

# Dialects in NLP

CS 685, Spring 2021

Advanced Topics in Natural Language Processing

<http://brenocon.com/cs685>

[https://people.cs.umass.edu/~brenocon/cs685\\_s21/](https://people.cs.umass.edu/~brenocon/cs685_s21/)

Brendan O'Connor

College of Information and Computer Sciences

University of Massachusetts Amherst

# Presentations next Monday!

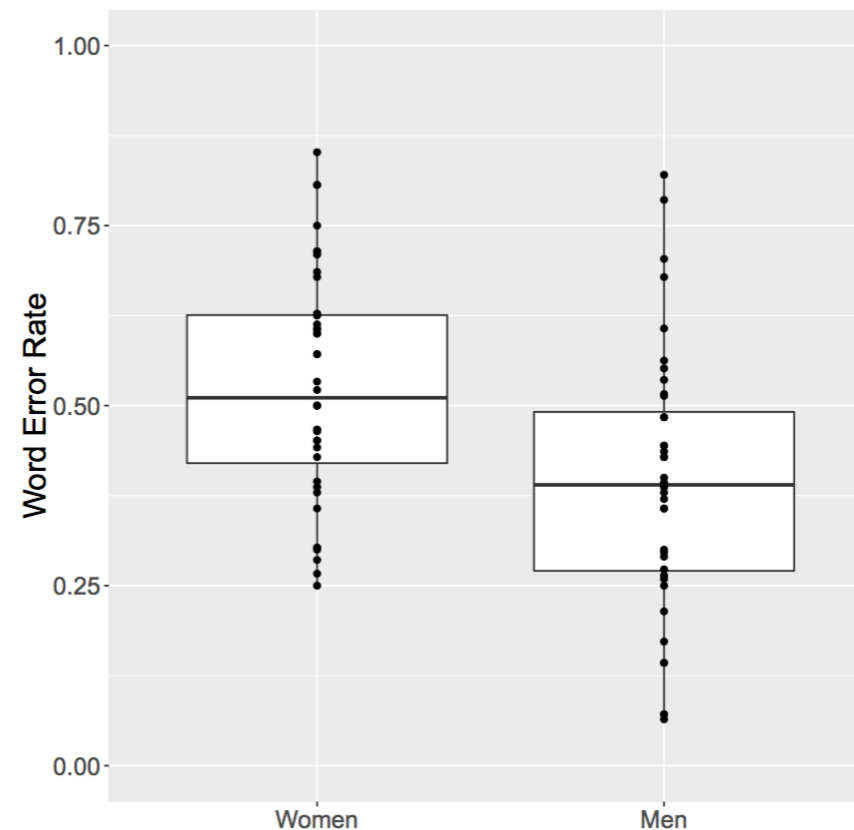
- One slide, no build/animations please - will send out a Google Slides URL
- 90 seconds per group
  
- Final report due **May 12** - last day of finals
  - Will also post optional "HW3" extra credit questions - due same day

- Language is socially situated
  - **By** and **for** communicators
    - Today - focus on *dialect*
  - and **about** people

# Author-conditional NLP disparity

- Language technologies analyze the linguistic behavior of people
- Language is affected by social context and attributes

Example: gender bias in  
YouTube aut captions  
*[Tatman 2017]*



- What information can a user access?
- Whose voices are heard?

# Variation in language

- Social factors drive language change and correlate with language varieties: by geography, ethnicity, gender, class...

- e.g. socioeconomic class: Labov (1966) finds more (r)-usage at more expensive New York department stores

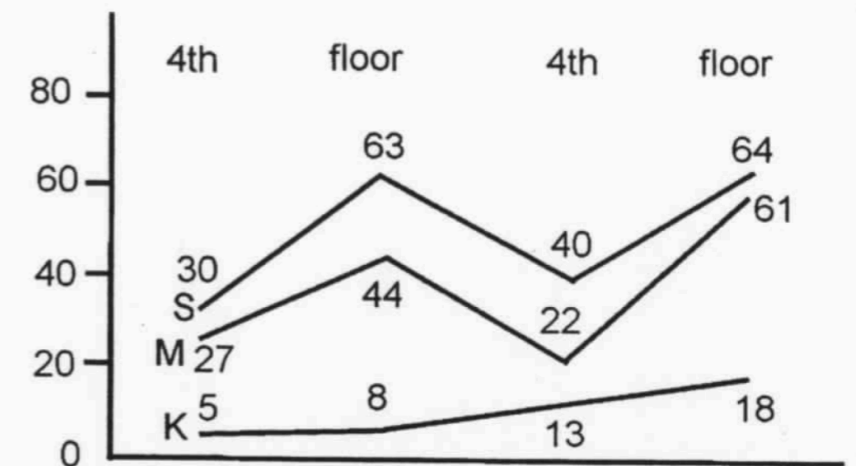


Figure 13.2: Percentage of all (r-1) by store for four positions (S = Saks, M = Macy's, K = Kleins)

# Variation in language

- Social factors drive language change and correlate with language varieties: by geography, ethnicity, gender, class...
- Arbitrariness of language designations - "dialect is a language with an army and navy"
  - Cantonese, Moroccan Arabic ...
  - Swedish, Danish...
- Distinct language varieties often associated with social groups and segregated communities
- We all know NLP domain adaptation is hard. How does this affect NLP performance?

# Dialectal NLP

- 1. Identify the dialect & make a corpus
  - Through author-level metadata
  - Through in-text linguistic features
- 2. Evaluate NLP performance on the dialect
- 3. Adapt systems to work well on the dialect
  - Turns out to be tough!

