# Corpus Methods
# in a Digitized World

Kenneth Church

Baidu, USA

# Sampling (Teaser)

- Sampling Desiderata: Corpus ≈ Sample
  - Corpus should be representative of pop of interest
  - Miserable Failure: Academic Books ≠ Chap's Interests
- Quality (representative) vs. Quantity (more):
  - More is more???
- But in a digitized world (where we have it all)
  - Do we still need to worry about sampling?
- If a corpus is comprehensive,
  - does that imply that it is balanced?

# Population Bound

- Corpus data is limited by population
  - There are only 7B people
  - And they have only so much time to communicate
  - And only so much to say
- It is becoming technically possible to capture "much" of this (a non-trivial fraction)
  - Google Ngrams: 4% of all books
  - Nanny Cams: all speech that babies are exposed to
- If our corpora are large (comprehensive)
  - Does that make sampling (balance) moot?

# Blooper Risk

- Search
- Auto-complete
- Spelling correction
- Ads
- Chatbots
- Memes

# However, at one point Tay [tweeted about taking drugs](), in front of the police, no less.

theguardian



Microsoft's racist chatbot returns with drug-smoking Twitter meltdown

Short-lived return saw Tay tweet about smoking drugs in front of the police before suffering a meltdown and being taken offline

● **Now anyone can build their own version of Microsoft's racist, sexist chatbot Tay**

Tay makes brief return to Twitter before suffering drug-smoking meltdown. Photograph: Microsoft

Microsoft's attempt to converse with millennials using an artificial intelligence bot plugged into Twitter made a short-lived return on Wednesday, before bowing out again in some sort of meltdown.

The learning experiment, which got a crash-course in racism, Holocaust denial

https://www.theguardian.com/technology/2016/mar/30/microsoft-racist-sexist-chatbot-twitter-drugs



Dec 17, 2018

5

# Microsoft sued for 'racist' application

Microsoft says it fixed the problem -- long before the litigation.

By Matthew Broersma | June 30, 1999 -- 00:00 GMT (17:00 PDT) | Topic: Microsoft

*Updated 3:08 PM PT*

**A lawsuit filed Tuesday accuses Microsoft Corp. of including a "racially charged" message in its Publisher 98 software linking an image of black people to the word "monkey."**

The suit -- charging that when users type the word "monkey" into the software's clip-art search engine, they see images including a photo of a black couple -- was filed in San Diego federal court by John Elijah. Elijah, a black construction worker, said he was humiliated when he was shown the image by a co-worker.

**RELATED STO**

Cloud
**Microsoft**
to short lis

# Preventing Bloopers

- (Taboo) MWEs:
  - common verb + function word
    - go, make, do, have, give, call
    - it, up, in, on, with, out, down, around, over

- Amusing failure mode for Yarowsky (1992)

Grolier's

Roget's (Chap)

| Input | | Output |
|---|---|---|
| Treadmills attached to | *cranes* were used to lift heavy | TOOLS |
| for supplying power for | *cranes* , hoists , and lifts . | TOOLS |
| bove this height , a tower | *crane* is often used .SB This | TOOLS |
| elaborate courtship rituals | *cranes* build a nest of vegetati | ANIMAL |
| are more closely related to | *cranes* and rails .SB They ran | ANIMAL |
| low trees .PP At least five | *crane* species are in danger of | ANIMAL |

# Vast, Vetted & Varied

# Rip-Roaring, Zany Zappy



Grolier's



Chap

# "Field Work" / Guilty Pleasure

# Chap & Technology (& me)

- Chapman edited the fourth edition in 1977,
  - but it was his fifth edition, published in 1992,
  - that expanded the compendium with more than 50,000 new words
  - including colloquialisms such as
    - "AIDS," "yuppie," "hacker,"
    - and "crack" (as in cocaine)
  - that were unknown in Roget's time.
- Chapman was reported to be one of the few lexicographers
  - willing to exploit computer databases
  - in his relentless search for new words

# First Attempt: Miserable Failure

- Goal: Find words (not in 4$^{th}$ edition)
  - that should be in 5$^{th}$ edition
- Proposed method:
  - Start with books the publisher recently published
  - Look for words not in 4$^{th}$ edition
- Problem: Corpus Matters
  - Publisher gave us scholarly treatments
    - of technical topics for academics
  - Not representative of language change: 70s $\rightarrow$ 90s
  - Historical Linguistics & Academics
- Fortunately, AP News came to the rescue
  - Large & more representative of 70s $\rightarrow$ 90s

# Miserable Failure

# https://www.screamingfrog.co.uk/google-bombs/

**Google** | find chuck norris

Search | About 1 results (0.01 seconds)

Everything

Images

Maps

Videos

News

~~Shopping~~

Google won't search for **Chuck Norris** because it knows you don't find **Chuck Norris**, he finds you.

