

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

Text Processing Pipeline



1. Tokenization

Segment text into terms.

Remove stop words? *a, an, the, of, to, be*

Numbers and symbols? *#gocard, @stanfordfball, Beat Cal!!!!!!!!!!*

Entities? *San Francisco, O'Connor, U.S.A.*

2. Stemming

Group together different forms of a word.

Porter stemmer? *visualization(s), visualize(s), visually* → *visual*

Lemmatization? *goes, went, gone* → *go*

3. Ordered list of terms

Bag of Words Model

Ignore ordering relationships within the text

A document \approx vector of term weights

Each dimension corresponds to a term (10,000+)

Each value represents the relevance

For example, simple term counts

Aggregate into a document-term matrix

Document vector space model

Document-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

Tag Clouds



- Strength
 - Can help with initial query formation.
- Weaknesses
 - Sub-optimal visual encoding (size vs. position)
 - 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

Keyword Weighting

Term Frequency

$tf_{td} = \text{count}(t) \text{ in } d$

Can take log frequency: $\log(1 + tf_{td})$

Can normalize to show proportion: $tf_{td} / \sum_t tf_{td}$

Keyword Weighting

Term Frequency

$$tf_{td} = \text{count}(t) \text{ in } d$$

TF.IDF: Term Freq by Inverse Document Freq

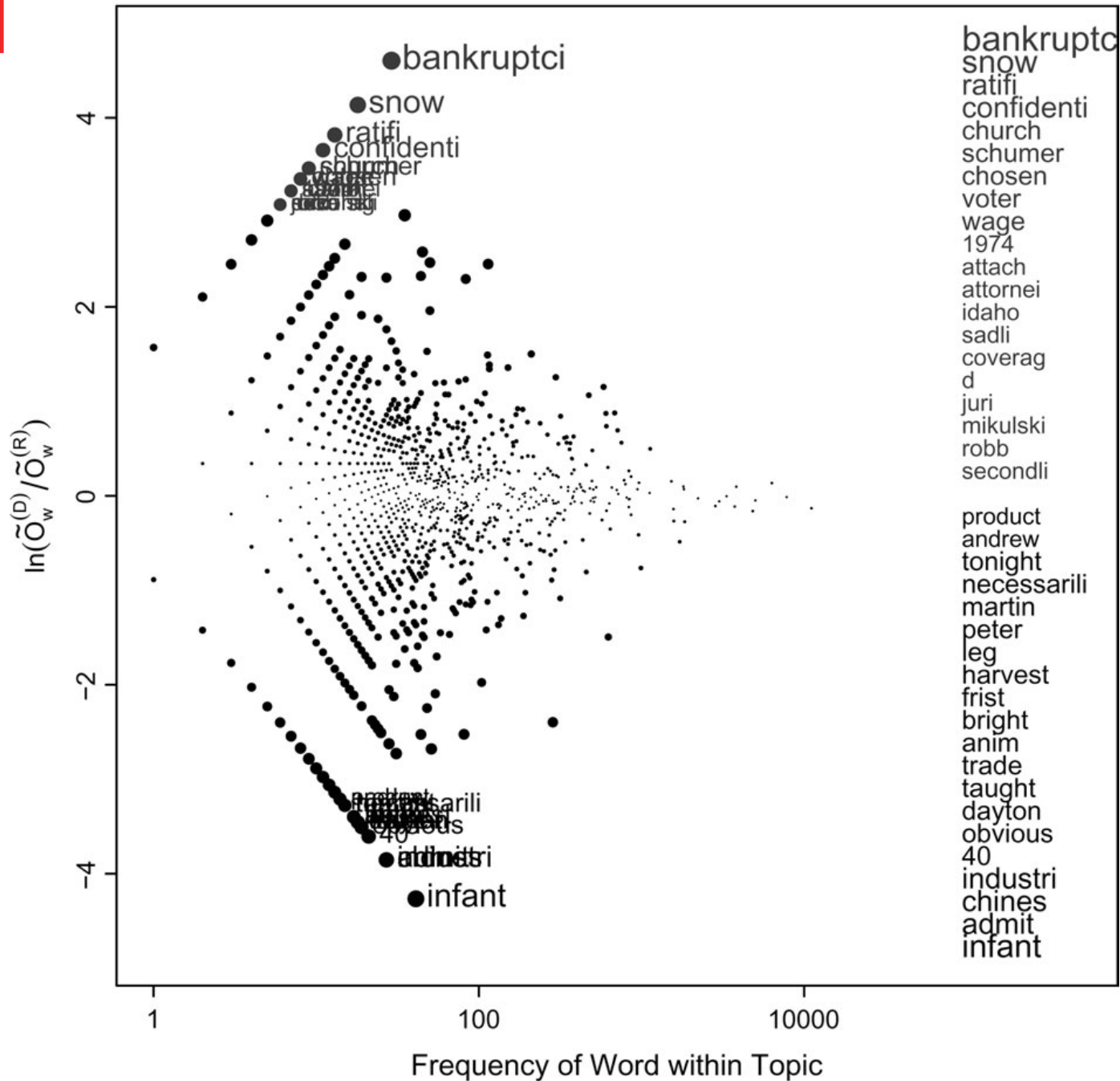
$$tf.idf_{td} = \log(1 + tf_{td}) \times \log(N/df_t)$$