# Dialect in social media

SAE:

*he is woke af*



he woke af smart af educated af daddy af  
coconut oil using af GOALS AF & shares food af

RETWEETS	LIKES	
3	42	

1:08 AM - 8 Jul 2016

← ↻ 3 ❤️ 42 ⋮

- **Why is social media different?**

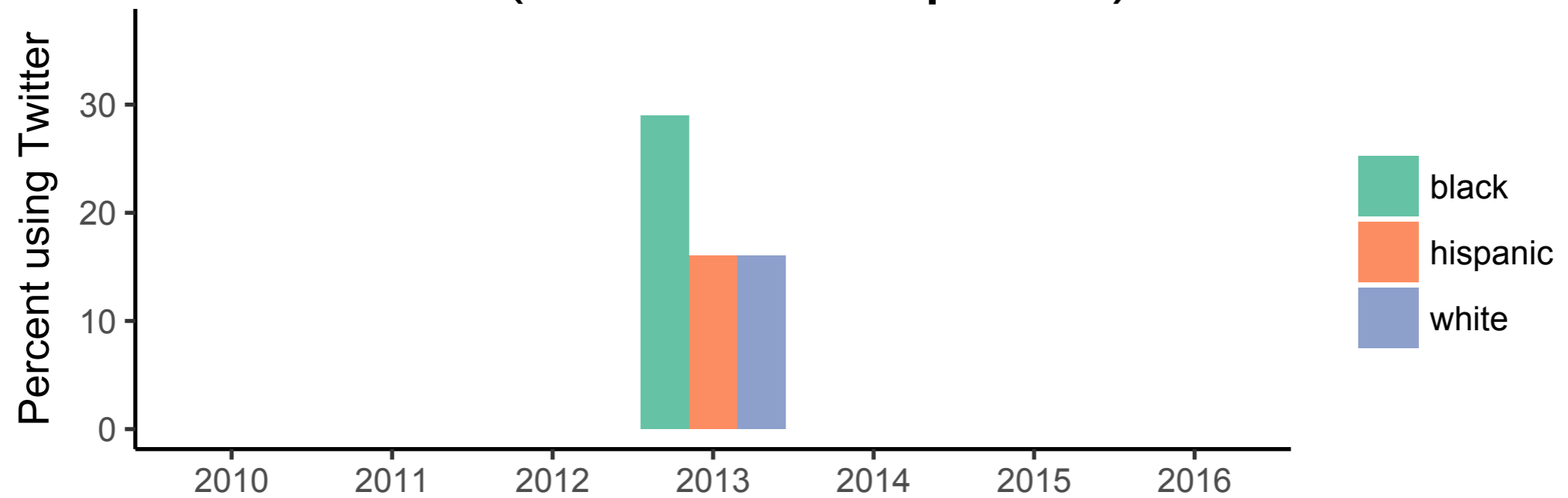
- Internet speech?
- Pre-existing dialectal English?
  - Geographic patterns of word usage often reveal relationships to race, ethnicity etc.
  - African-American English in Twitter  
*[Eisenstein 2013, Jorgensen et al. 2015, Jones 2015]*



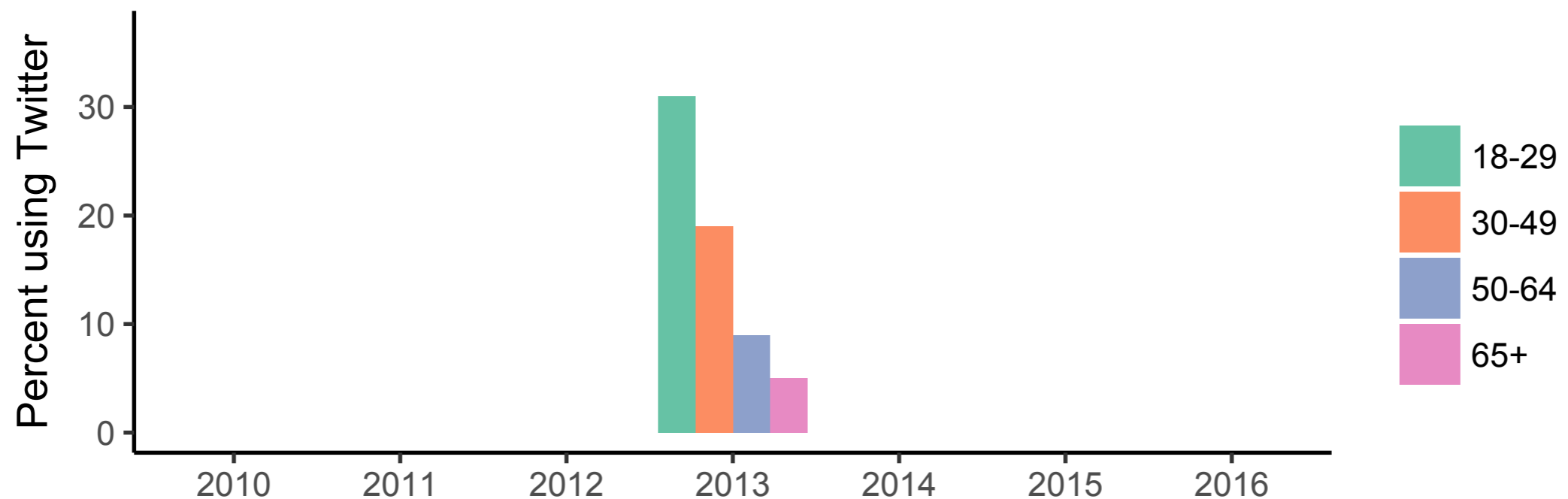
# Youth, minorities on Twitter

[Pew Research]

## P(use twitter | race)



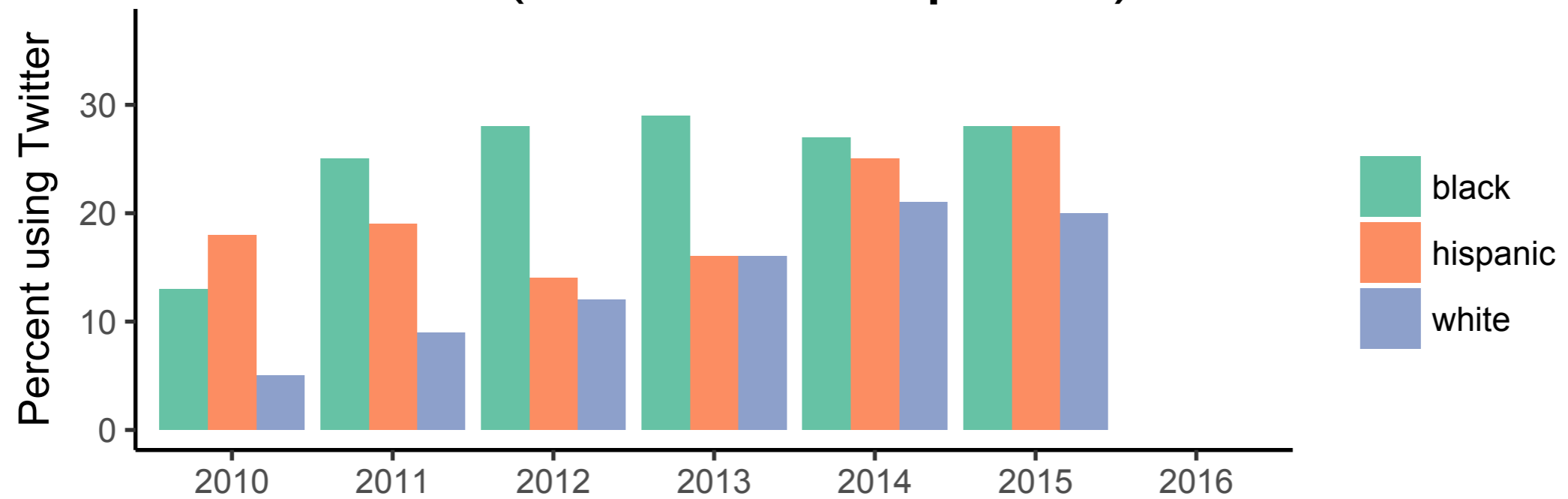
## P(use twitter | age)