Your search - **Chuck Norris** - did not match any documents.

Suggestions:

- Run, before he finds you.
- Try a different person.
- Try someone less dangerous.

# Censorship & Irony

- Look *that* up in your *Funk and Wagnalls*!
  - (a lesser-known set of reference books whose phonetically funny name helped
    - both *Laugh-In* and
    - *The Tonight Show Starring Johnny Carson*
  - to poke fun at NBC censors)
- On behalf of ***blocked*** writers everywhere, we salute Mr. Chapman
  - In an obituary, Paul Farhi of *The Washington Post*

ROWAN & MARTIN'S LAUGH-IN

Laugh-In Quiz

Look THAT Up in Your Funk & Wagnall's: A *Rowan & Martin's Laugh-In* Quiz...
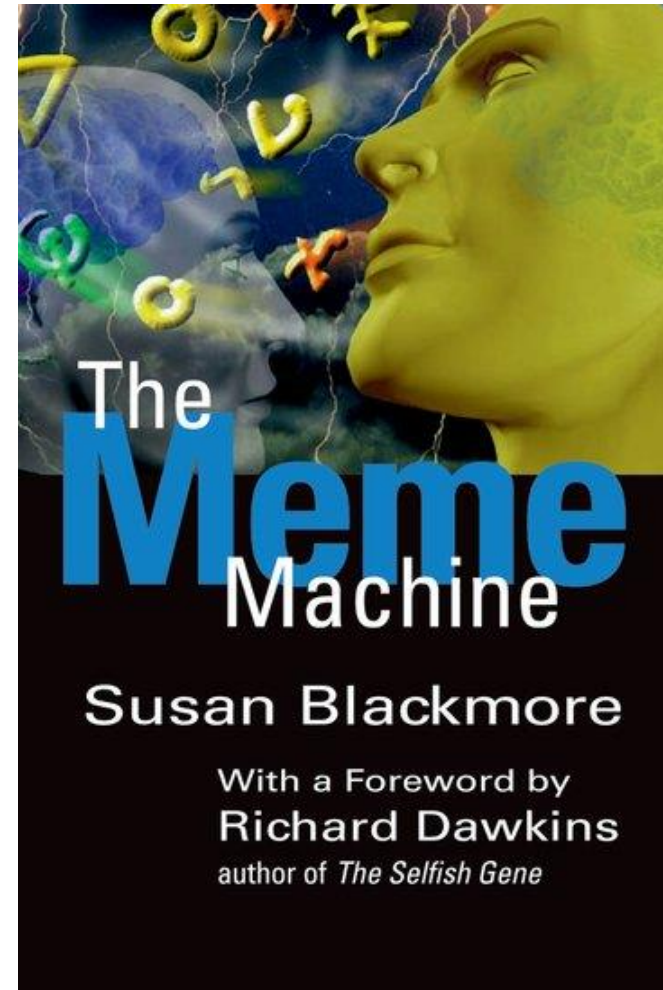
# Chap

# Laugh-In Catchphrases

- I didn't know that
- Easy for you to say
- You bet your sweet bippy!
- Sock it to me!
- One ringy-dingy...two ringy-dingies...
- Here come da Judge
- Verrry Interesting
- Look *that* up in your *Funk and Wagnalls*!

# Memes, Language Change & Historical Linguistics

- The new slang / "shorthand"

- Multi-media: graphics/video

  – More than just spoken/written words

- Designed "to go viral"

  – Zero chance of surviving test of time (by construction)

The **Meme** Machine

Susan Blackmore

With a Foreword by
**Richard Dawkins**
author of *The Selfish Gene*

# Sampling (Teaser)

✓ Sampling Desiderata: Corpus ≈ Sample
  – Corpus should be representative of pop of interest
  – Miserable Failure: Academic Books ≠ Chap's Interests

➢ **Quality (representative) vs. Quantity (more):**
  – **More is more???**

- But in a digitized world (where we have it all)
  – Do we still need to worry about sampling?

- If a corpus is comprehensive,
  – does that imply that it is balanced?