$$df_t = \# \text{ docs containing } t; \quad N = \# \text{ of docs}$$

Partisan Words, 106th Congress, Abortion
(Log-Odds-Ratio, Smoothed Log-Odds-Ratio)



VIT®



Keyword Weighting

Term Frequency

$$tf_{td} = \text{count}(t) \text{ in } d$$

TF.IDF: Term Freq by Inverse Document Freq

$$tf.idf_{td} = \log(1 + tf_{td}) \times \log(N/df_t)$$

$$df_t = \# \text{ docs containing } t; N = \# \text{ of docs}$$

G²: Probability of different word frequency

$$E_1 = |d| \times (tf_{td} + tf_{t(C-d)}) / |C|$$

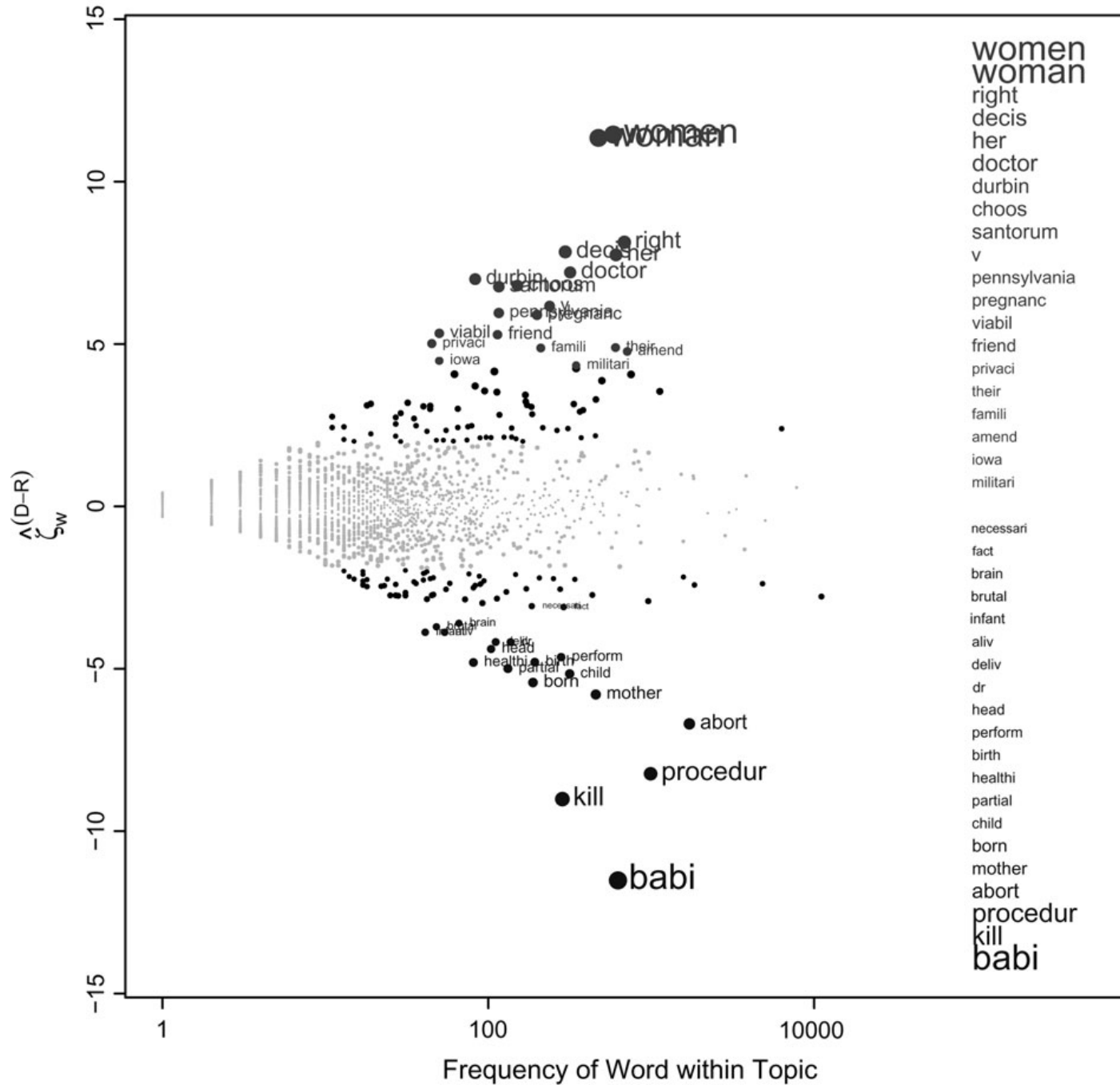
$$E_2 = |C-d| \times (tf_{td} + tf_{t(C-d)}) / |C|$$