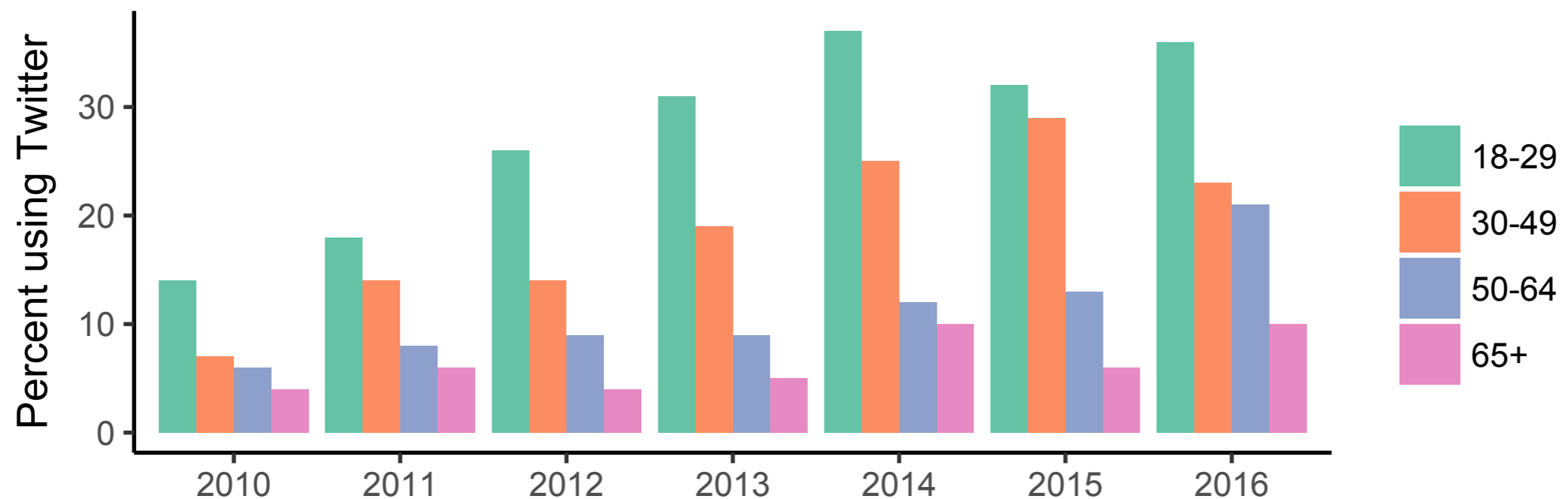
# Youth, minorities on Twitter

[Pew Research]

## P(use twitter | race)



## P(use twitter | age)



# (Immediate?) future auxiliaries

gonna gunna gona gna guna gnna ganna qonna gonnna gana  
qunna gonne goona gonnaa g0nna goina gonnah goingto  
gunnah gonaa gonan gunnna going2 gonnna gunnaa gonny  
gunaa quna goonna qona gonns goinna gonnae qna gonnaaa  
gnaa

tryna gon finna bouta trynna boutta gne fina gonn tryina  
fenna qone trynaa qon boutaa funna finnah bouda boutah  
abouta fena bouttah boudda trinna qne finnaa fitna aboutta  
goin2 bout2 finnna trynah finaa ginna bouttaa fna try'na g0n  
trynn tyrna trna bouto finsta fna tranna finta tryinna finnuh  
tryingto boutto

- finna ~ “fixing to”
- tryna ~ “trying to”
- bouta ~ “about to”

# African American English

## A Linguistic Introduction

Lisa J. Green

### 2.4

#### Preverbal markers: *finna*, *steady*, *come*

Markers *finna*, *steady* and *come* have been identified in AAE, but they have not been analyzed to the extent that markers such as aspectual *be*, remote past *BIN* and *dən* have been analyzed. There are some descriptions of them in the literature, which will be cited in the summary of each preverbal marker. Also, note that lexical entries for *steady* and *come* are given in chapter 1.

#### **Finna**

*Finna* (including variants *fixina*, *fixna* and *fitna*) indicates that the event is imminent; it will happen in the immediate future. It precedes non-finite verbs, which are not marked for tense and agreement. Sentences in which this marker occurs are given below:

- (77) a. I don't know about you, but I'm **finna leave**.  
'I don't know about you, but I'm getting ready/about to leave'
- b. Y'all **finna eat**?  
'Are you getting ready/about to eat?'
- c. She was **finna move** the mattress herself when I got there.  
'She was getting ready/about to move the mattress when I got there'
- d. Oh-oh they pulling they coats off. That mean they **fixna kill** us or something.  
(attested)

