# Challenges for Using Automatic Speech Recognition Systems to Detect Children's Reading Errors

Jaiden Magnan, Liliana Lusvardi, Olivia Schwartz, Alejandra Casillas, Walter Leite
Virtual Learning Lab: vll@coe.ufl.edu

## Introduction

- Underperformance in reading impacts emotional and academic areas, including science and math.
- Schools lack resources for one-on-one tutoring; reading apps can help but currently lack tailored feedback.
- No existing apps provide real-time, individualized feedback on reading errors.
- Current Automatic speech recognition (ASR) research focuses on typical development, with little focus on error patterns in struggling readers.
- This project aims to develop ASR-based models for detecting reading errors and creating a reading app intervention with just-in-time feedback.

## Methods

- Two pre-trained ASR models are used: Whisper (transformer-based encoder-decoder) and XLSR-53 (WAV2VEC architecture with CTC loss).
- Models are tested for accuracy in detecting reading errors.
- Audio recordings of first to third graders' reading practice.
- ASR - generated transcripts are being used to train and evaluate models and app feedback system.

## Whisper

- Whisper is Open AI's ASR model that was trained on 680,000 hours of audio from a variety of different languages
- Whisper offers five models: base, tiny, small, medium, and large. Tiny is the most efficient but least accurate, while large is the least efficient but most accurate.
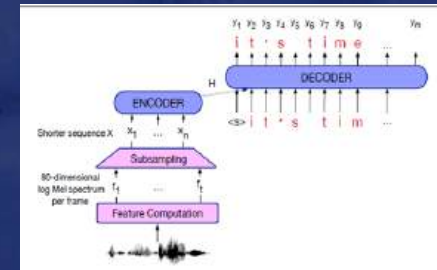
## Wav2Vec

- Wav2Vec is a model developed by Facebook AI that converts speech to text.
- The XSL-53 is designed to handle multilingual scenarios, likely helping with decoding children's speech.

## Limitations

- Children's recordings had background noise, which impacted the accuracy of the ASR models
- The whisper model had issues generating the correct punctuation, which caused the accuracy to drop.



Above are audio transcriptions from Whisper. The left side displays the results from the recording with no errors, and the right side shows it with errors. The green highlights show when either the model accurately identified a mispronunciation or generated nonsensical text because it could not parse input.



ufdatastudio.com/posts/2023-08-20-ASR-Blog/

## Results

| Model | % Accuracy of transcription of text without errors | % Accuracy of transcription of text with errors |
|---|---|---|
| Whisper | 90.7% | 62.7% |
| Wav2Vec | 48 % | 5 % |

References