$$G^2 = 2 \times (tf_{td} \log(tf_{td}/E_1) + tf_{t(C-d)} \log(tf_{t(C-d)}/E_2))$$

Partisan Words, 106th Congress, Abortion (Weighted Log-Odds-Ratio, Informative Dirichlet Prior)



VIT®



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

A “bag of words” ignores additional information

- Grammar / part-of-speech
- Position within document
- Recognizable entities

How do people describe text?

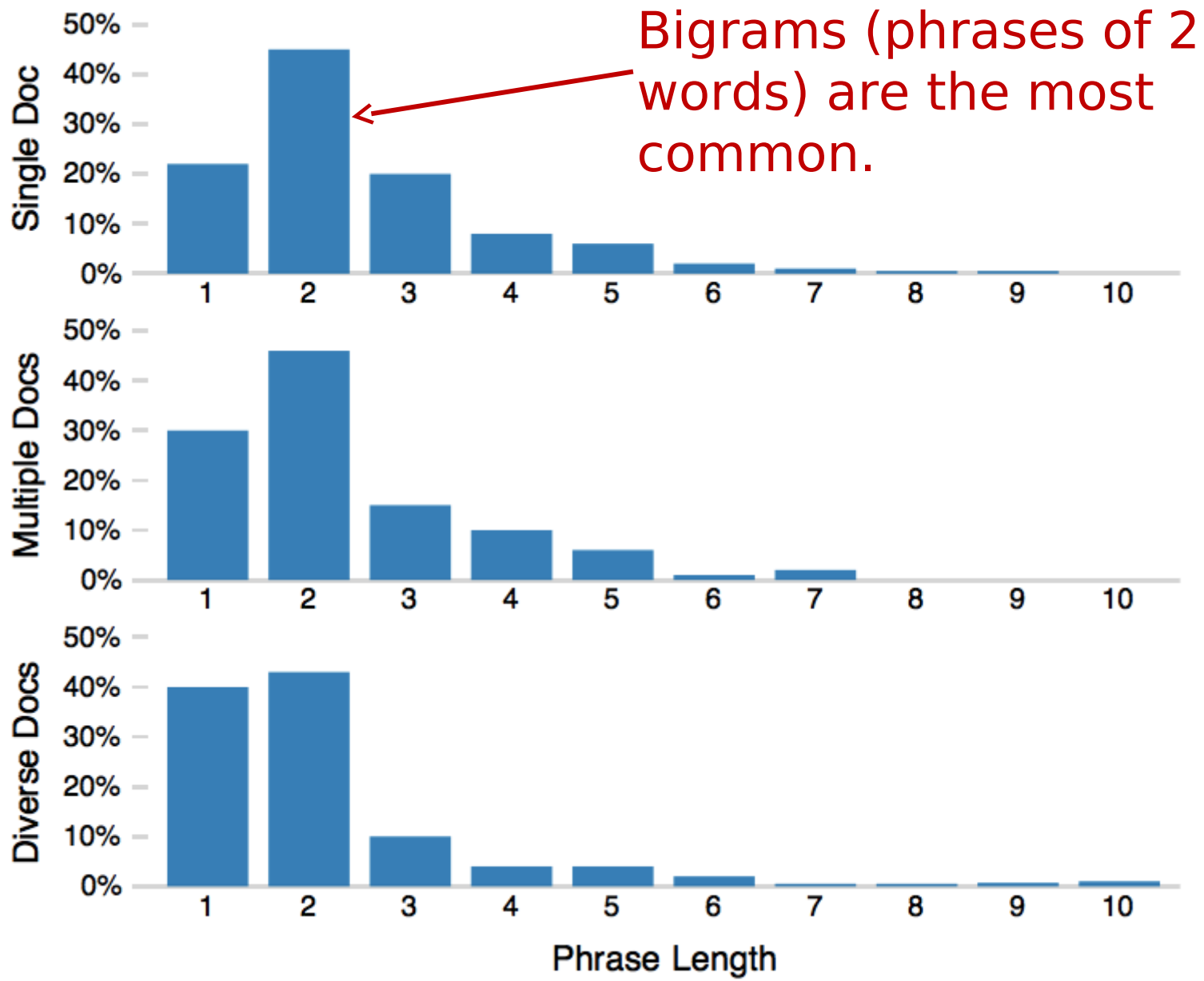


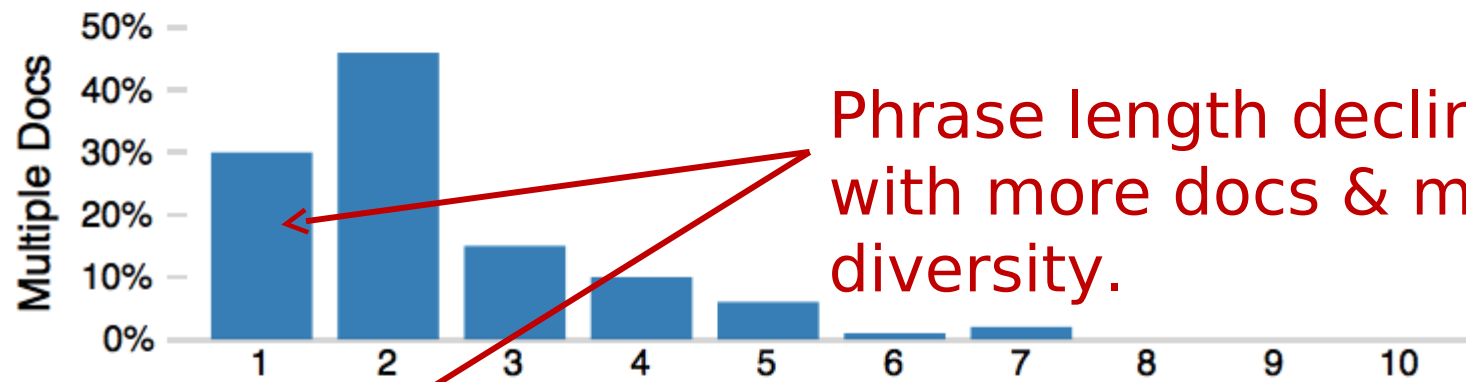
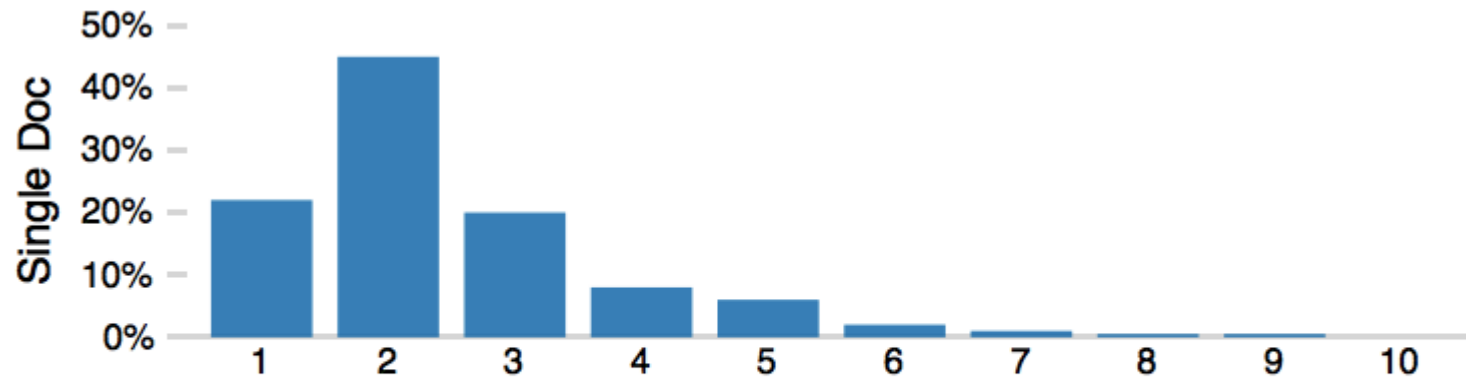
We asked 69 subjects (graduate students) to read and describe dissertation abstracts.

Students were given 3 documents in sequence; they then described the collection as a whole.

