

Computational Linguistics I

CMSC 723 / LING 723 / INST 725

Hal Daumé III [he/him]

27 Aug 2019

(many slides c/o Marine Carpuat)



Welcome!

Note-taker volunteer needed!

Do you take well-organized, comprehensive notes? Do you have good penmanship or do you currently type your notes? Why not get **paid** to share your notes with classmates who are eligible to receive course lecture notes?

If you are interested in providing this much needed service to a fellow student, please go to <https://go.umd.edu/adsNoteTakers> to apply. If you are selected by an eligible student, the Accessibility and Disability Service (ADS) will compensate you with a one-time payment at the end of the semester.

Staff at ADS are available to answer any questions you may have. Feel free to contact us at adsnotetaking@umd.edu.

Course information

- Main stuff: <https://github.com/hal3/cl1f19umd>
- Discussion/grades: <https://umd.instructure.com/courses/1267356>

What is language?

Wikipedia:

“Language is the ability to acquire and use complex systems of communication, particularly the human ability to do so, and a language is any specific example of such a system. The scientific study of language is called linguistics.”

- Computational Linguistics (CL)
 - The science of doing what linguists do with language, but using computers
- Natural Language Processing (NLP)
 - The engineering discipline of doing what people do with language, but using computers
- Speech/Language/Text processing
- Human Language Technology

NLP State of the Art

Still a challenging problem!

AI's Language Problem

“Machines that truly understand language would be incredibly useful. But we don't know how to build them.”

MIT Technology Review

Will Knight, Aug 9, 2016

Many useful applications already exist



What makes a language a *natural* language?

Nuxati Kishelēmienkw, kehëla
wanishi tìlìch nkàski nipai yukwe
ènta kishkwik.
Kèxaptun nkata kèku luwe.

HAI
CAN HAS STDIO?
I HAS A VAR
IM IN YR LOOP
UP VAR!!1
VISIBLE VAR
IZ VAR BIGGER THAN 10? KTHX
IM OUTTA YR LOOP
KTHXBYE

If you were to design a language,
what would you need to do?

wals.info/feature

[wals.info/language/lect/wals code eng](http://wals.info/language/lect/wals_code_eng)

[wals.info/feature/83A#2/20.2/152.8](#)

[wals.info/feature/85A#2/20.3/153.1](#)

