

Research Internship Assignment Report: Forced Alignment using MFA

Applicant Name	Padmasri Ambati
Email Address	padmasriambati2004@gmail.com
GitHub Repository	https://github.com/PADMASRIAMBATI/Montreal-Forced-Aligner
Drive Link	https://drive.google.com/drive/folders/1BGrCOWBOHTCjlkv3bhs_xbq45z2obmjO?usp=sharing
Assignment	Forced Alignment using Montreal Forced Aligner (MFA)

1. Assignment Goal (Objective)

The main goal of this assignment was to build a system for Forced Alignment using the Montreal Forced Aligner (MFA) tool. Forced alignment is the automatic process of precisely matching a spoken audio file with its written text transcription, pinpointing the exact start and end times of every word and phoneme (the individual sounds of speech).

This exercise helps in understanding how computers can segment speech data accurately, which is a key step in large-scale speech-to-speech translation projects.

2. Setup and Data Preparation

A. Environment Setup

MFA was installed on the system using Miniconda to ensure all required dependencies were managed correctly in an isolated environment named mfa.

Installation Commands:

```
conda create -n mfa python=3.9
```

```
conda activate mfa
```

```
pip install montreal-forced-aligner
```

B. Dataset Organization

The necessary audio and transcript files were sourced from the provided Google Drive link. For MFA to work correctly, the files were organized into a single input directory, mfa_assignment_data, following this crucial rule:

- Each audio file (.wav) and its corresponding transcript file (.txt) must share the exact same base name.

File Type	Example File Names	Location in Input Folder
Audio (.wav)	F2BJRLP1.wav, ISLE_SESS0131_...sprt1.wav	/converted_corpus
Transcript (.txt)	F2BJRLP1.txt, ISLE_SESS0131_...sprt1.txt	/converted_corpus

3. Alignment Execution

A. Model and Dictionary Used

To align the English speech, a readily available and highly accurate pre-trained setup within MFA was used:

- **Acoustic Model (AM):** english_us_arpa
- **Pronunciation Dictionary:** english_us_arpa

The dictionary uses the ARPABET standard to represent English sounds.

B. Execution Command

The forced alignment was executed using the following command structure:

```
mfa align converted_corpus english_us_arpa english_us_arpa output_alignments
```