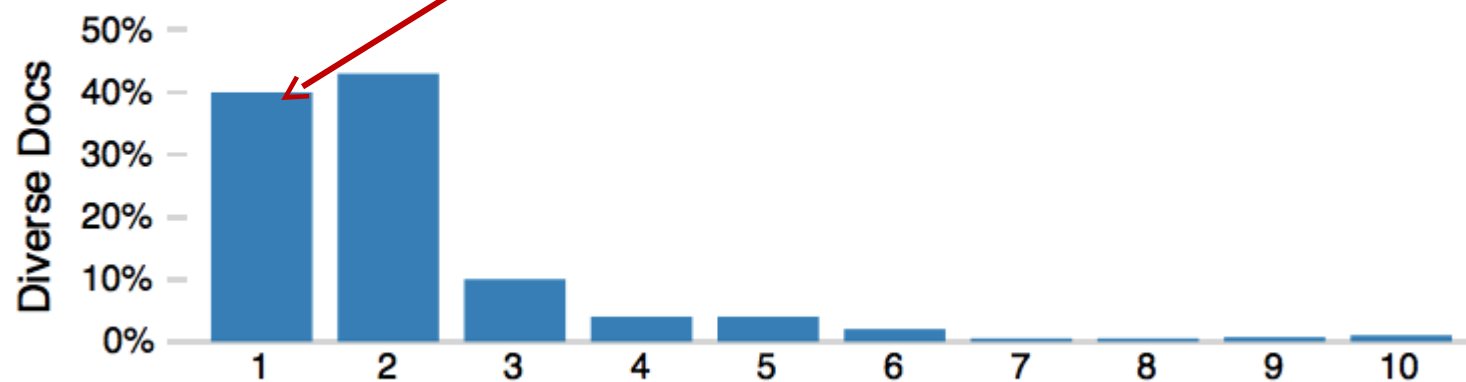
Students were matched to both *familiar* and *unfamiliar* topics; *topical diversity* within a collection was varied systematically.

[Chuang, Heer & Manning, 2010]





Phrase length declines with more docs & more diversity.



Phrase Length

Term Commonness

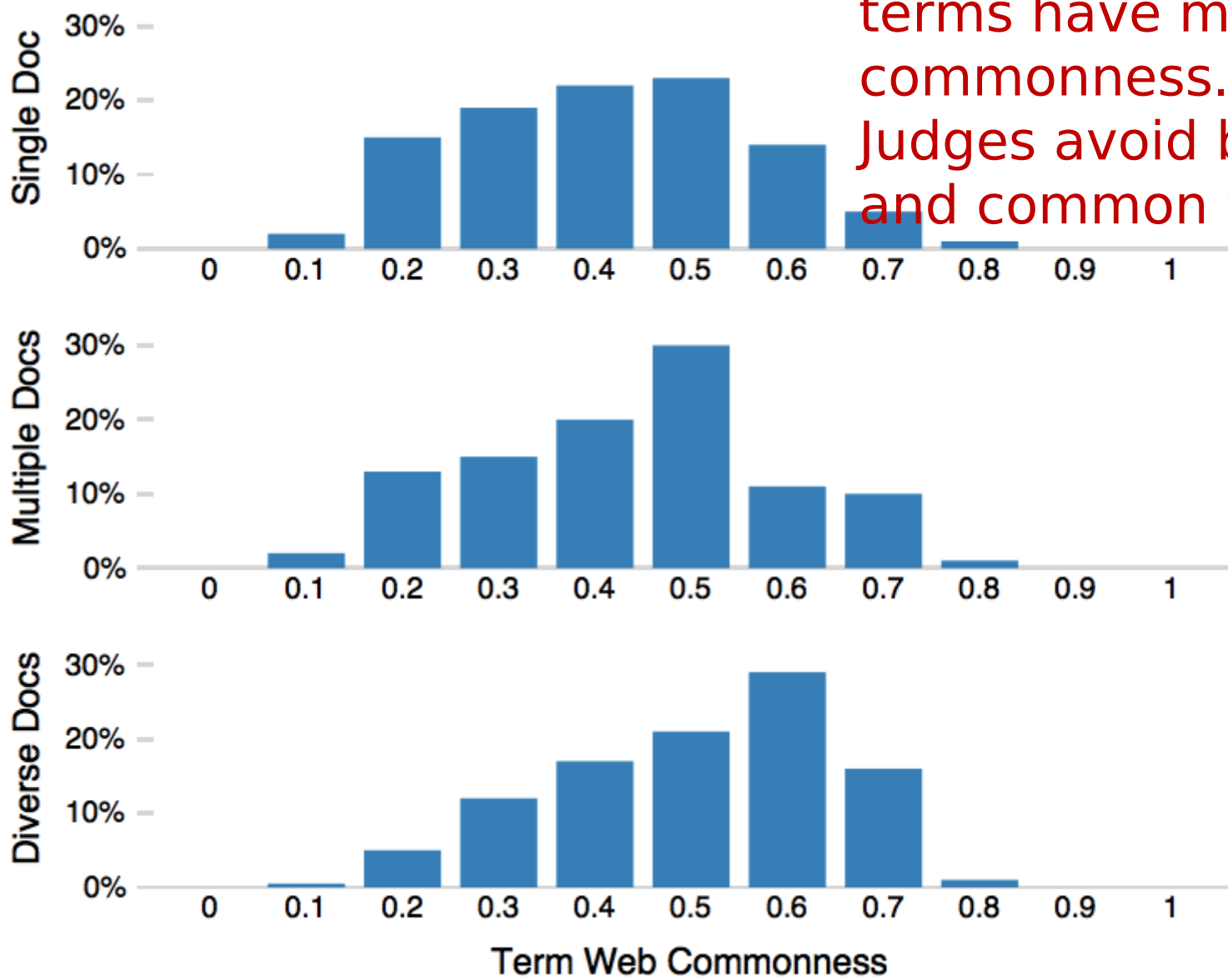
$$\log(\text{tf}_w) / \log(\text{tf}_{\text{the}})$$