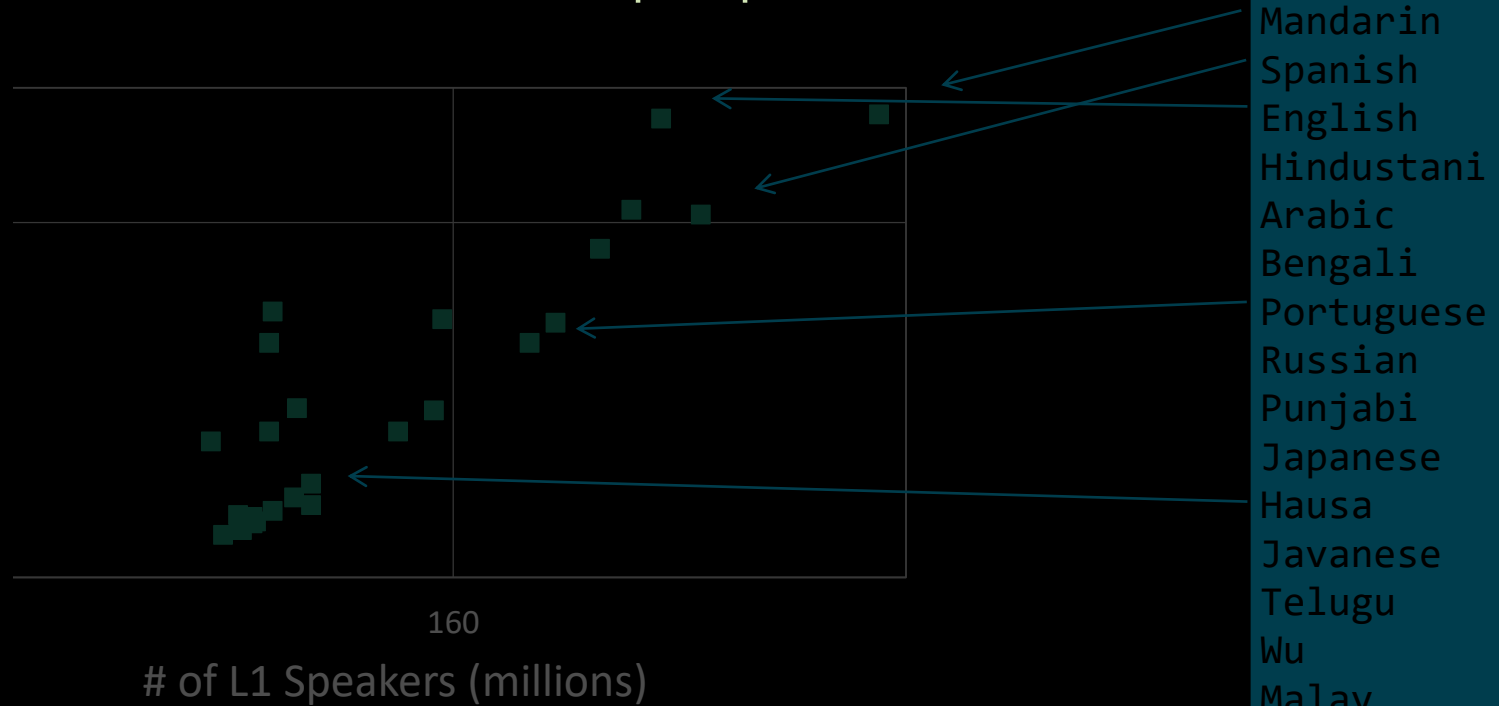
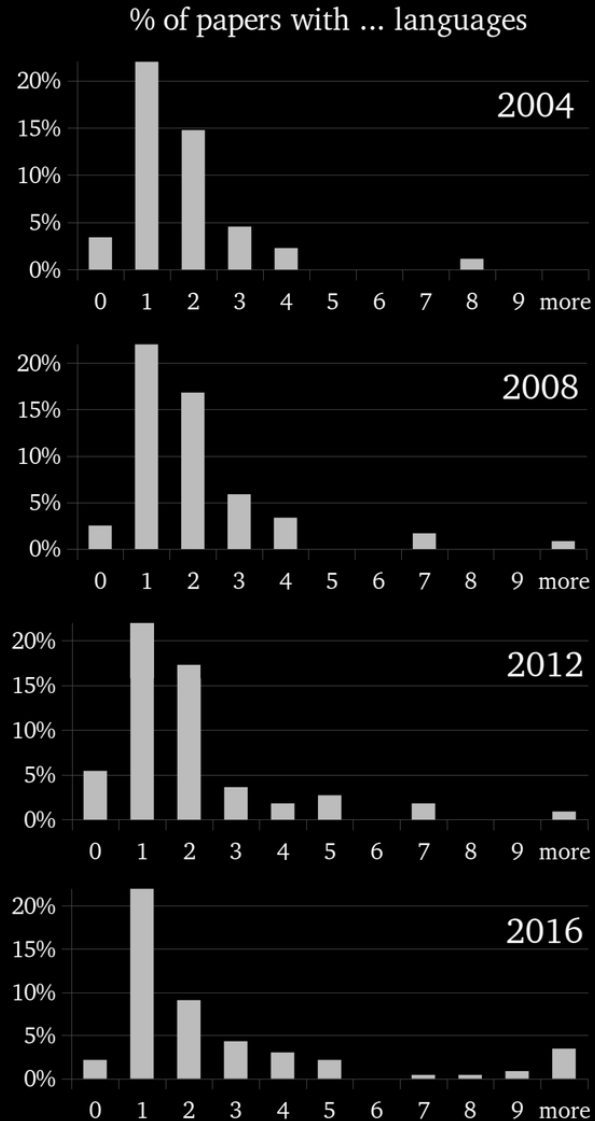
[wals.info/feature/87A#2/20.9/153.3](#)

Language sounds

“th” sound (voiced/voiceless dental fricative)