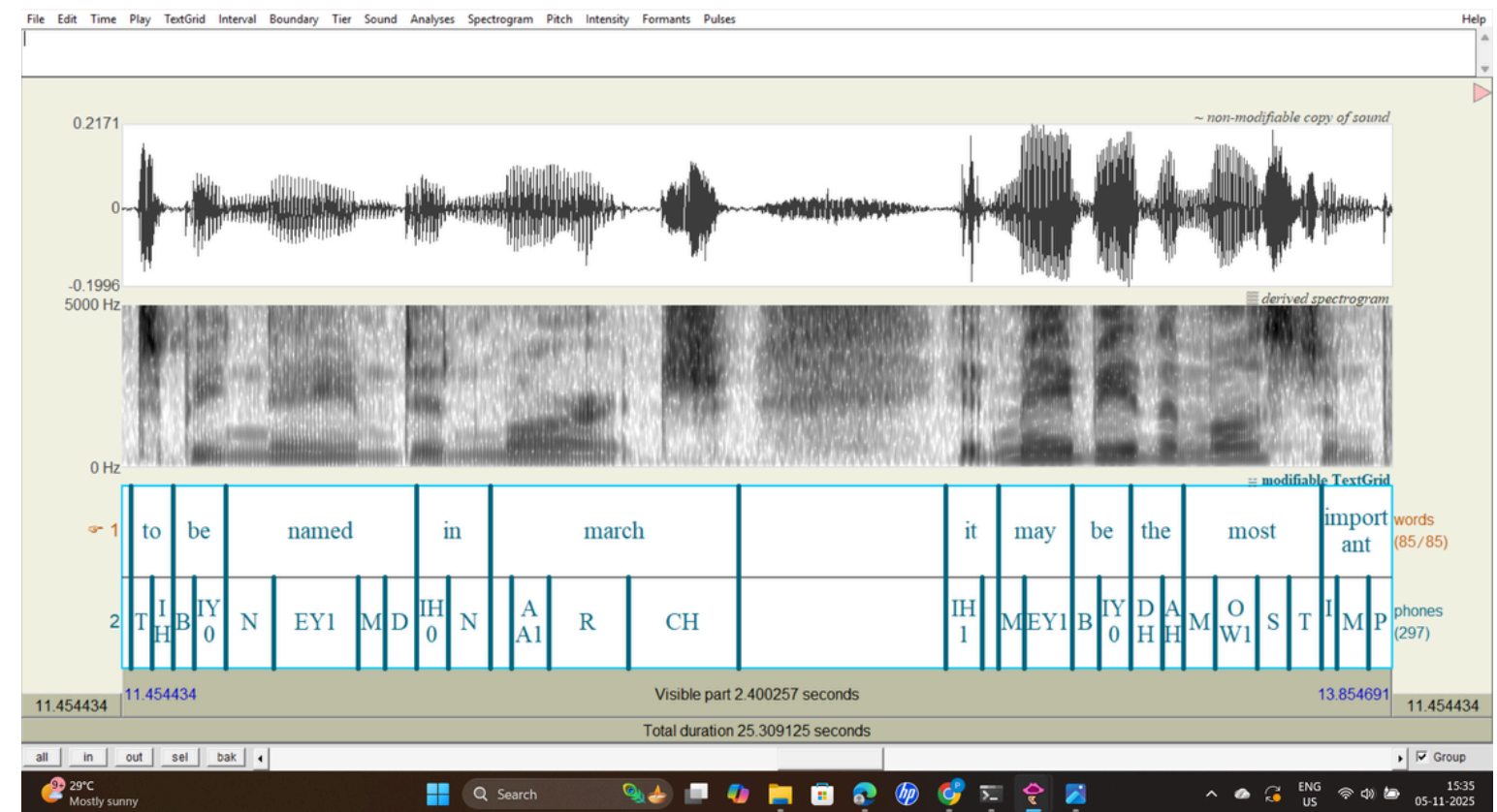
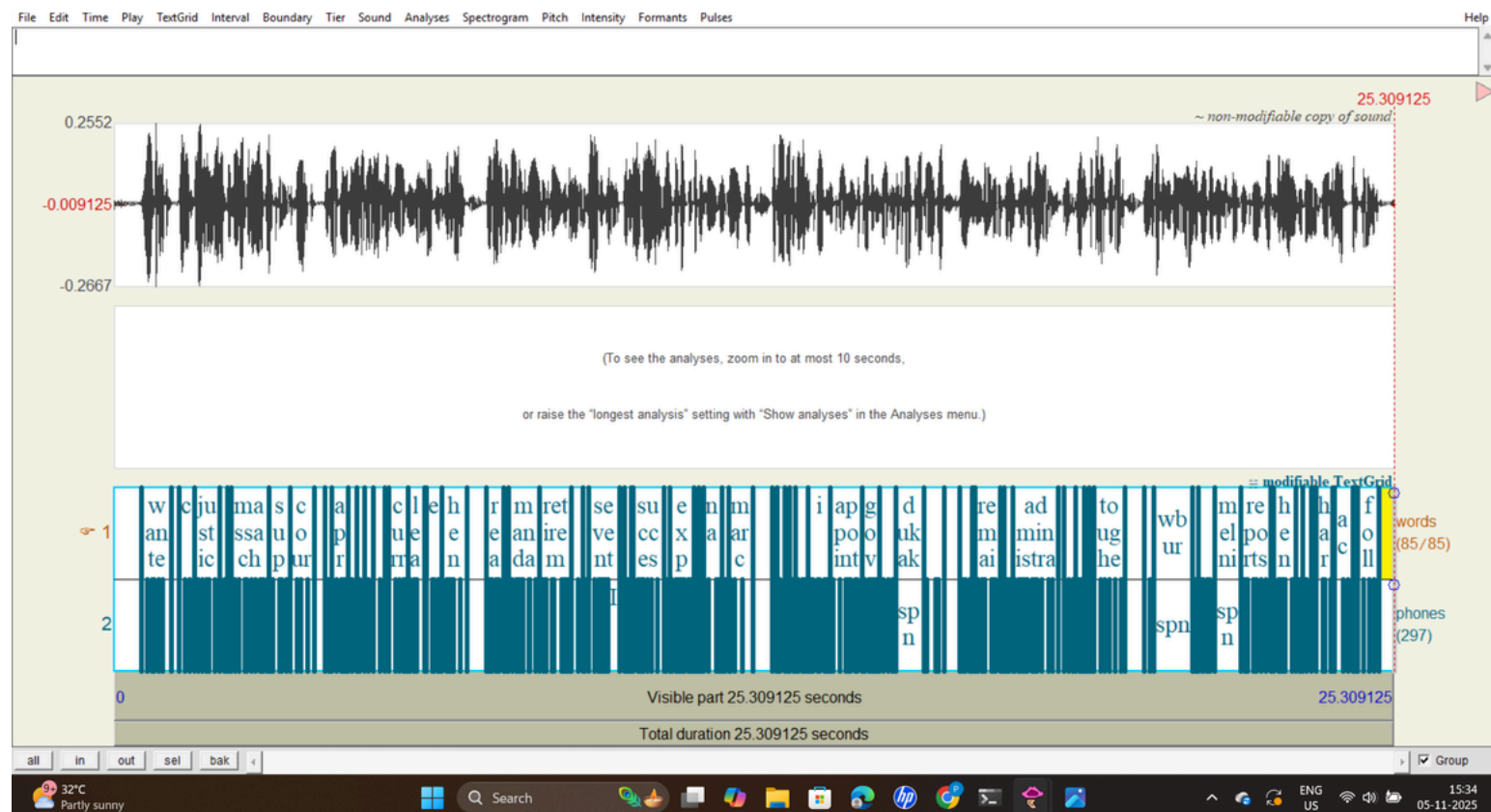
C. Output