# Associating geolocated tweets with demographics

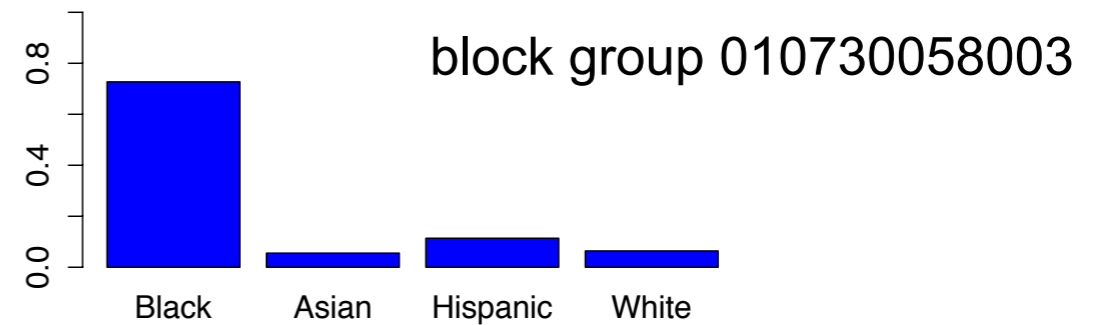
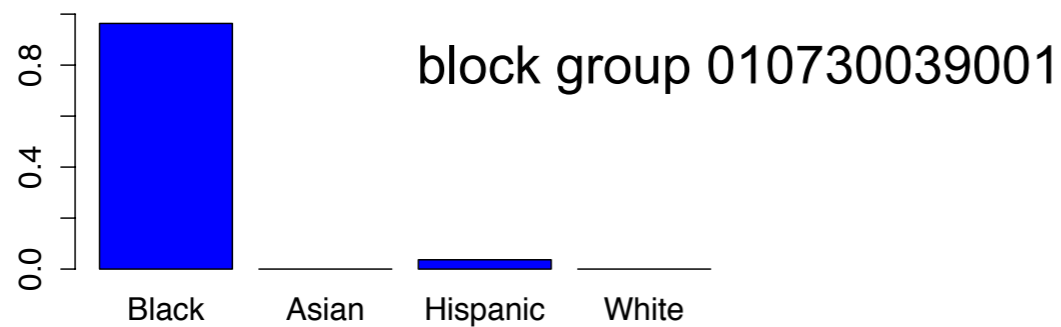


**King Me**  
@tmac\_datboi

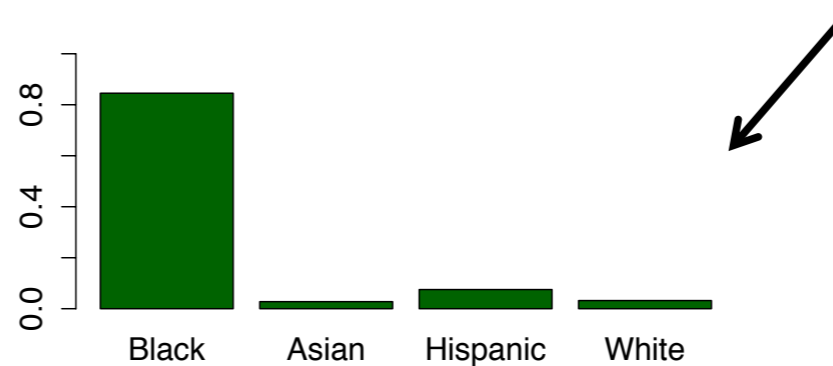


he woke af smart af educated af daddy af  
coconut oil using af GOALS AF & shares food af

Bored af den my phone finna die!!!



$$\pi_{\text{user}} =$$

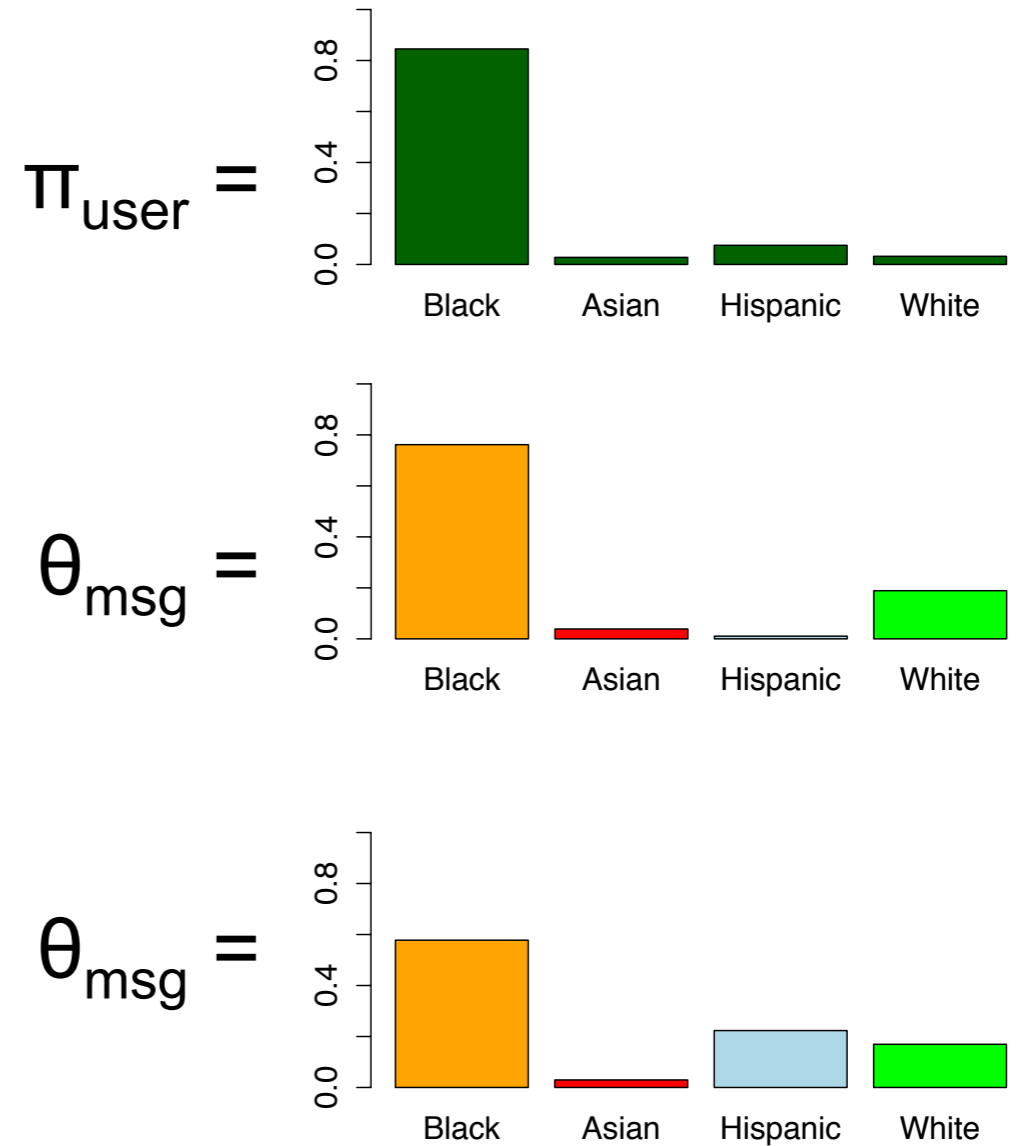


# Mixed membership model

$$\theta_{\text{msg}} \sim \text{Dir}(\alpha\pi), \quad z \sim \theta_{\text{msg}}, \quad w \sim \phi_z$$

he woke af smart af educated af daddy af  
 coconut oil using af GOALS AF & shares food af

Bored af den my phone finna die!!!