“r” sound

the world vs research papers

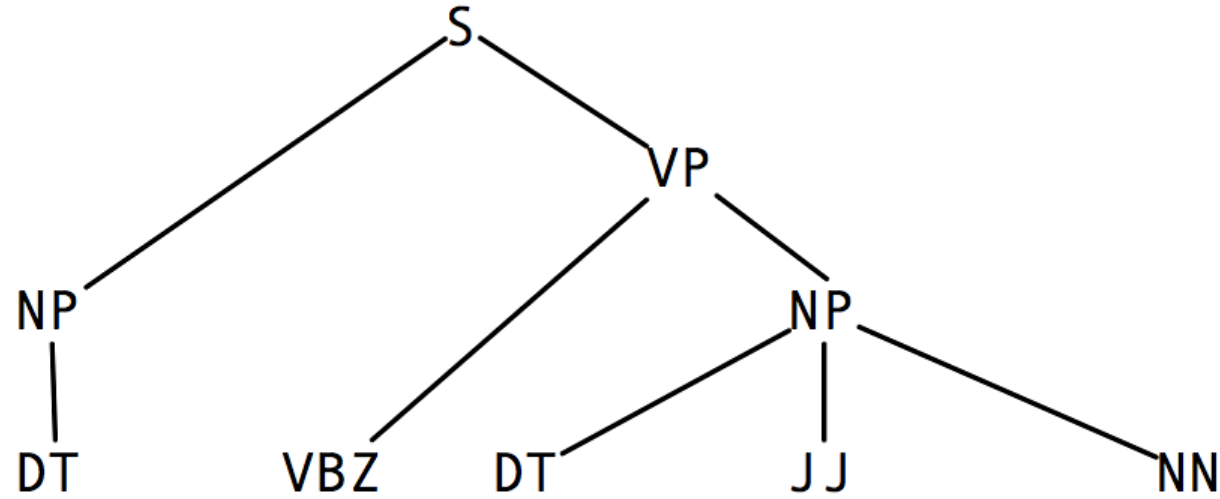


src: Sebastian Mielke
see also: Emily Bender

langscape.umd.edu

Why is NLP hard?

This is a simple sentence



SYNTAX

PART OF SPEECH

WORDS

MORPHOLOGY

SEMANTICS

DISCOURSE

CONTRAST

This is a simple sentence

be
3sg
present

SIMPLE1
having
few parts

SENTENCE1
string of words
satisfying the
grammatical rules
of a language

But it is an instructive one.

Ambiguity

At the word level

- Part of speech
 - [V Duck]!
 - [N Duck] is delicious for dinner.
- Word sense
 - I went to the bank to deposit my check.
 - I went to the bank to look out at the river

Ambiguity

At the syntactic level

- PP Attachment ambiguity
 - I saw the man on the hill with the telescope
- Structural ambiguity
 - I cooked her duck
 - Visiting relatives can be annoying
 - Time flies like an arrow

Ambiguity

- Quantifier scope
 - Everyone on the island speaks two languages.
- Hard cases require world knowledge, understanding of speaker goals
 - The city council denied the demonstrators the permit because they advocated violence
 - The city council denied the demonstrators the permit because they feared violence

Some newspaper headlines

- Iraqi Head Seeks Arms
- Teacher Strikes Idle Kids
- Kids Make Nutritious Snacks
- Stolen Painting Found by Tree
- Local HS Dropouts Cut in Half
- Enraged Cow Injures Farmer with Ax
- Hospitals are Sued by 7 Foot Doctors
- Ban on Nude Dancing on Governor's Desk
- Scientists study whales from space

Despite ambiguity, language is predictable

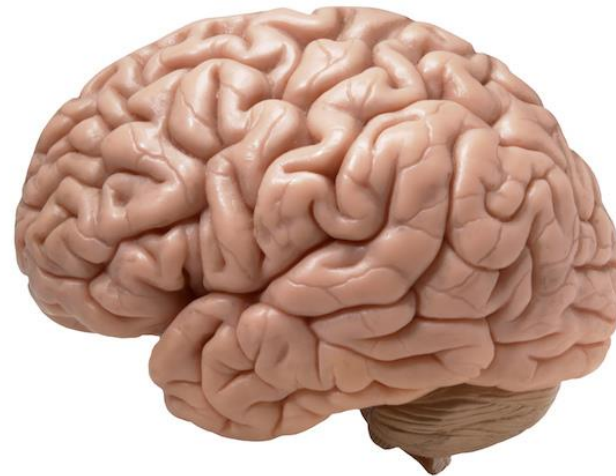
I like my coffee with cream and asparagus

This is crummy weather for Collocation

➤ The brain uses this information!

