

Problem Set #1



# Equitability in Voice Assistants

Aparna Ayyah & Matt Ho (Team 77)

# Methodology

## 1. PROBLEM CONTEXT

Current Voice Assistant limitations

## 2. DATASETS OVERVIEW

Indic TTS + CommonVoice

## 3. CUSTOM MODEL PERFORMANCE

RNN Model + Spectrograms

## 4. TEXT UNDERSTANDING

N-Gram Model Improvements

## 5. CONCLUSIONS

Summary of Improvements

# Inequity in Voice Assistants

## Dataset

Historically drawn from  
standard, “White” dialects



## Prediction Error

Word error rate up to 2x  
higher for accented  
speakers vs White  
speakers



## Positive Feedback Loop

Skews future models  
because of current error

# Key Underserved Factors



## Accented Speech

Those with accents may pronounce or phrase commands differently



## Speech Impairment

Voice Assistants have a long way to go with accessibility



## Children

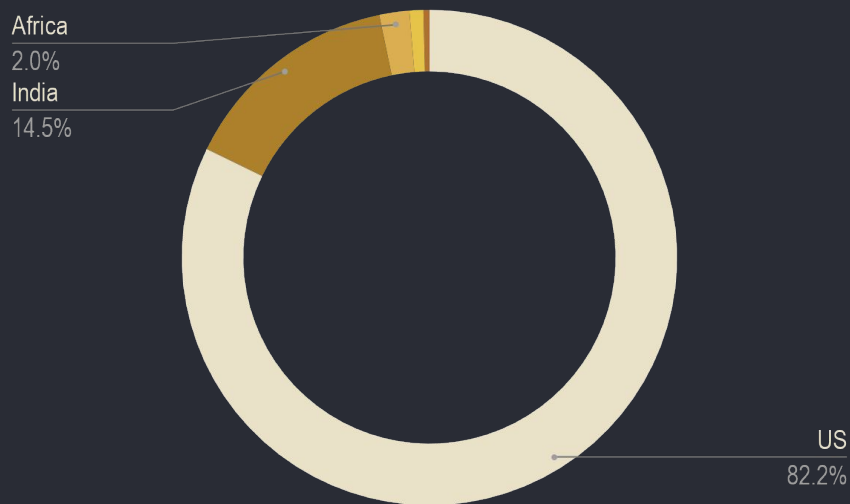
Children interact with and arrange words differently than adults

Scanlon, Dr. Patricia. "Voice Assistants Don't Work for Kids: The Problem with Speech Recognition in the Classroom."

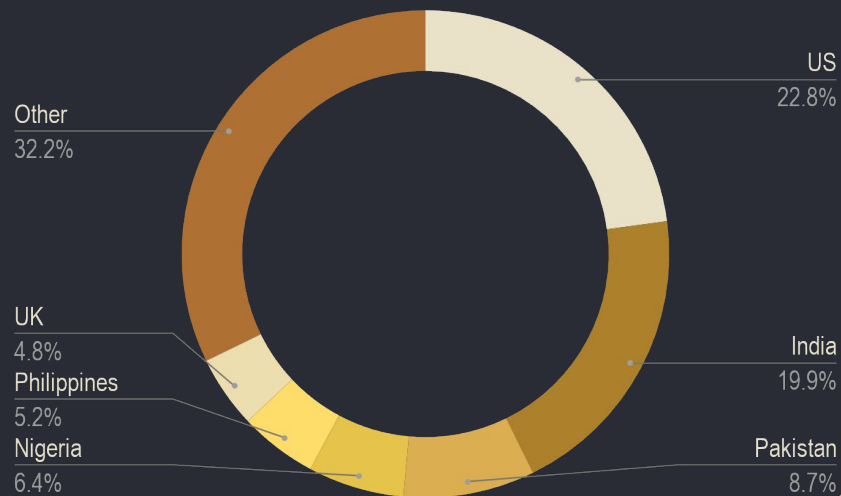
Masina, Fabio, et al. "Investigating the Accessibility of Voice Assistants With Impaired Users: Mixed Methods Study."