# Validation: Phonology

- For every word in vocabulary  $w$  and topic  $k$ , calculate

$$r_k(w) = \frac{p(w|z = k)}{p(w|z \neq k)}$$

- Calculate  $r_{AA}(w)$  for 31 phonological variants illustrated through nonstandard spellings
- For 30/31 variants:  $r \geq 1$

<b>AAE</b>	<b>Ratio</b>	<b>SAE</b>
sholl	1802.49	sure
iont	930.98	I don't
wea	870.45	where
talmbout	809.79	talking about
sumn	520.96	something

Model also has well-known syntactic phenomena in AAE (e.g. null copulas)

# Validation: Syntax

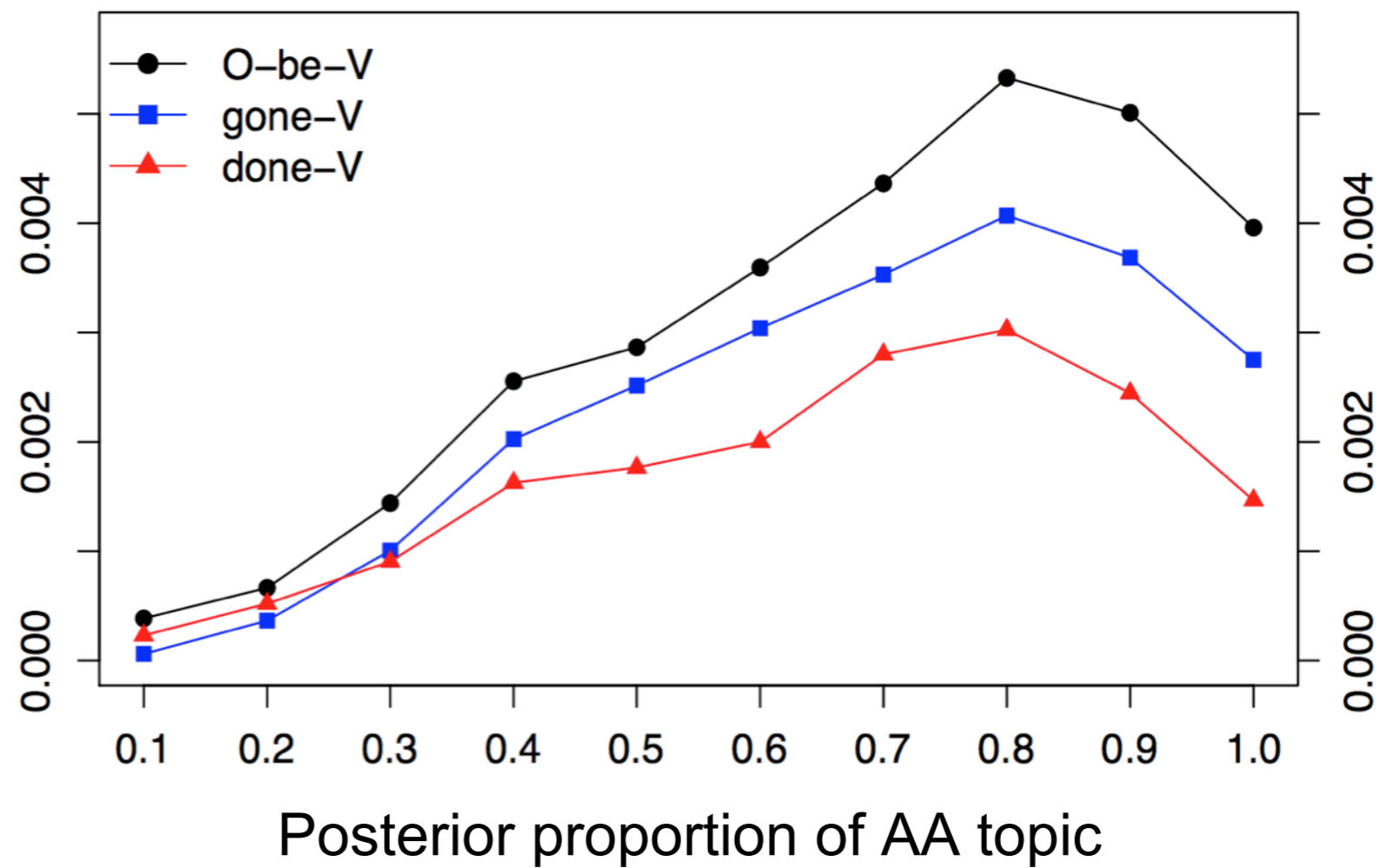
- Select 3 well-known AAE verbal markers
- Search for sequences of unigrams and POS tags

<b>Construction</b>	<b>Example</b>
<i>O-be/b-V</i>	<i>I be tripping bruh</i>
<i>gone/gne/gon-V</i>	<i>Then she gon be single Af</i>
<i>done/dne-V</i>	<i>I done laughed so hard that I'm weak</i>