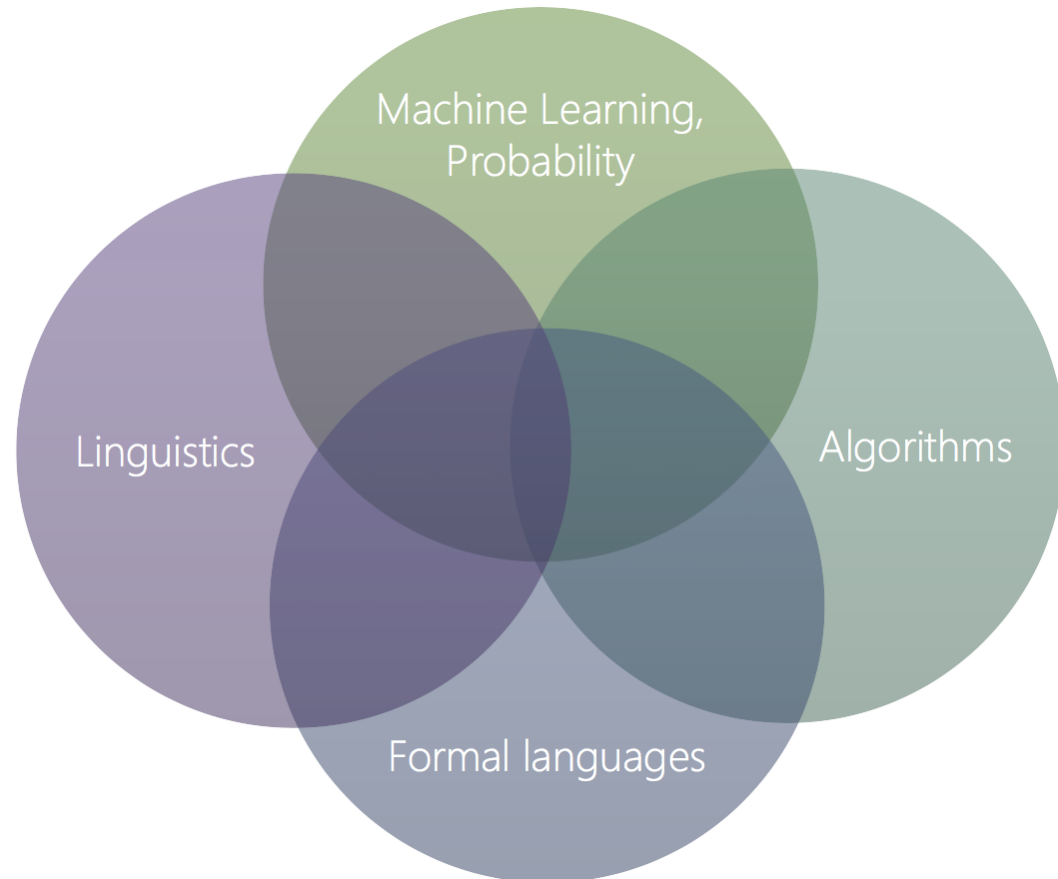
➤ Can we use predictability
to make decisions *before*
all of the input is observed?

YES!!!



Ambiguity

- NLP challenge: how can we model ambiguity, and choose the correct analysis in context?
- Approach: learn from data



