Final Project
# Topic Voice Synthesis: Generating Human Like Speech from Text
Okoronkwo, Alfred

CSCI E-89 Deep Learning, Fall 2025
**Harvard University Extension School**
Prof. Zoran B. Djordjević & prof. Rahul Joglekar
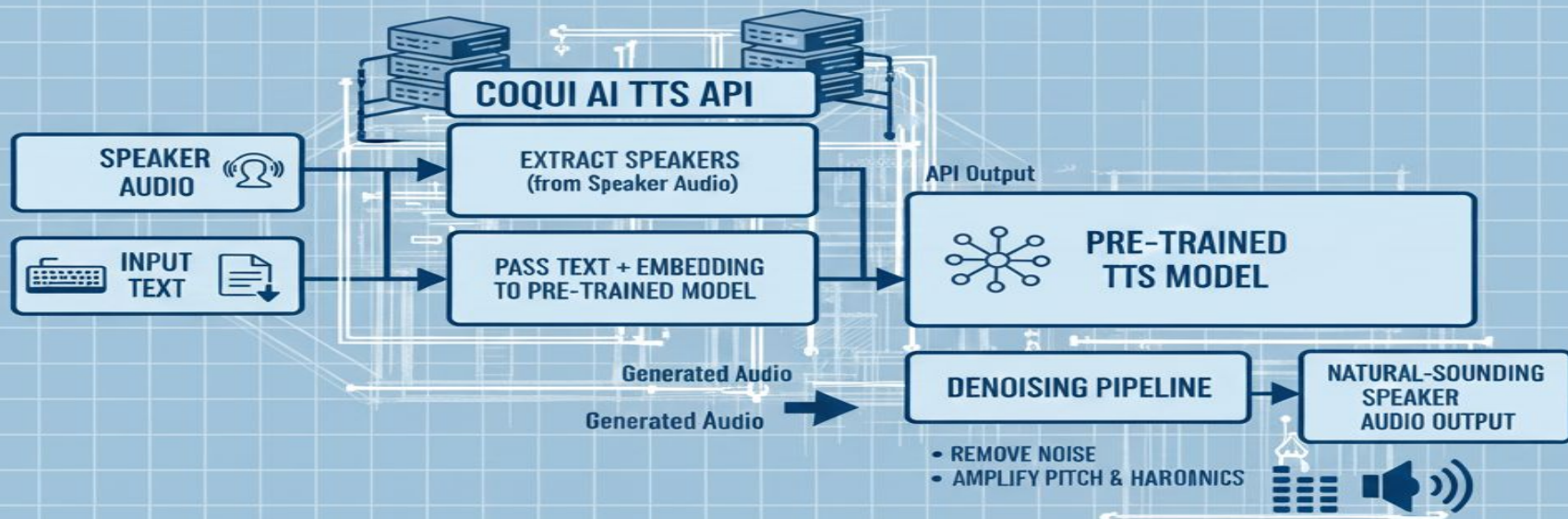
# Introduction