The normalized term frequency relative to the most frequent n-gram, e.g., the word “the”.

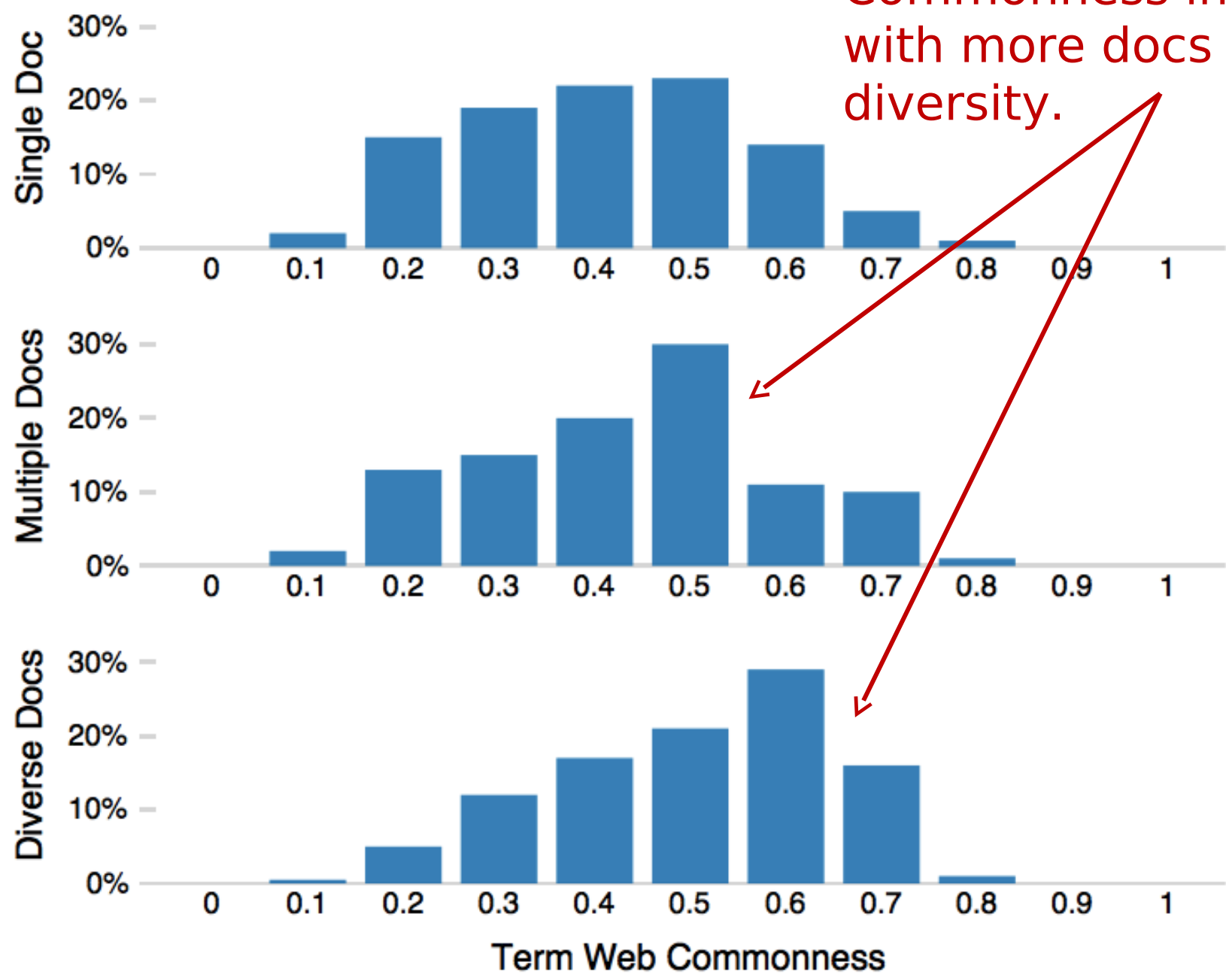
Measured across an entire corpus or across the entire English language (using Google n-grams)