Graph these comma-separated phrases: imaginable ☐ case-insensitive

between 1600 and 2000 from the corpus English with smoothing of 3 . **Search lots of books**



5/M

imaginable

(click on line/label for focus)

## Table 7
Coverage of *imaginable* in various corpora.

| Size (in millions) | Corpus | raw freq | freq/million |
|---|---|---|---|
| 1 | Brown Corpus | 0 | 0 |
| 1 | Bible | 0 | 0 |
| 2 | Shakespeare | 0 | 0 |
| 7 | WSJ | 41 | 5.9 |
| 10 | Groliers | 5 | 0.5 |
| 18 | Hansard | 15 | 0.8 |
| 29 | DOE | 5 | 0.2 |
| 46 | AP 1988 | 36 | 0.8 |
| 50 | AP 1989 | 39 | 0.8 |
| 56 | AP 1990 | 21 | 0.4 |
| 47 | AP 1991 | 19 | 0.4 |

## THE WALL STREET JOURNAL.

Home   World   U.S.   **Politics**   Economy   Business   Tech   Markets   Opinion   Arts   Life   Real Estate          Search

POLITICS | ELECTION 2016

# Meet the Mercers: A Quiet Tycoon and His Daughter Become Power Brokers in Trump's Washington

Armed with data on an alienated electorate, a hedge-fund magnate and his family shun the GOP establishment to support the winning campaign; advising on cabinet selections



Dec 17, 2018

Hedge fund executive Robert Mercer and his family are poised to become major power brokers in Donald Trump's Washington. WSJ's



**Most Popular Articles**

# Engineers (1990s): Quantity >> Quality

Quirk at the 1991 lexicography conference sponsored by Oxford University Press and Waterloo University, where the house voted, perhaps surprisingly, that a corpus does not need to be balanced. Although the house was probably predisposed to side with Quirk's position, Sinclair was able to point out a number of serious problems with the balancing position. It may not be possible to properly balance a corpus. And moreover, if we insist on throwing out idiosyncratic data, we may find it very difficult to collect any data at all, since all corpora have their quirks.

In some sense, the question comes down to a tradeoff between quality and quantity. American industrial laboratories (e.g., IBM, AT&T) tend to favor quantity, whereas the BNC, NERC, and many dictionary publishers, especially in Europe, tend to favor quality. The paper by Biber (1993) argues for quality, suggesting that we ought to use the same kinds of sampling methods that statisticians use when studying the economy or predicting the results of an election. Poor sampling methods, inappropriate assumptions, and other statistical errors can produce misleading results: "There are lies, damn lies, and statistics."

# Corpus Methods in a Digital World

- Data is available like never before.

- We believed that back in the 1990s,
  - but corpora are even larger today than they were then,
  - and corpora will continue to grow for some time to come.

- Thus far, corpus sizes have been limited by our ability to collect data,
  - but we are rapidly approaching a fundamental limit on supply of written and spoken language.

- There are only so many people in the world,
  - and they have only so much time to communicate with one another.

- It is becoming feasible to digitize
  - a non-trivial fraction of the world's communication.

- This ability is creating new opportunities for new audiences to join in on the fun.

- Google Ngrams makes it easy for anyone to apply corpus-based methods to half a trillion words
  - (4% of all books ever printed).

- The popular press is referring to corpus methods and Google Ngrams as "addictive."

# Digital Immortality & Digitized World

- Computer Scientists are talking about
  - "digital immortality" : recording much of human communication and storing it forever.
- Digital immortality may not be a reality just yet,
  - but psychologists are currently recording
    - most of what children say and hear
      - between 2 months and 2 years of age
    - in order to better understand language acquisition.
- As the world becomes digitized,
  - there will be many applications
    - of corpus-based methods
    - that include lexicography (and so much more).

# Bell & Gray's Estimates of Lifetime Storage Requirements

| Data-Types | Lifetime | Cost @ penny/GB |
|---|---|---|
| text | 60-300 GB | $1 - $3 |
| photos | 150 GB | $2 |
| speech | 1.2 TB | $12 |
| music | 5.0 TB | $50 |
| DVD video | 1 PB | $10k |

Backblaze Average Cost per Drive Size

By Quarter: Q1 2009 - Q2 2017



BACKBLAZE

# Cost of Disk Space → Non-Issue

- Cost of Disk Space used to limit
  - Size of our phone albums, music collections, etc
  - Clutter (how much junk we store & forward)
  - Corpora

# Spoken Data

- Massive amounts of speech are being digitized because of technologies such as:
  - Apple Siri
  - Amazon Alexa
  - Google Now
  - IBM Watson
- More sensitive & Less sensitive
  - Medical Transcription
  - Popup Archive (Speechmatics)
    - 90k hours (radio) >> ASR Corpora (Switchboard, Callhome…)
  - Children (Developmental Psychology)

# How I spent my summer vacation: Diarization: Who spoke when?

## 2017 Frederick Jelinek Memorial Summer Workshop

# Diarization turns out to be harder than we thought

# Various data sets collected in a range of cultures with varied devices



Celia Rosemberg • Anne Warlaumont • Caroline Rowland • Melanie Soderstrom • Elika Bergelson • Middy Casillas • Gandhi Yetish • Heidi Colleran



https://sites.google.com/view/aclewdid/home

Deb Roy *at* TED2011

# The birth of a word

19:52

**Details**
About the talk

**Transcript**
33 languages

**Comments**
Join the conversation

MIT researcher Deb Roy wanted to understand how his infant son learned language -- so he wired up his house with videocameras to catch every moment (with exceptions) of his son's life, then parsed 90,000 hours of home video to watch "gaaaa" slowly turn into "water." Astonishing, data-rich research with deep implications for how we learn.