# Data Composition

Composition of CommonVoice Dataset

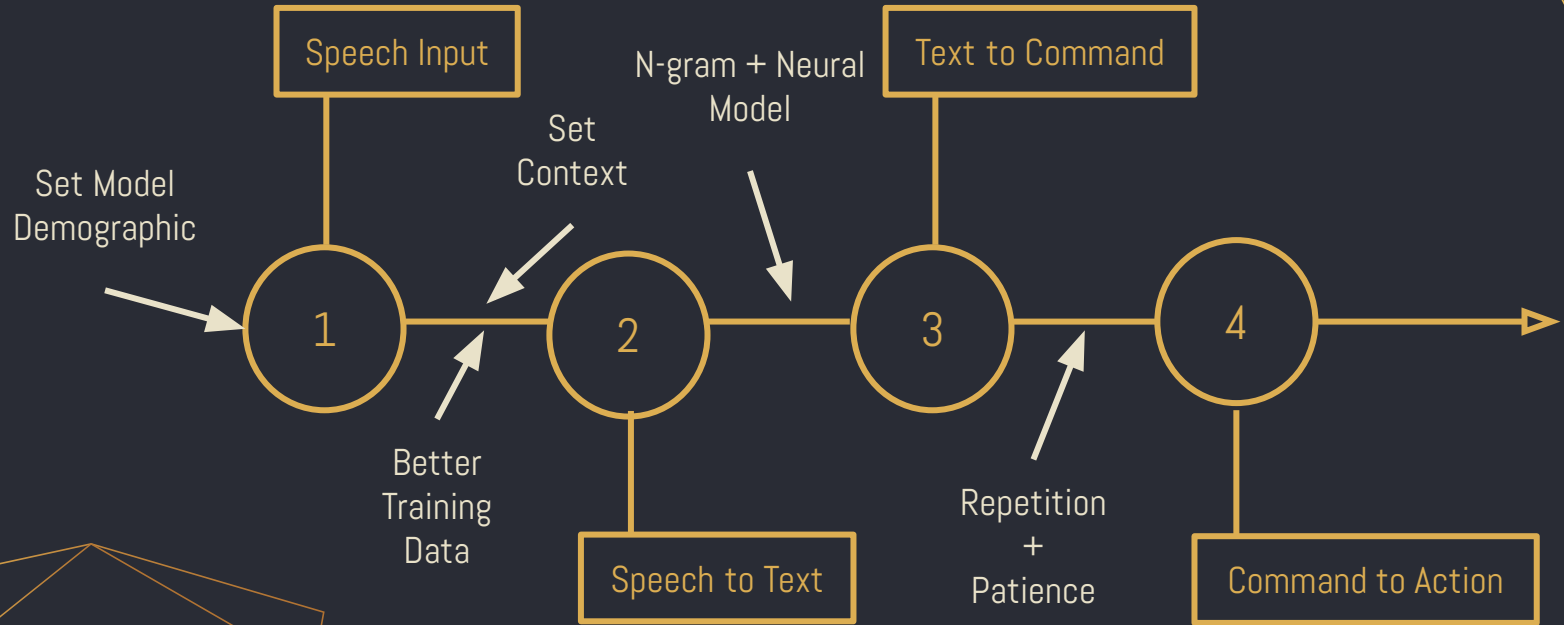


English Speakers Worldwide



**The Training Data is not representative of Users.**

# Voice Assistant Sequence



# Improvements in Speech to Text

# Datasets

## Indic TTS

Open Source Speech Dataset:

- Created by Indian Government
- Has extensive audio of Indian accents + phrasing
- Recorded by researchers
- Used for our **Custom Model**