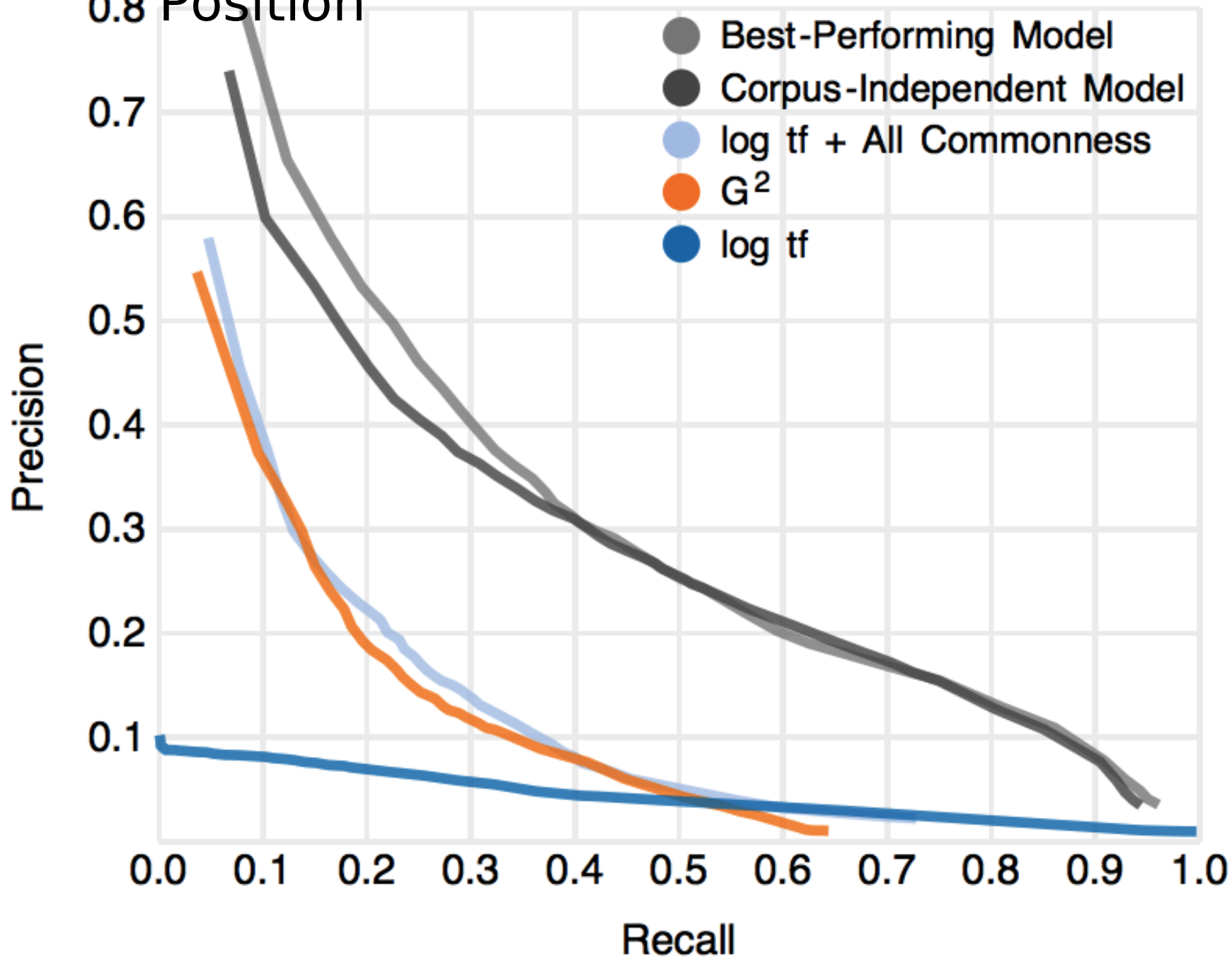
Selected descriptive terms have medium commonness. Judges avoid both rare and common words.



Commonness increases with more docs & more diversity.



Scoring Terms with Freq, Grammar & Position



A fighter jet rain check

Story and video by [Chamila Jayaweera](#)

Have you ever thought about what it takes to make sure that sea-based fighter jets stay dry?

When it comes to the F/A-18 Super Hornet, Boeing engineers in St. Louis use a special process called the Water Check Test to rule out areas where moisture could seep into the aircraft and its electronics suite.

