# Racial Disparities in Automated Speech Recognition
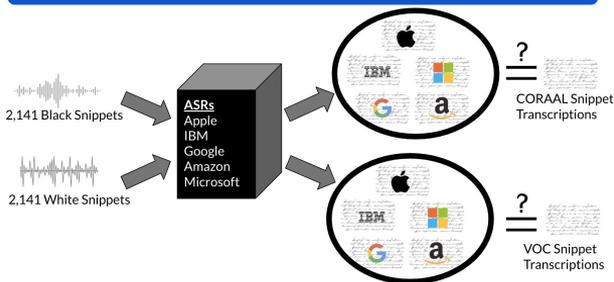
Allison Koenecke
{koenecke@stanford.edu}

## Motivation
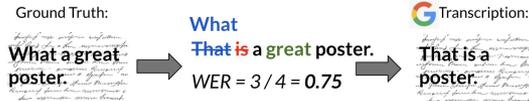


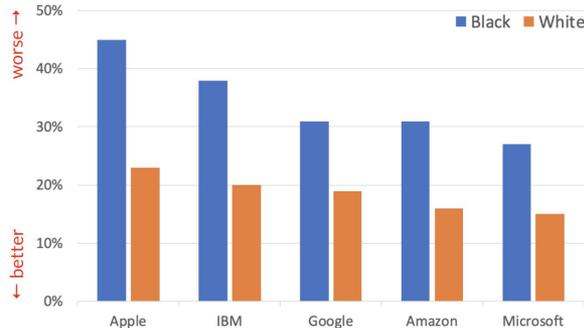## Methodology



2,141 Black Snippets

2,141 White Snippets

ASRs
Apple
IBM
Google
Amazon
Microsoft

? = CORAAL Snippet Transcriptions

? = VOC Snippet Transcriptions

## Metric

Word Error Rate = (**Substitutions** + **Deletions** + **Insertions**) / # Ground Truth Words

Ground Truth:
What a great poster.

**What** ~~That~~ **is** a **great** poster.

Transcription:
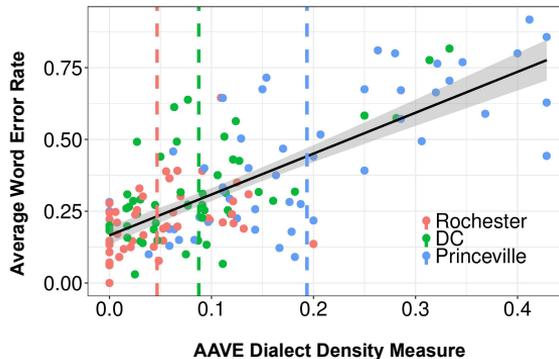That is a poster.

WER = 3 / 4 = **0.75**

## Results

- Speech-to-text transcription error rates are twice as high for Black speakers versus white speakers
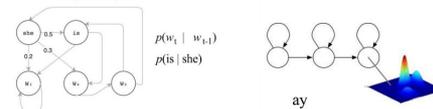


- High word error rates are correlated with high usage of African American Vernacular English linguistic features



## Models

- ASR = Language Model + Acoustic Model



- Controlling for for Black and white speakers uttering the same words, we find the acoustic model alone yields racial disparities
- Surprising result: language model is *not* the driver of disparity, despite GPT-2 analysis:

| | AAVE Perplexity | SE Perplexity |
|---|---|---|
| He**'s** a pastor. | 305 | 67 |
| We**'re** going to the ark. | 190 | 88 |
| We**'re** able to fight for the cause. | 54 | 51 |
| Where **are** they from? | 570 | 20 |
| Have you decided what you**'re** going to sing? | 106 | 25 |

## Call to Action

- Invest in resources to ensure inclusivity in ASR systems and institutions building them
- Collect more diverse training data: for AAVE and other non-standard varieties of English
- Regularly assess and publicly report ASR progress in fairness over time
- Study technical and regulatory progress made in other domains (computer vision)