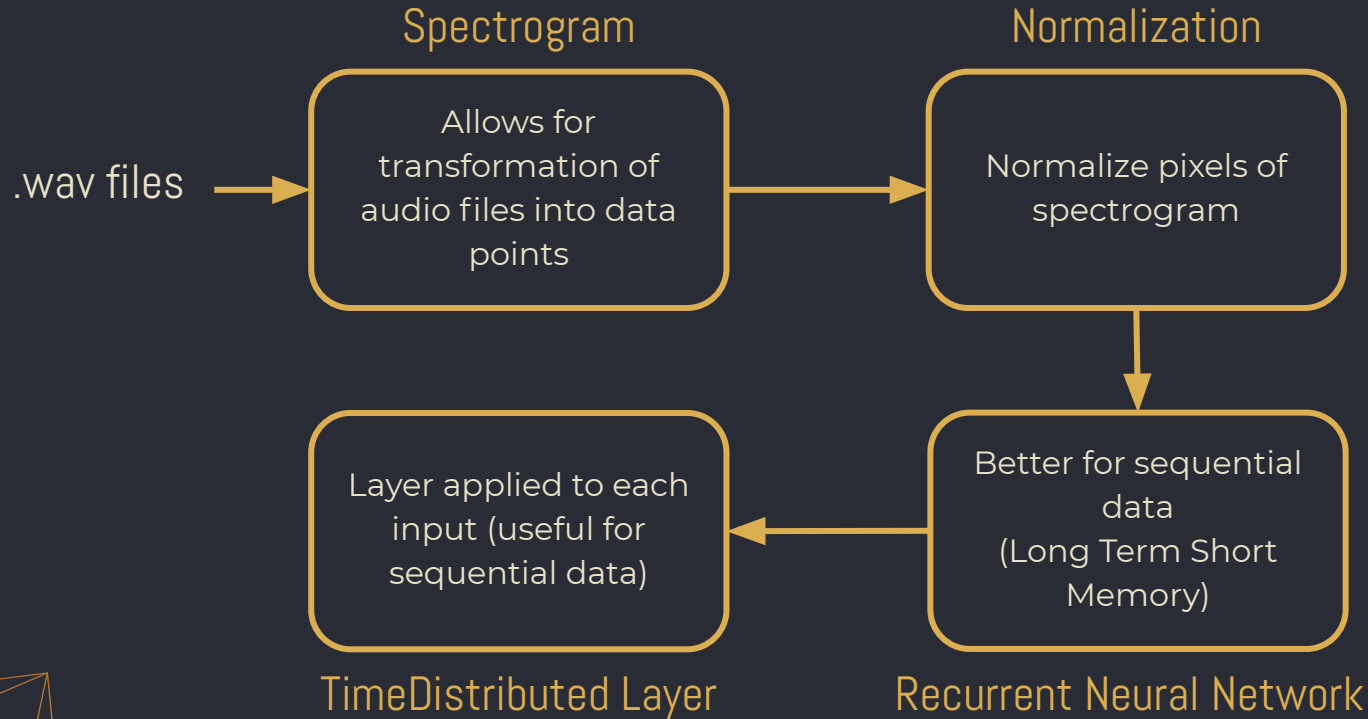
## CommonVoice

Open Source Speech Dataset:

- Created by Mozilla
- Similar to the (proprietary) datasets used by Amazon, Google, etc.
- User and Volunteer Generated
- Used by **DeepSpeech** (Voice Recognition Engine)



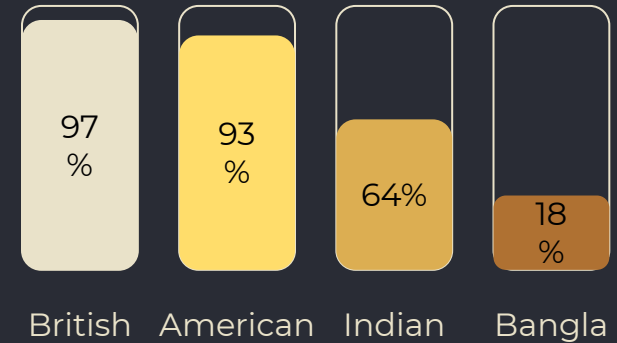
# Custom Model Components



# Comparison of Models

	Generic DeepSpeech	Trained Custom Model
Training Data	Standard	Indian Accents
Epochs	20	20
Word Error Rate	0.44	<b>0.16</b>
Word Accuracy	0.56	<b>0.84</b>

Prediction Confidence of DeepSpeech



# K-Fold Cross Validation



# Improvements in Text to Action

# N-Gram Model

## N-Gram Method

- Used by conventional models
- Calculate probability of the next word given the past  $n-1$  words
- $N = 4$  (for large scale, working models)