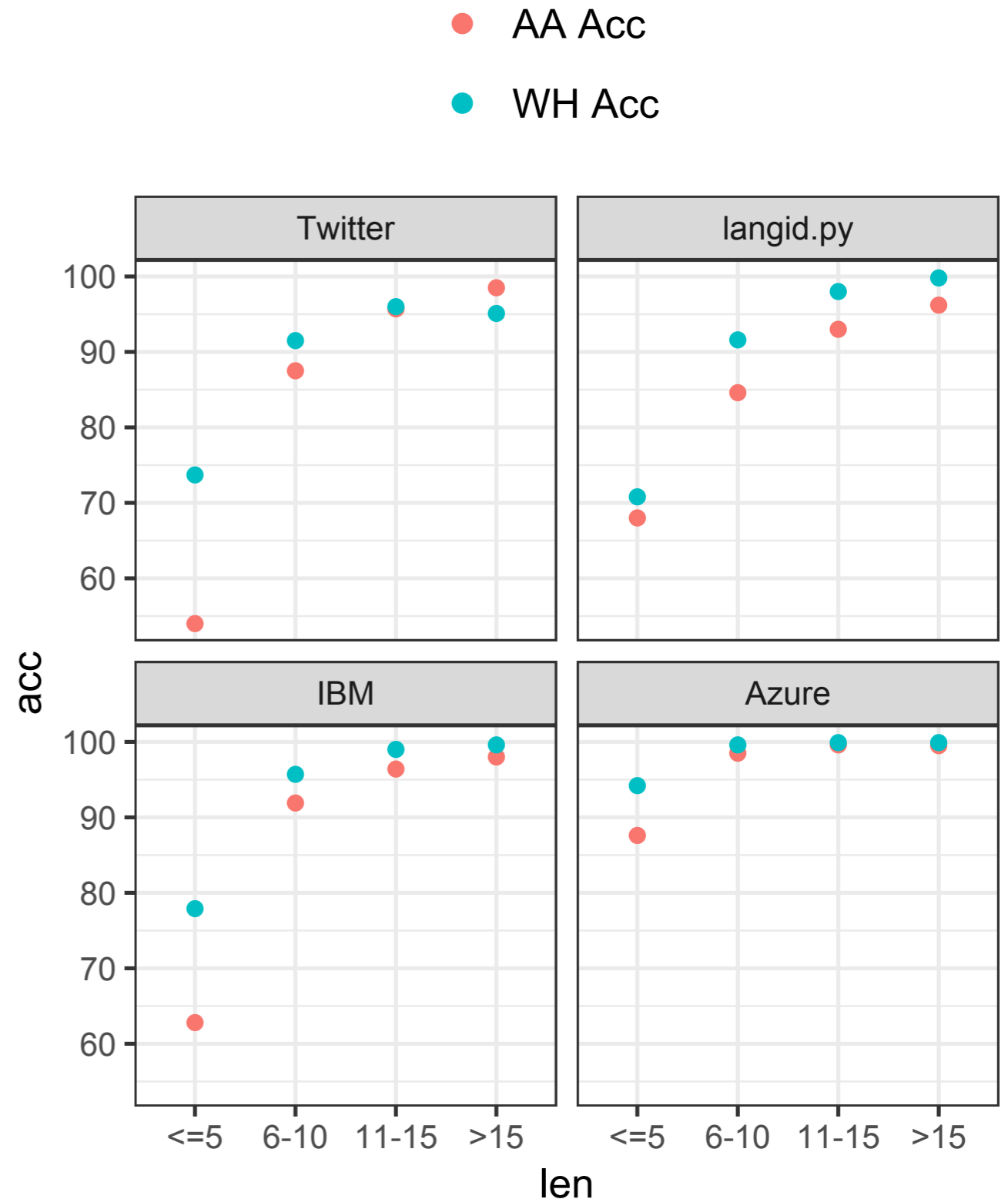
# Validation: Syntax

Proportion of tweets  
with construction

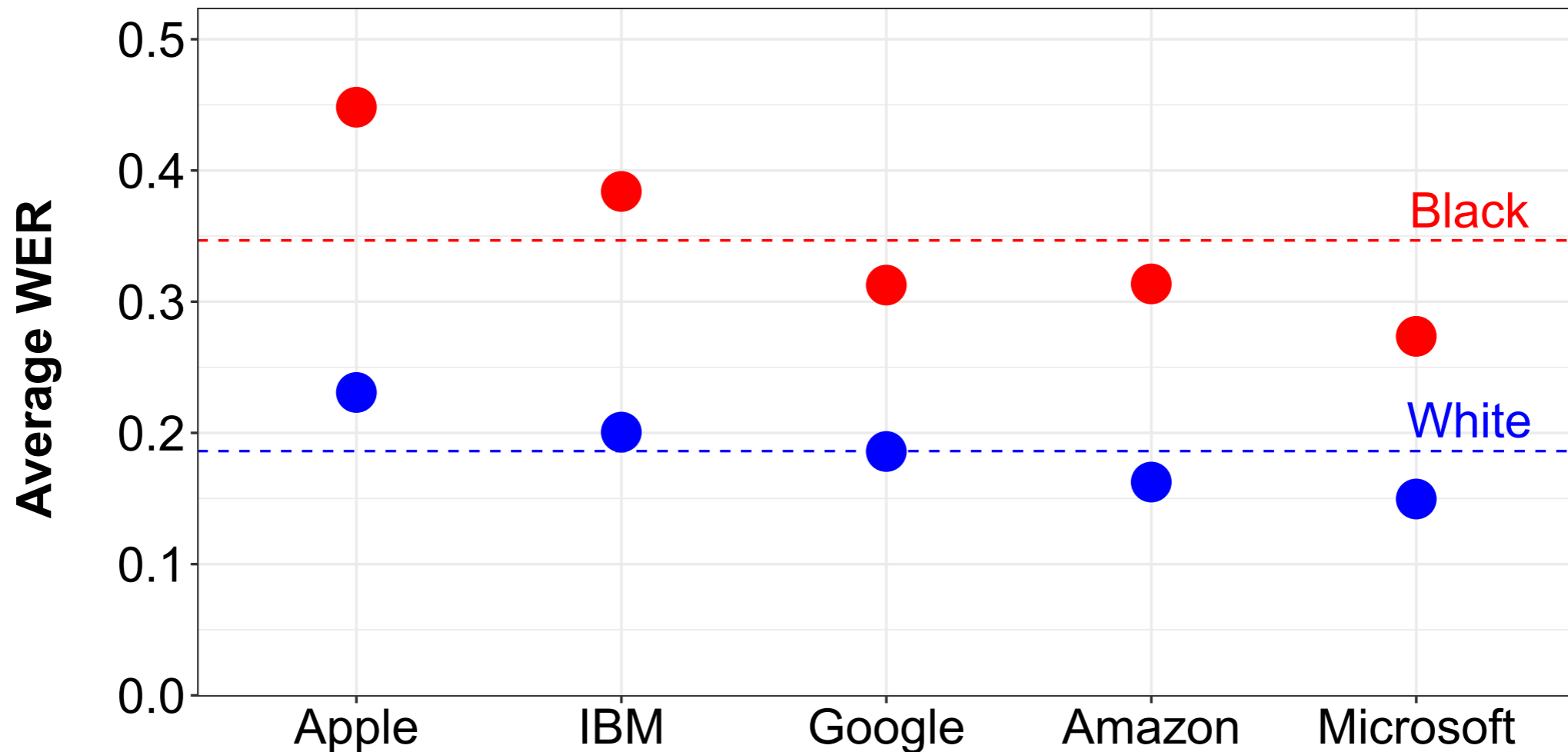


# Racial disparity in accuracy

- $p(\text{correct} \mid \text{Wh})$  vs  $p(\text{correct} \mid \text{AA})$
- Assess disparity in
  - langid.py: popular open-source system [Lui and Baldwin, 2012]
  - Twitter (in metadata)
  - IBM, Microsoft