## Goal:

develop a system that can clone a target speaker's voice and generate speech from arbitrary text in the speaker's voice.
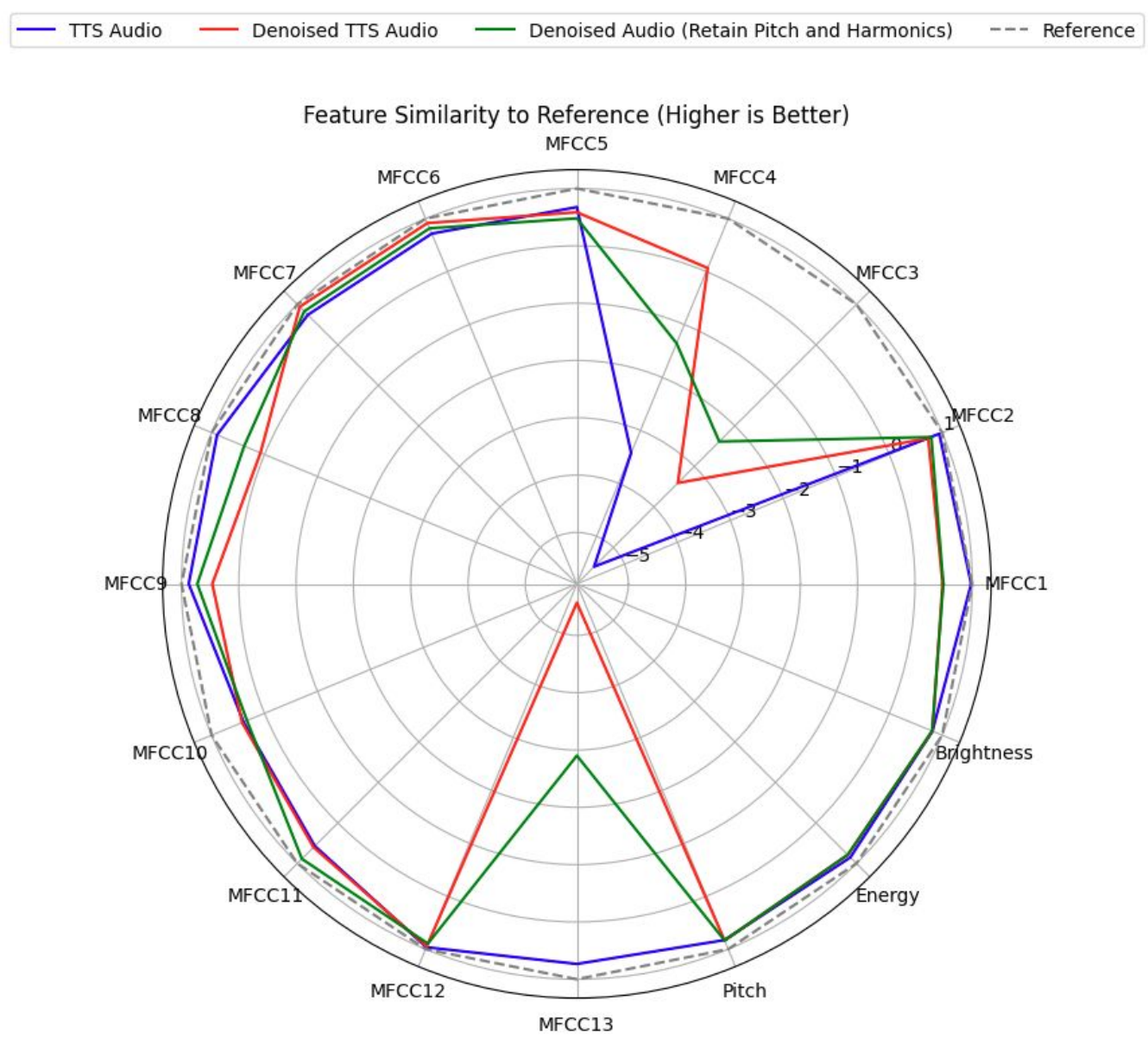
## Model:

the project uses the tts_models/multilingual/multi-dataset/your_tts pretrained model

# Zero-Shot Voice Cloning

# Audio Similarity Comparison



Feature Similarity to Reference (Higher is Better)

Legend: TTS Audio, Denoised TTS Audio, Denoised Audio (Retain Pitch and Harmonics), Reference

# Conclusion

By preserving the speaker's harmonics and pitch, we were able to improve the TTS-generate audio. Listening to the audio confirms that the resulting speech is much closer to the target speaker, sounding significantly better than both the original TTS-generated audio and the initially denoised TTS output.

# YouTube URLs, Last Page

- 2 minute (short): https://youtu.be/7qUkhFdp_Js
- 15 minutes (long):