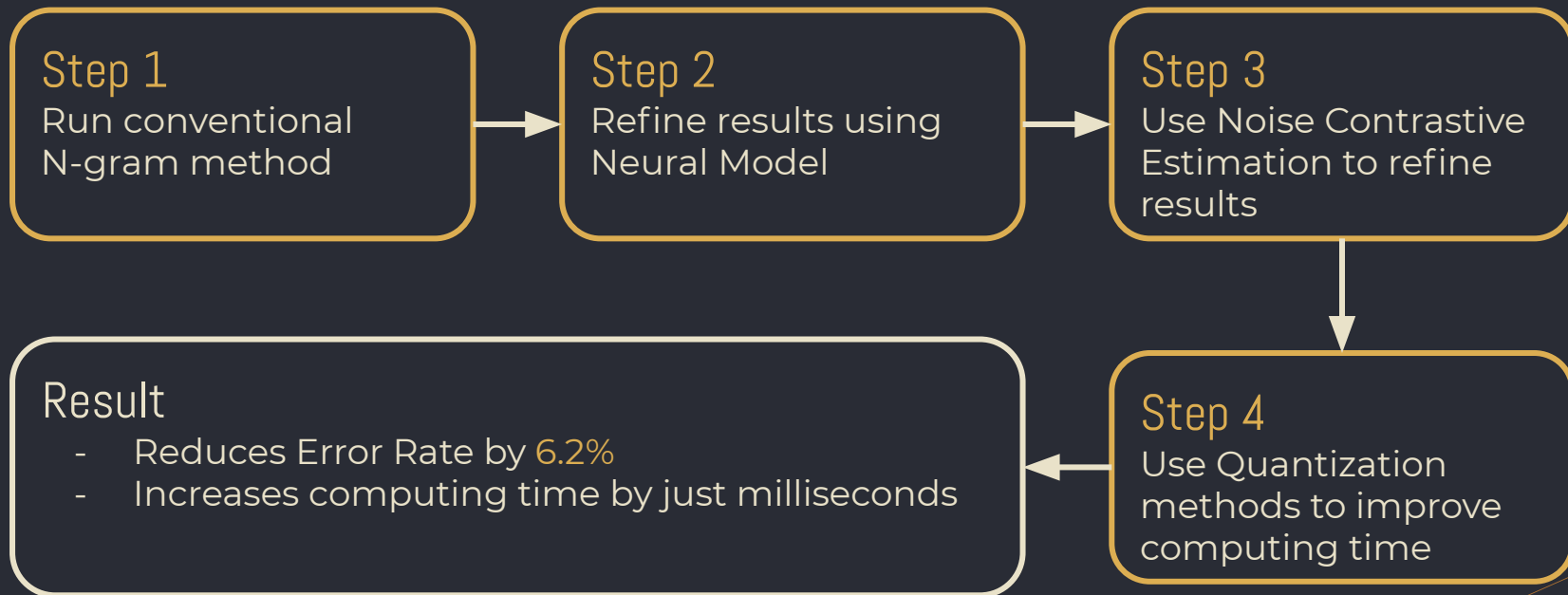
## Advantages

- Uses some context in prediction (N amount)

## Limitations

- Misses longer range dependencies

# Modified N-Gram Model





# Additional Improvements

## Setting Context

- Classifying by Intent:
  - Saying “Temperature” for temperature related commands
  - Adds to time taken, but increases accuracy

## Repetition + Learning

- Ask again if unsure:
  - “Sorry, could you repeat that?”
- Track frequent errors

# Conclusions

## Speech to Text

- Currently training datasets are not robust enough
  - Natural Variations (like accent) can and will occur
- Mirroring population composition with training data improves model performance
  - Users can input preferred model ("select dialect") to improve voice assistant performance