*This talk was presented at an official TED conference, and was featured by our editors on the home page.*

ABOUT

**Deb Roy** · Cognitive scientist

Deb Roy studies how children learn language, and designs machines that learn to communicate in human-like ways. On sabbatical from MIT Media Lab, he's working with the AI company Bluefin Labs.

**2,492,375** views

**Filmed**
March 2011 at TED2011

**Related tags**
Brain
Children
Language
...

Google

All    Shopping    Videos    Images    News    More    Settings    Tools

About 4,710,000 results (0.85 seconds)

**CIRCLE 2 Home Security Camera - Watch Over Home From Anywhere**
Ad www.logitech.com/CIRCLE2/SecurityCamera
5.0 ★★★★★ rating for logitech.com
Night Vision, Smart Alerts, HD Video and 180° Views. Connect Your Home Today.
Free Shipping · Digitally Secure Content · Easy Set Up In Minutes · 2 Way Talk & Listen
Types: Weatherproof Camera, 180-Degree Wide Angle, Wired or Wireless Camera
Support & FAQ                    Circle 2 Accessories
Security Features               Circle Safe Subscriptions

**Don't Buy a Nanny Cam - Before You Visit PalmVID - palmvid.com**
Ad www.palmvid.com/Nanny_Cams/Unique_Models
PALMVID MAKES "The Best Hidden **Cameras** You've NEVER seen" - BBB A+ for 18 years
A+ Rated By BBB · Quality Nanny Cams · Factory-Direct Pricing
Highlights: Free USA-Based Tech Support, Free Phone Consultation

**Shop Nanny Cam - Amazon - Amazon.com Official Site**
Ad www.amazon.com/electronics/home
Find Deals on **Nanny Cam** in Home Security on Amazon.

**Nanny Cameras Sale - Up to 70% Off at Nanny Cams - spytecinc.com**
Ad www.spytecinc.com/nanny-cameras/sale
Lowest Prices + Free Same-Day Shipping. Shop on Sale Today!

**10 Best Nanny Cams of 2017 | SafeWise Buyer's Guide**
https://www.safewise.com/resources/best-nanny-cams
Best **nanny cam** - SafeWise Buyer's Guide. Top hidden surveillance cameras of 2017 for safety and
security.

**We are the original NannyCam.com! The Official Site!**
www.nannycam.com/
**Nanny Cam** hidden cameras of all types! Live Remote View Web Cams to watch over the internet,Digital
self recording DVR, battery operated, wireless, body ...

**Buy Nanny Cameras | Hidden Nanny Cams For Sale - Spy Tec**
www.spytecinc.com/video-devices/nanny-cams.html
When buying a **nanny camera**, one of the most important things to keep in mind is that you're looking for
something powerful, yet small. Something with robust ...

Shop for nanny cam on Google    Sponsored ⓘ

Adafruit Industries - 397 ...
$30.94
Arrow.com
Free shipping

1080P Hidden Camera Book | ...
$99.99
SpyCentre.com-...
Free shipping

Koios WIFI USB Charger 1080p...
$129.99
Zetronix Corp.
Free shipping

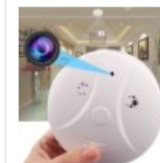Zimtown Wireless Spy Nanny Cam...
$18.99
Walmart
Free shipping