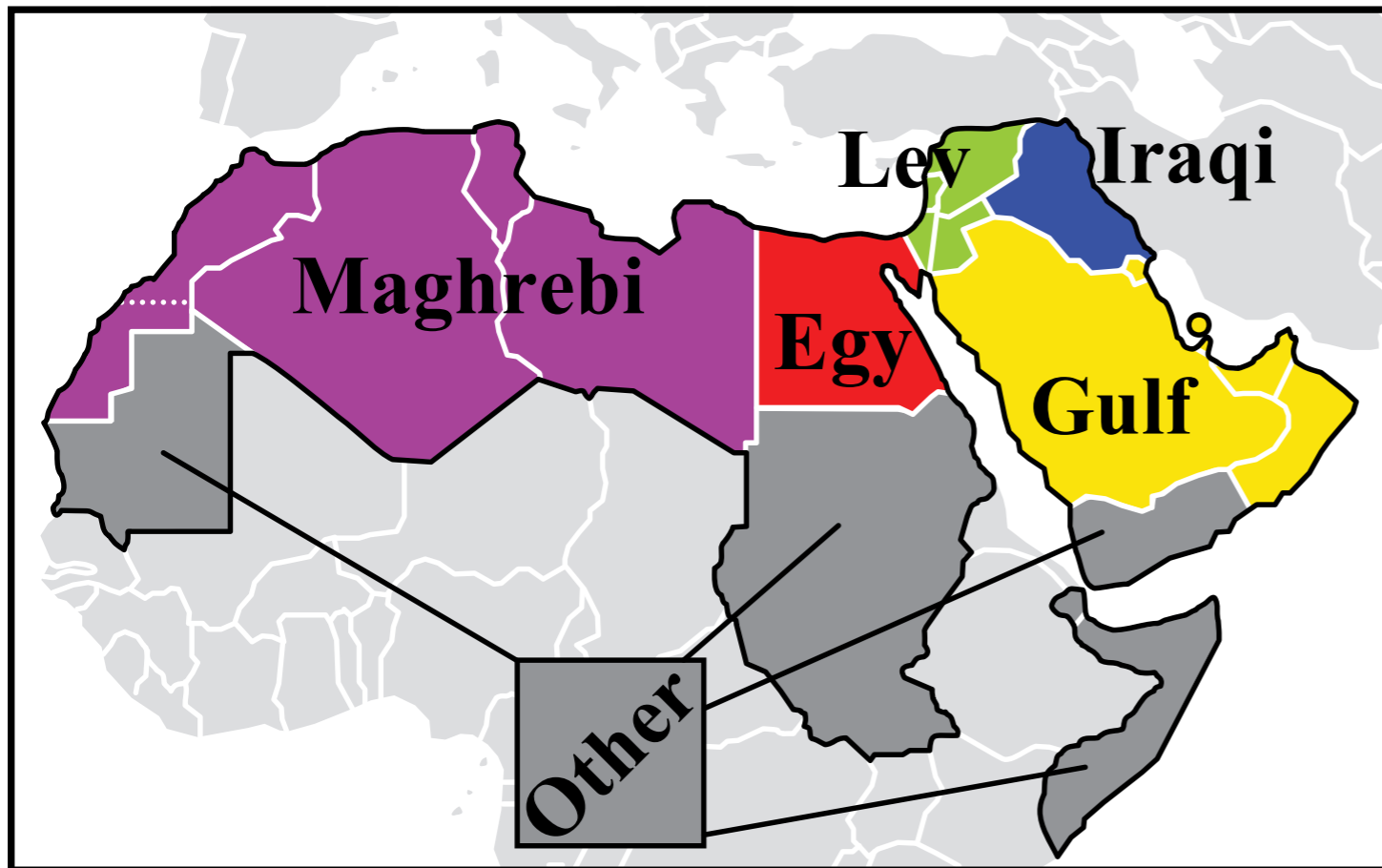
# ASR and AAE



- From the Corpus of African-American Language (CORAAAL): audio recordings of interviews

# Direct identification of dialects

- Dialectal Arabic - annotated & classified from web text [Zaidan and Callison-Burch, 2006]



English	MSA	LEV	GLF	EGY
Book	<i>ktAb</i>	<i>ktAb</i>	<i>ktAb</i>	<i>ktAb</i>
Year	<i>snħ</i>	<i>snħ</i>	<i>snħ</i>	<i>snħ</i>
Money	<i>nqwd</i>	<i>mSAry</i>	<i>flws</i>	<i>flws</i>
Come on!	<i>hyA!</i>	<i>ylA!</i>	<i>ylA!</i>	<i>ylA!</i>
I want	<i>Aryd</i>	<i>bdy</i>	<i>Abγý</i>	<i>ςAyz</i>
Now	<i>AlĀn</i>	<i>hlq</i>	<i>AlHyn</i>	<i>dlwqt</i>
When?	<i>mtý?</i>	<i>Aymtý?</i>	<i>mtý?</i>	<i>Amtý?</i>
What?	<i>mAA?</i>	<i>Ayš?</i>	<i>wš?</i>	<i>Ayh?</i>
I drink	<i>Āšrb</i>	<i>bšrb</i>	<i>Ašrb</i>	<i>bšrb</i>
He drinks	<i>yšrb</i>	<i>bšrb</i>	<i>yšrb</i>	<i>byšrb</i>
We drink	<i>nšrb</i>	<i>bnšrb</i>	<i>nšrb</i>	<i>bnšrb</i>

# Dialect ID from ling. features

- Linguistic knowledge-driven approach: identify sentences with particular *linguistic features*
- Supervise BERT fine-tuning with *minimal pairs*
- Application: identify Indian English

## 176. Deletion of copula *be*: before NPs

Feature area: Agreement

Typical example: He  $\emptyset$  a good teacher.