Program experts douse the jet with simulated rain at a 15-inch-per-hour rate for about 20 minutes inside an enormous hangar in St. Louis.

"Our ultimate customers are U.S. Navy fighter pilots, and we want to ensure their safety in flight and on the ground, and water-tight integrity of the aircraft also helps increase their effectiveness," said Boeing's Rich Baxter, F/A-18 Super Hornet final assembly manager.

To find out more about how the process works and watch the action unfold, click above to see the video story.



CHAMILA JAYAWEERA/BOEING

The Water Check team rolls in a large metal frame, which they affectionately call their "spray tree," over a Super Hornet inside a St. Louis hangar.



G^2

Regression Model

fighter

F/A

Hornet

Super

Boeing

-18

rain

St.

jet

Louis

15-inch-per-hour

douse

hangar

water-tight

Check

Baxter

sea-based

aircraft

Rich

seep

click

Navy

sure

Water

moisture

watch

enormous

stay

want

Super Hornet

F/A -18

fighter jet

Boeing engineers

special process

rain check

electronics suite

Program experts

simulated rain

ultimate customers

enormous hangar

water-tight integrity

Rich Baxter

15-inch-per-hour rate

video story

aircraft

U.S. Navy fighter pilots

Super Hornet final assembly manager

U.S.
Navy fighter
fighter pilot
sea-based fighter

Yelp: Review Spotlight [Yatani 2011]



VIT®



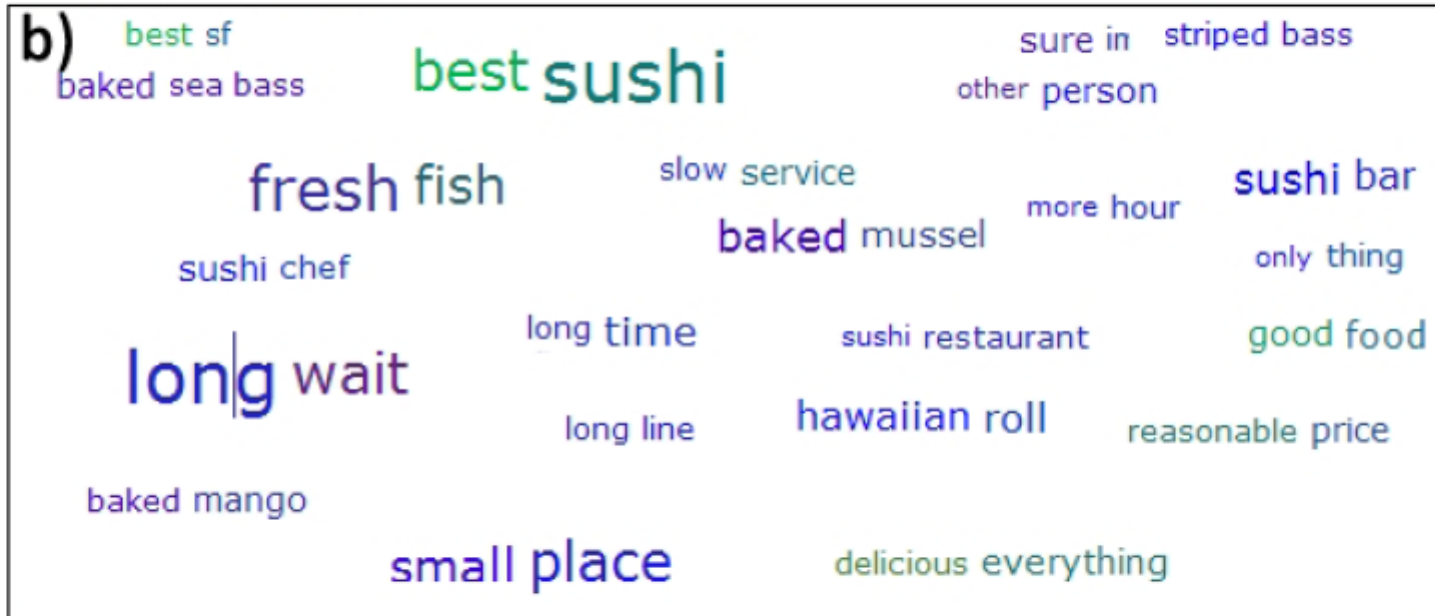
“long wait” or “no wait”?

what type of sushi roll?

Yelp: Review Spotlight [Yatani 2011]



VIT®



Mentioned 63 times

possess sage of the halos wisdom , and know in advance **sushi zone** only accepts cash and the waits will be **long** and arduous .

yes , its a **long** wait , learn the master of zen if you want to eat here .



Tips: 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

Show keyphrases within source text



Thank You