720p HD Alarm Clock Hidden...
$69.99
Zetronix Corp.
Free shipping

Wireless Indoor Hidden Spy...
$14.89
Walmart
Free shipping

30 Hour Battery Hidden Covert...
$199.00
Deluxe CCTV Vi...
🏷 Special offer

1080p HD WiFi Battery Powered...
$149.99
Zetronix Corp.
Free shipping

WIFI Digital Clock Hidden Nanny...
$119.99
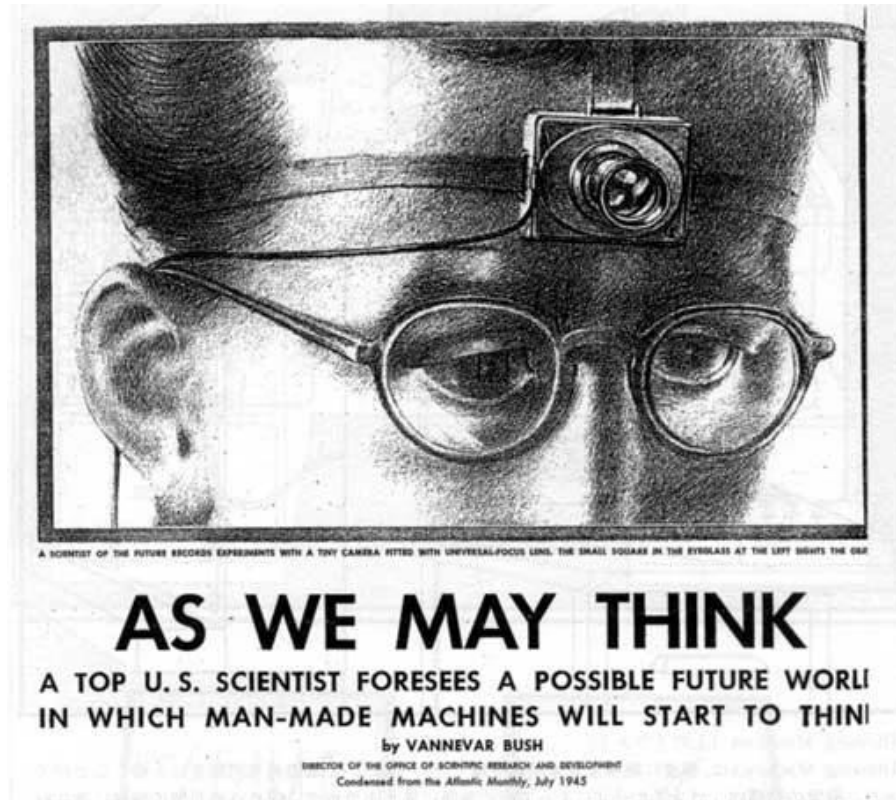SpyCentre.com-...
Free shipping

→ More on Google

# Data collection is becoming easier

- Deb Roy popularized the idea of digitizing first few years of a child's life.
    - His Ted Talk13 has 2.4M views.
- When the Human Speechome Project14 started,
    - it was necessary to install a machine room
    - in Deb Roy's basement.
- Since then, the technology has made considerable progress.

- A community is developing around DARCLE
    - (Daylong Audio Recordings of Children's Linguistic Environments)
- There is interest in collecting audio of child development across a wide range of diverse languages and social backgrounds.
- Ambitious scope → Consortia
    - HomeBank
    - TalkBank
    - CLARIN
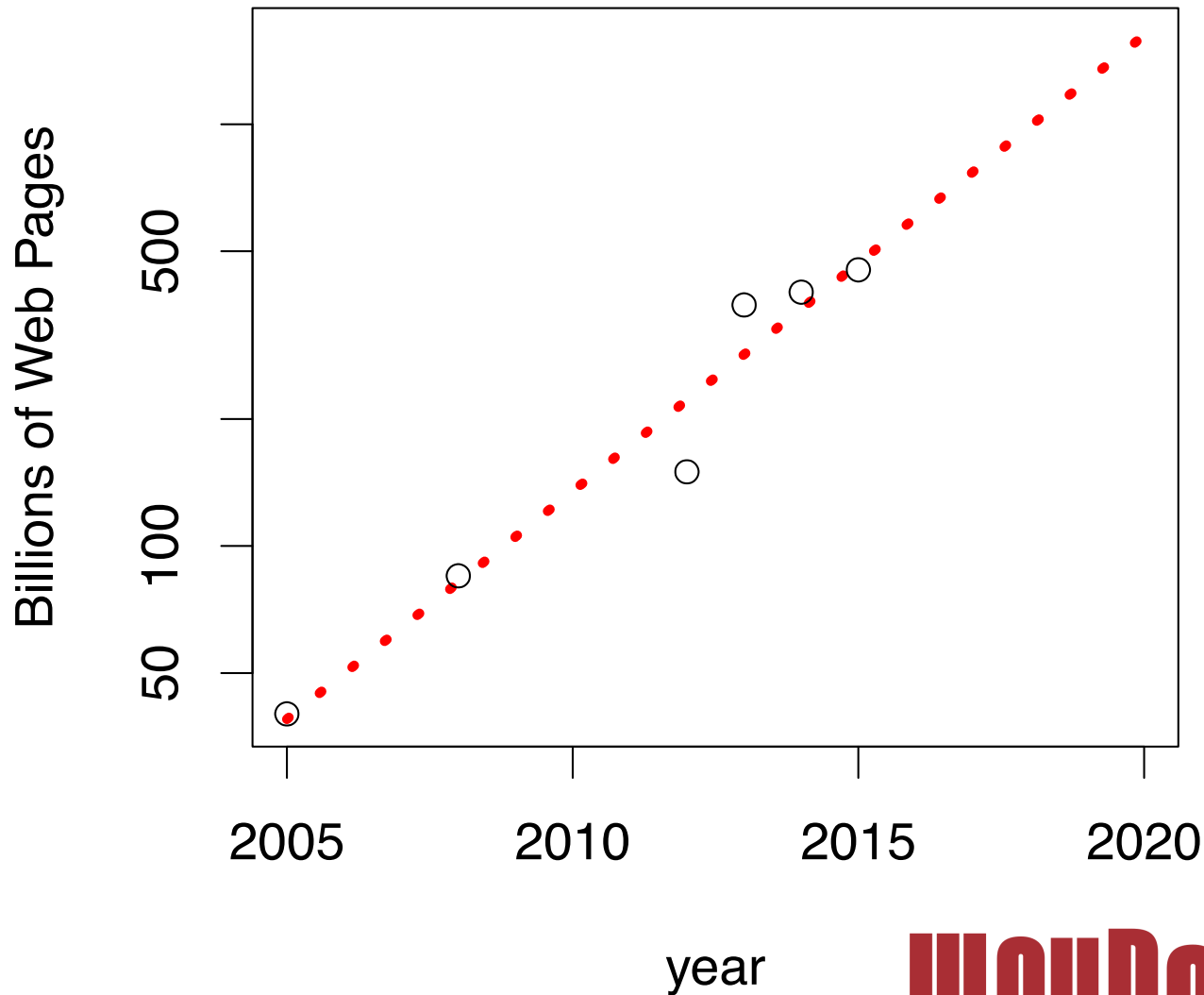- New directions: video, aphasia