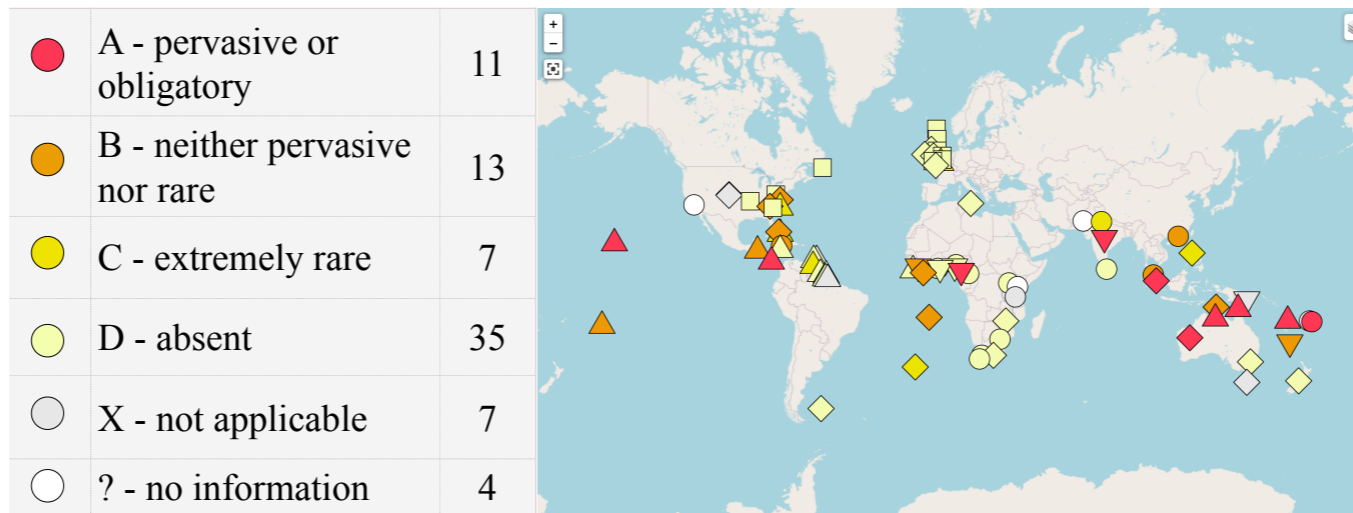


Figure 1: An example dialect feature from the Electronic World Atlas of Varieties of English (eWAVE).<sup>1</sup>

Feature	Example	Count of Instantiations	
		Lange (2012)	Our data
ARTICLE OMISSION	<i>(the) chair is black</i>		59
DIRECT OBJECT PRO-DROP	<i>she doesn't like (it)</i>		14
FOCUS <i>itself</i>	<i>he is doing engineering in Delhi <u>itself</u></i>	24	5
FOCUS <i>only</i>	<i>I was there yesterday <u>only</u></i>	95	8
HABITUAL PROGRESSIVE	<i>always we <u>are giving</u> receipt</i>		2
STATIVE PROGRESSIVE	<i>he <u>is having</u> a television</i>		3
LACK OF INVERSION IN WH-QUESTIONS	<i>what <u>you are</u> doing?</i>		4
LACK OF AGREEMENT	<i>he <u>do</u> a lot of things</i>		23
LEFT DISLOCATION	<i><u>my father</u>, he works for a solar company</i>	300	19
MASS NOUNS AS COUNT NOUNS	<i>all the <u>musics</u> <u>are</u> very good</i>		13
NON-INITIAL EXISTENTIAL	<i>every year inflation <u>is there</u></i>	302	8
OBJECT FRONTING	<i><u>minimum one month</u> you have to wait</i>	186	14
PP FRONTING WITH REDUCTION	<i><u>(on the) right side</u> we can see a plate</i>		11
PREPOSITION OMISSION	<i>I went (to) another school</i>		17
INVERSION IN EMBEDDED CLAUSE	<i>I don't know what <u>are they</u> doing</i>		4
INVARIANT TAG ( <i>isn't it, no, na</i> )	<i>the children are outside, <u>isn't it?</u></i>	786	17
EXTRANEOUS ARTICLE	<i>she has <u>a</u> business experience</i>		25
GENERAL EXTENDER <i>and all</i>	<i>then she did her schooling <u>and all</u></i>		7
COPULA OMISSION	<i>my parents (are) from Gujarat</i>	71	
RESUMPTIVE OBJECT PRONOUN	<i>my old life I want to spend <u>it</u> in India</i>	24	
RESUMPTIVE SUBJECT PRONOUN	<i>my brother, <u>he</u> lives in California</i>	287	
TOPICALIZED NON-ARGUMENT CONSTITUENT	<i><u>in those years</u> I did not travel</i>	272	

Table 1: Features of Indian English used in our evaluations and their counts in the two datasets we study.

# Minimal pairs

**ARTICLE OMISSION:** *chair is black* → *the chair is black*

**FOCUS *only*:** *I was there yesterday only* → *I was there just yesterday.*

**NON-INITIAL EXISTENTIAL:** *every year inflation is there* → *every year there is inflation.*

- Standard method of presentation in linguistics
- Demszky et al. use as supervision (one positive, one negative example) for BERT fine tuning to identify that particular feature

<b>Dialect feature</b>	<b>DAMTL Multihead</b>	
ARTICLE OMISSION	0.581	0.658
DIRECT OBJECT PRO-DROP	0.493	0.563
EXTRANEIOUS ARTICLE	0.546	0.465
FOCUS <i>itself</i> *	1.000	0.949
FOCUS <i>only</i> *	0.998	0.775
HABITUAL PROGRESSIVE	0.439	0.718
INVARIANT TAG	0.984	0.901
INVERSION IN EMBEDDED CLAUSE	0.719	0.884
LACK OF AGREEMENT	0.543	0.674
LACK OF INVERSION IN WH-QUESTIONS	0.649	0.660
LEFT DISLOCATION	0.758	0.820
MASS NOUNS AS COUNT NOUNS	0.443	0.465
NON-INITIAL EXISTENTIAL*	0.897	0.885
OBJECT FRONTING	0.722	0.789
PREPOSITION OMISSION	0.500	0.648
PP FRONTING WITH REDUCTION	0.655	0.697
STATIVE PROGRESSIVE	0.645	0.789
GENERAL EXTENDER <i>and all</i>	0.994	0.991
<b>Macro Average</b>	0.698	0.741

Table 5: ROC-AUC results on the extended feature set, averaged across five random seeds. Because labeled





# Social impact of NLP research

- Real-world language technology implementation
- Uses of language technologies
- Social biases in NLP models