Word counts

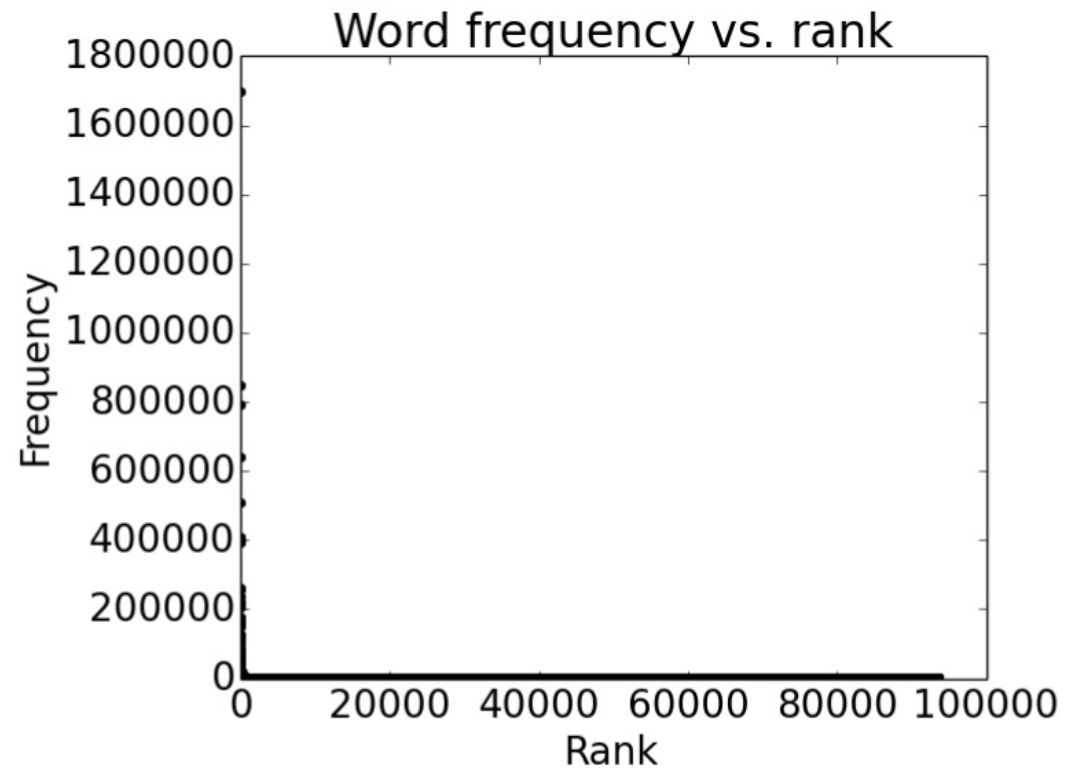
- Most frequent words in the English Europarl **corpus**
- (out of 24M word **tokens**)

any word		nouns	
Frequency	Token	Frequency	Token
1,698,599	the	124,598	European
849,256	of	104,325	Mr
793,731	to	92,195	Commission
640,257	and	66,781	President
508,560	in	62,867	Parliament
407,638	that	57,804	Union
400,467	is	53,683	report
394,778	a	53,547	Council
263,040	I	45,842	States