## Text to Action

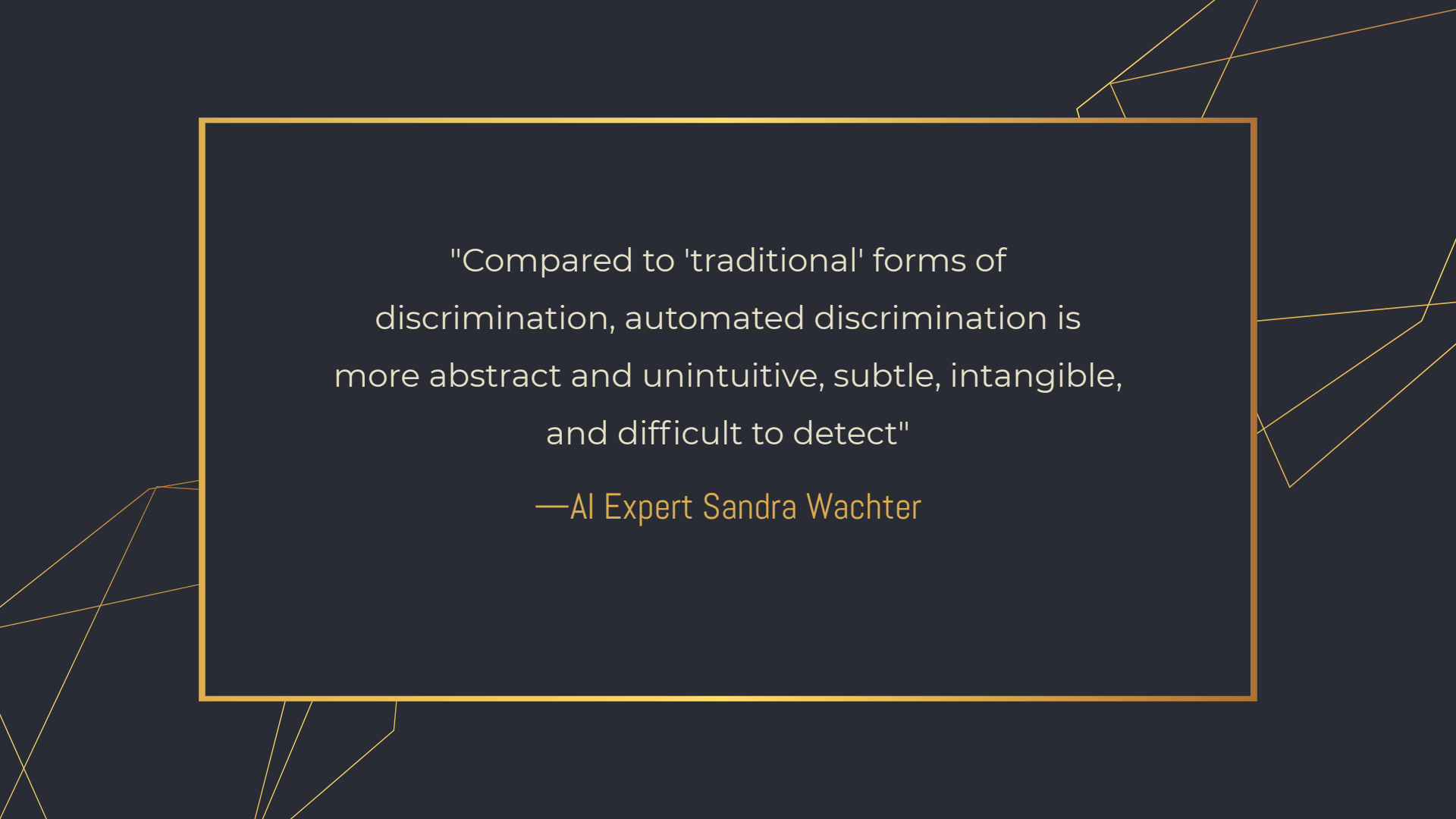
- Text to Action understanding can be improved
  - Understand long range context
- Allow User to set "context"
  - Helps reduce misunderstandings in word choice





# Thank You!



The background of the slide is a dark navy blue. It is decorated with thin, light gold lines that form various geometric shapes, including triangles and polygons, scattered across the frame. A prominent rectangular frame in the same light gold color encloses the central text.

"Compared to 'traditional' forms of  
discrimination, automated discrimination is  
more abstract and unintuitive, subtle, intangible,  
and difficult to detect"

—AI Expert Sandra Wachter