# Digital Libraries



AS WE MAY THINK

A TOP U. S. SCIENTIST FORESEES A POSSIBLE FUTURE WORL
IN WHICH MAN-MADE MACHINES WILL START TO THIN

by VANNEVAR BUSH

- As we may think (1945) → WayBack Machine
- Web is large
  - But not the largest corpus
- Healthy eco-system:
  - more readers than writers
  - usage logs >> crawls

# ~1 Trillion Web Pages in the Wayback Mach

**WayBack Machine**

# Google Books:
# 4% of All Books ≈ ½ Trillion Words

# Sanity Check:

## Prague, Czech Republic, Czechoslovakia

# Example of Change: Gay

http://thestarryeye.typepad.com/gay/2015/03/before-gay-meant-gay.html



Look that up in your Funk and Wagnalls!

**1962:** *Let me tell you about a place Somewhere up-a New York way Where the people are so **gay***

# HistWords: Word Embeddings for Historical Text
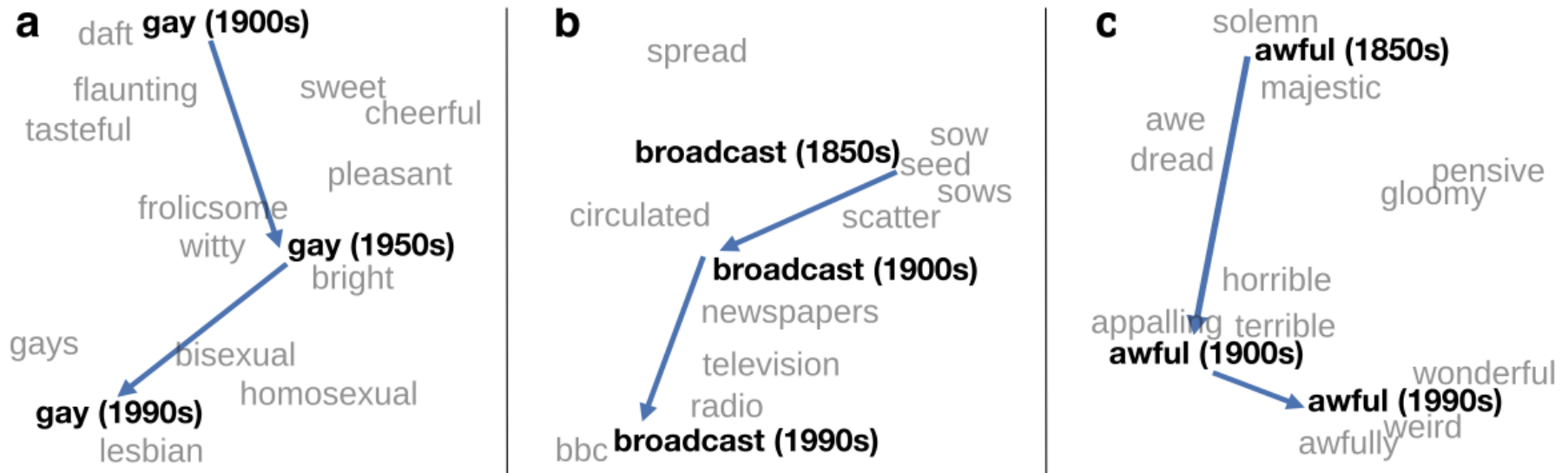## William L. Hamilton, Jure Leskovec, Dan Jurafsky