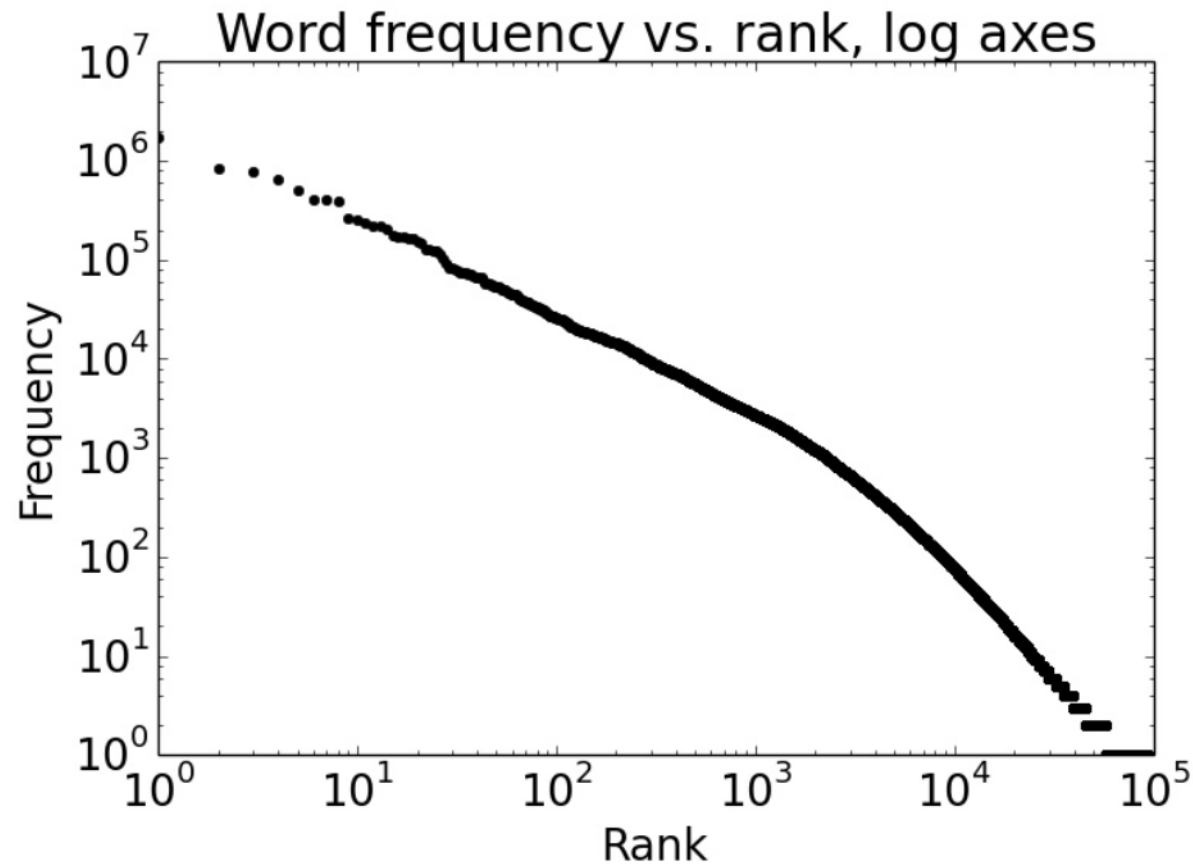
Word counts

- But also, out of the 93,638 distinct words (word **types**), 36,231 occur only once
 - cornflakes, mathematicians, fuzziness, jumbling
 - pseudo-rapporteur, lobby-ridden, perfunctorily,
 - Lycketoft, UNCITRAL, H-0695
 - policyfor, Commissioneris, 145.95, 27a

Plotting word frequencies



Plotting word frequencies (with log-log axes)



Zipf's law

$$f \times r \approx k$$

- f = frequency of a word
- r = rank of a word (if sorted by frequency)
- k = a constant

Zipf's law: implications

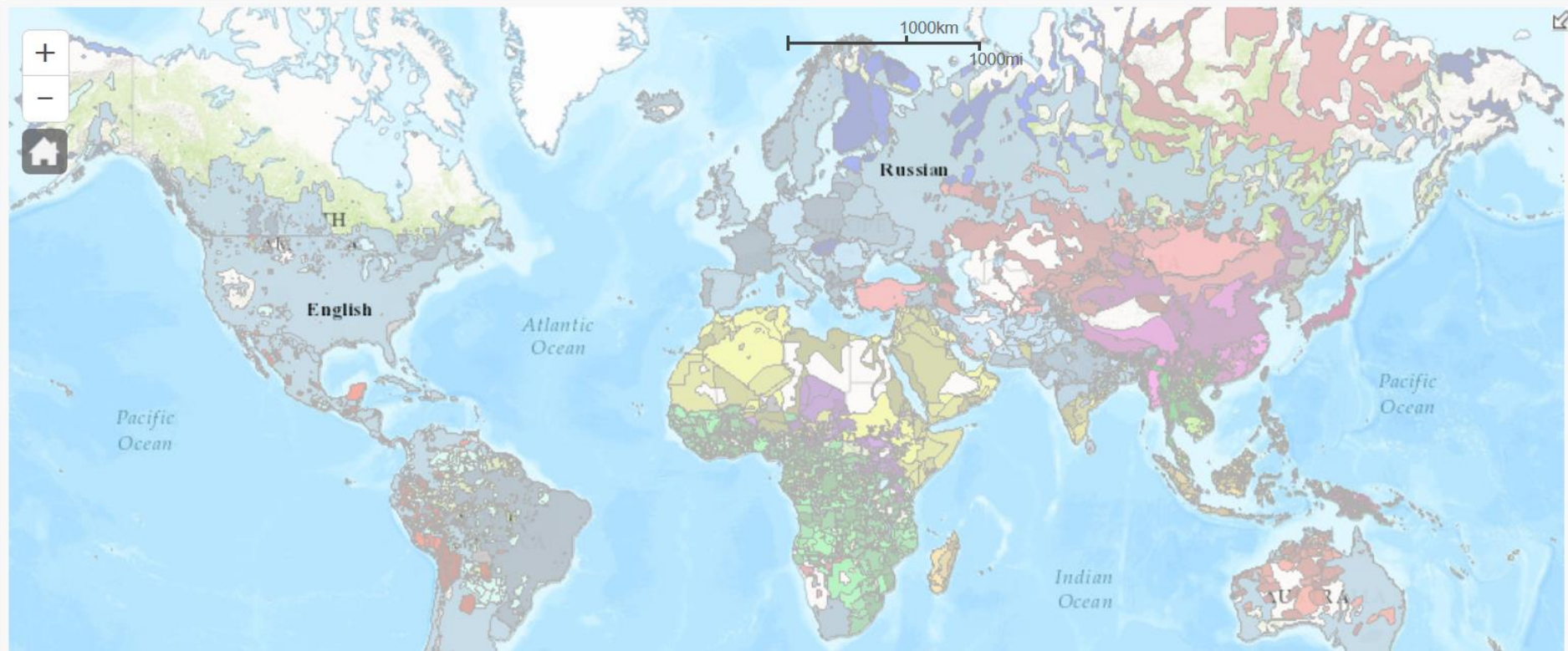
- Even in a very large corpus, there will be a lot of infrequent words
- The same holds for many other levels of linguistic structure
- Core NLP challenge: we need to estimate probabilities or to be able to make predictions for things we have rarely or never seen

Variation and Expressivity

- The same meaning can be expressed with different forms
 - I saw the man
 - The man was seen by me
- She needed to make a quick decision in that situation
- The scenario required her to make a split-second judgment



Search for a language, dialect name or major city...



6,800 living languages
600 with written tradition
100 spoken by 95% of population

even one "language" isn't

How do you pronounce *caramel*?

How would you address a group of two or more people?

What do you call the small gray bug that curls up

SFR 08:36 100 %

Hey Siri translate you're welcome in mandarin

Tap to Edit

I can't translate Canadian English yet. Sorry about that.

?



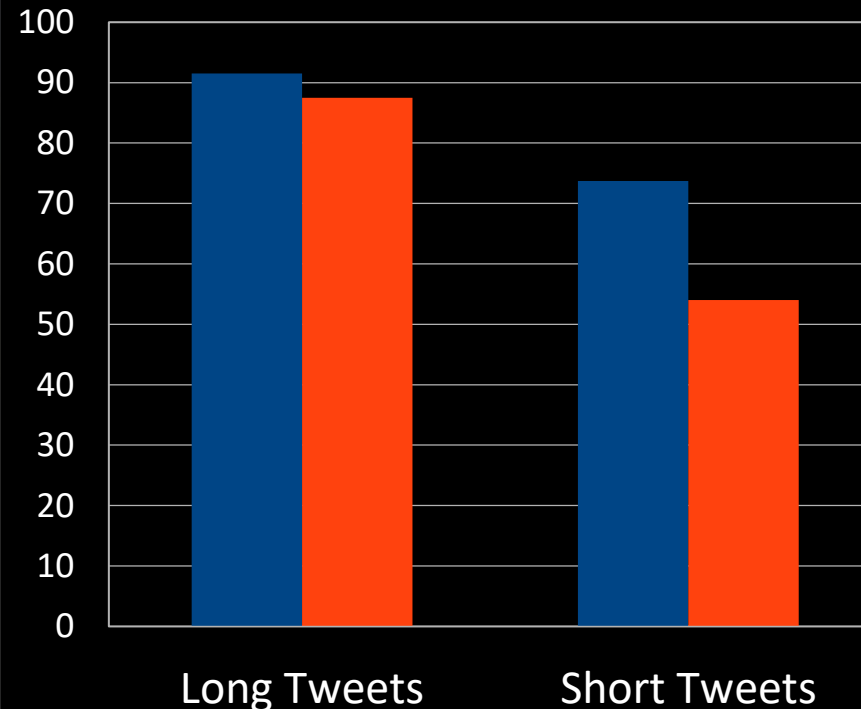
4. A greeting you might use?

Alre 6. A word meaning "good"?

Hoo Sound.

Alrig Barrie.