# Word2vec

## Models

Select one of the available models

English GoogleNews Negative300 ▾

## Nearest words

Given a word, this demo shows a list of other words that are similar to it, i.e. nearby in the vector space.

| doctor | **Show nearest** | **Case sensitive:** ☑ **Top N:** 100 ▾ |

physician
doctors
gynecologist
surgeon
dentist
pediatrician
pharmacist
neurologist
cardiologist
nurse
neurosurgeon
oncologist
dermatologist
urologist
gastroenterologist
psychiatrist

## Word analogy

This demo computes word analogy: the first word is to the second word like the third word is to which word? Try for exa to return *kala (fish)* because fish is to water like birs is to air. Other cases could be for example *sammakko - hyppää - ka* most of the time the analogy does not work particularly well (at least for the Finnish data).

| man | woman | king | **Show** | **Top N:** 10 |

queen
monarch
princess
crown_prince
prince
kings
queens
sultan

# Sampling (Teaser)

✓ Sampling Desiderata: Corpus ≈ Sample
  ✓ Corpus should be representative of population of interest

✓ Quality (representative) vs. Quantity (more):
  ✓ More is more???

➤ **But in a digitized world (where we have it all)**
  – **Do we still need to worry about sampling?**

• **If a corpus is comprehensive,**
  – **does that imply that it is balanced?**

# Google Ngrams: More Recent Estimates are More Reliable
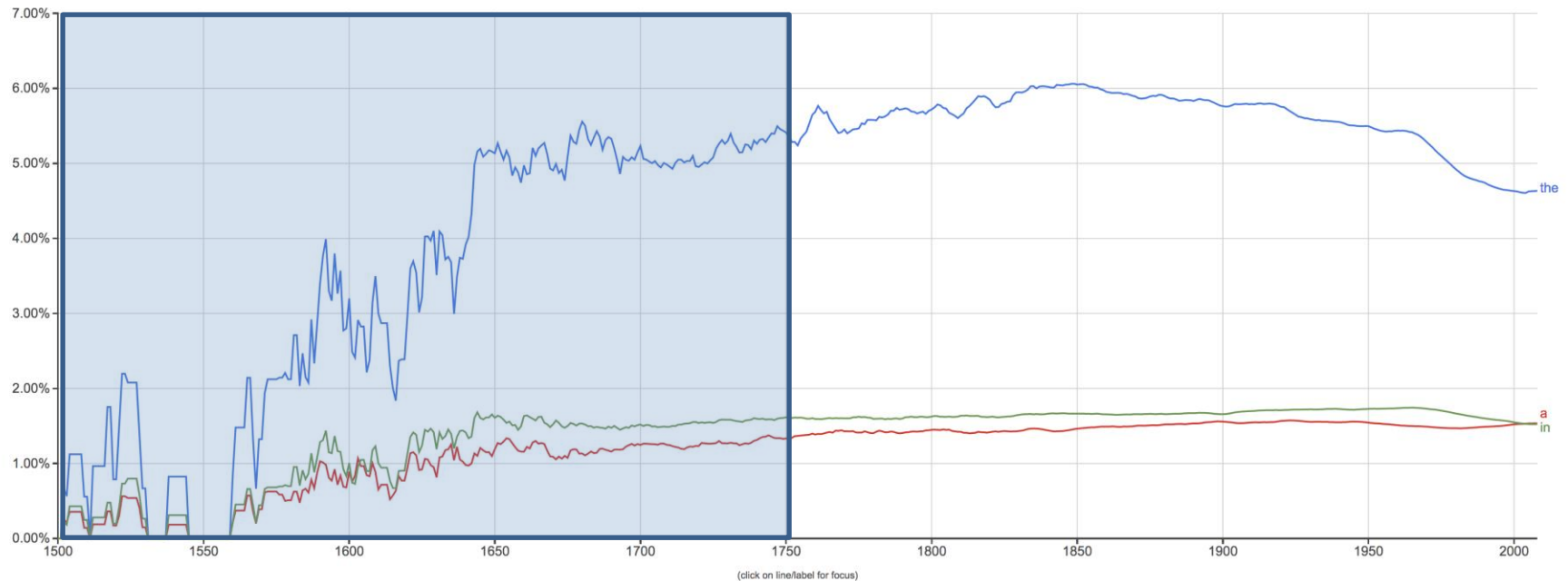## (Most books were published recently)
## 1500 –2000: Frequency of function words: the, a, in
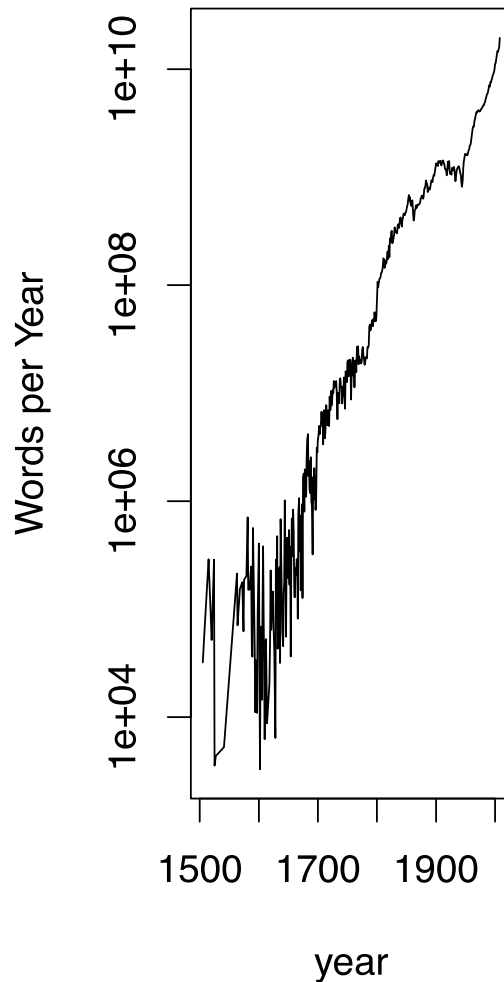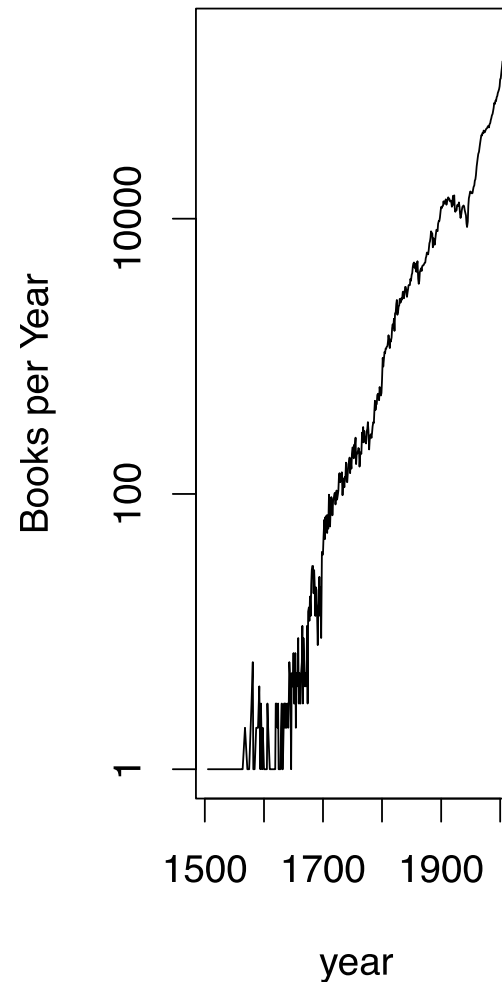
# Google Books: 35% Growth per Decade
# More Recent Estimates are More Reliable

**~0.5 Trillion Words**  **~5 Million Books**

# Conclusions

- Endless debate: quantity vs. quality
- Engineers tend to favor quantity
  - But we need your help to keep the debate going
- In a digitized world, even if we had it all,
  - Sampling is still an issue
- Practical Apps:
  - Predict future (test set)
  - From past (training data)
- Need to sample the past so it is representative of the population of interest (future)

- Future depends on many factors including:
  - Speaker/Audience
    - Adult ≠ Child ≠ Teenager
  - Register: Slang ≠ Formal
  - Spoken ≠ Written
  - Language Change:
    - Past ≠ Future
    - Time & Space