Accuracy of identifying Tweets as "English"

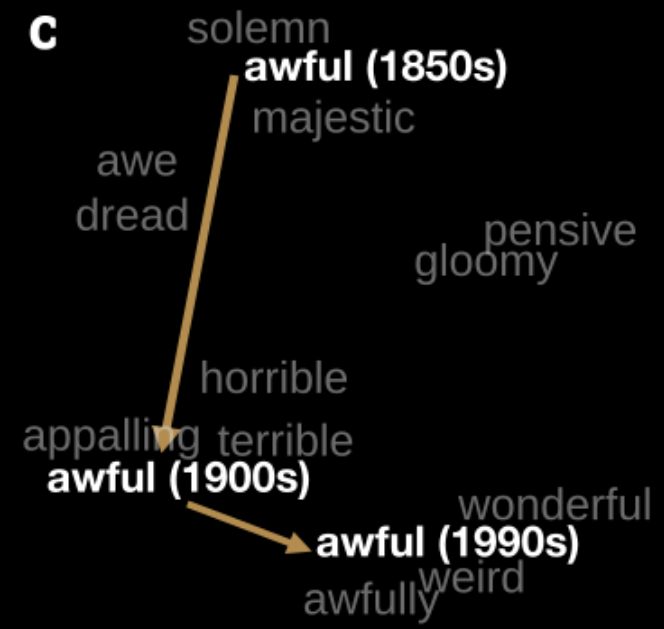
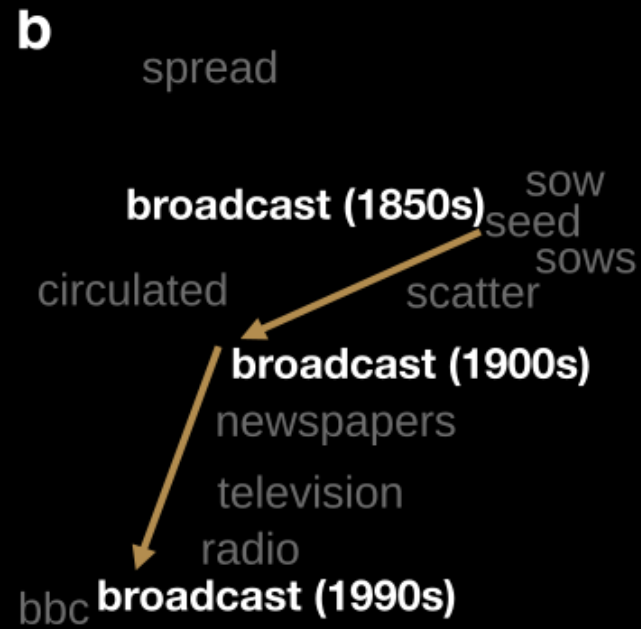
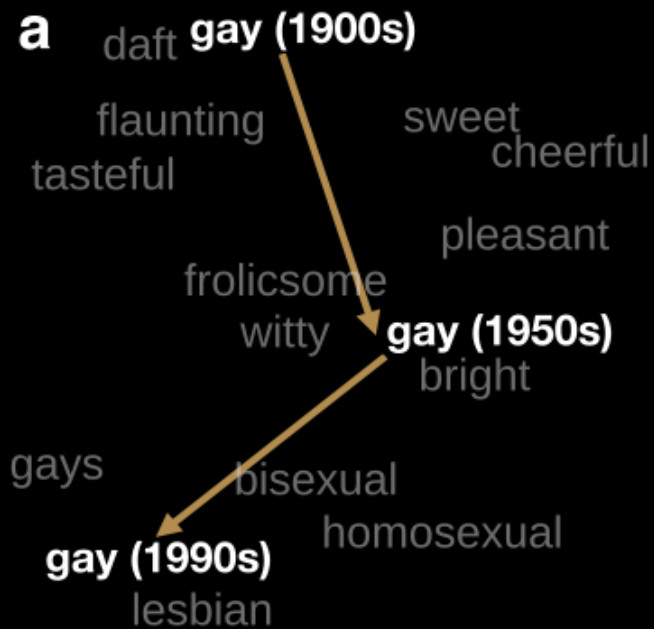


■ "White" English

■ African American English

Blodgett+
Green+
O'Connor
EMNLP 2016

and things don't stay put....



Social Impact

- NLP experiments and applications can have a direct effect on individual users' lives
- Some issues
 - Privacy
 - Exclusion
 - Overgeneralization
 - Dual-use problems

Today

- Levels of linguistic analysis in NLP
 - Morphology, syntax, semantics, discourse
- Why is NLP hard?
 - Ambiguity
 - Sparse data
 - Zipf's law, corpus, word types and tokens
 - Variation and expressivity
 - Social Impact

Before next class

- Read the syllabus
- Make sure you have access to canvas
- Read SLP3 6.2—6.5