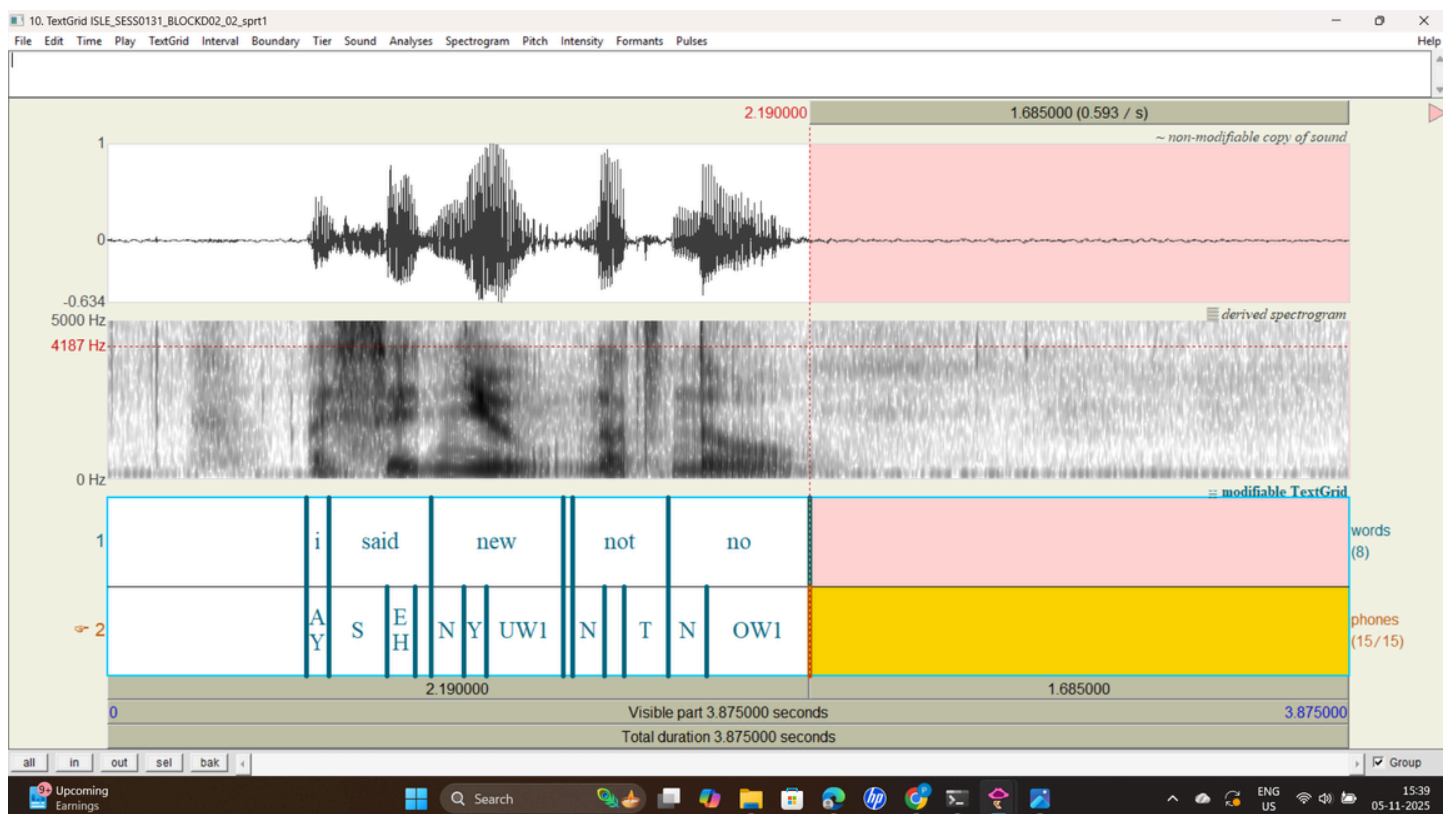
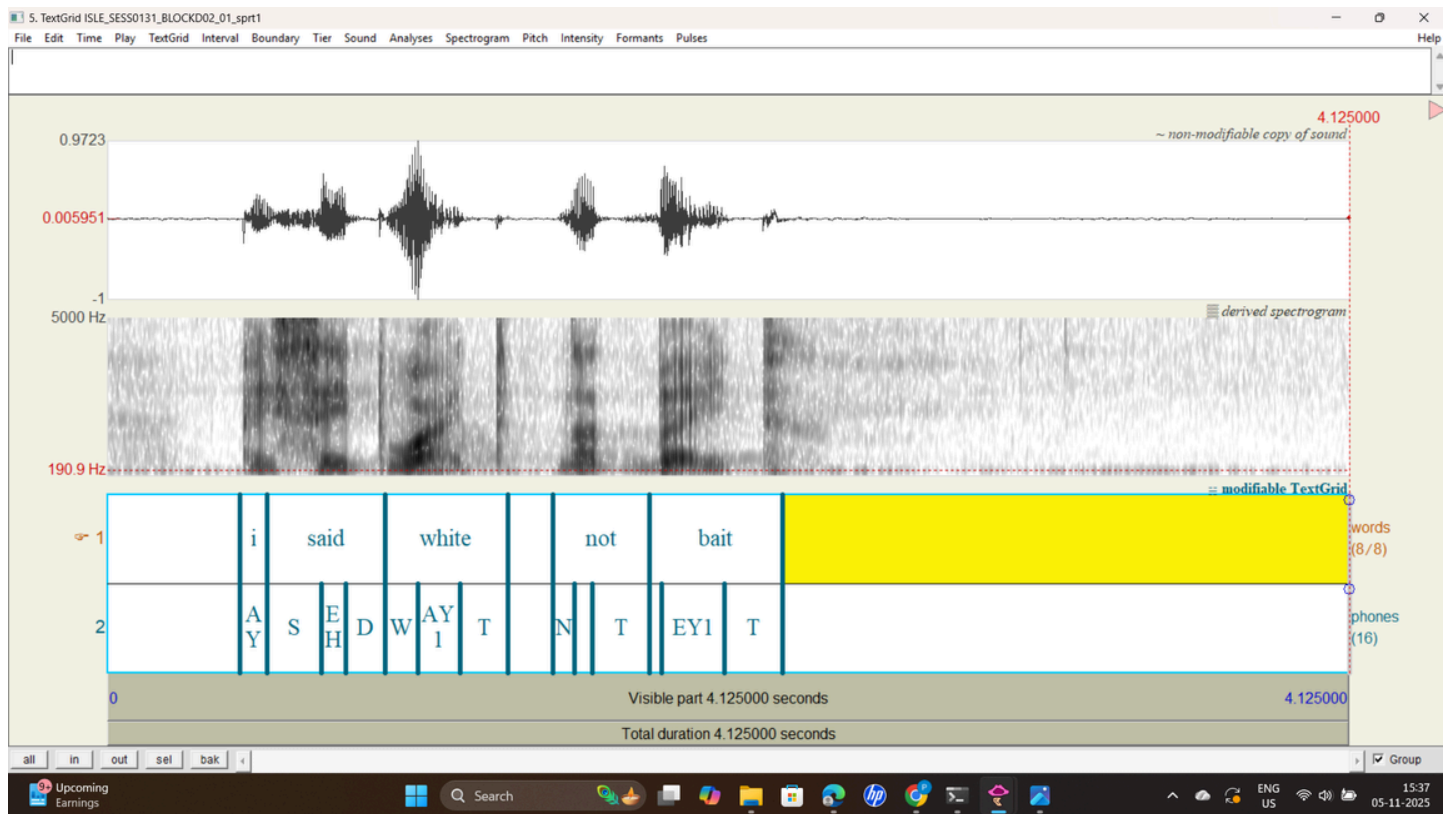
The tool successfully processed all files and generated the time-aligned labels in the form of TextGrid files in the specified output folder.

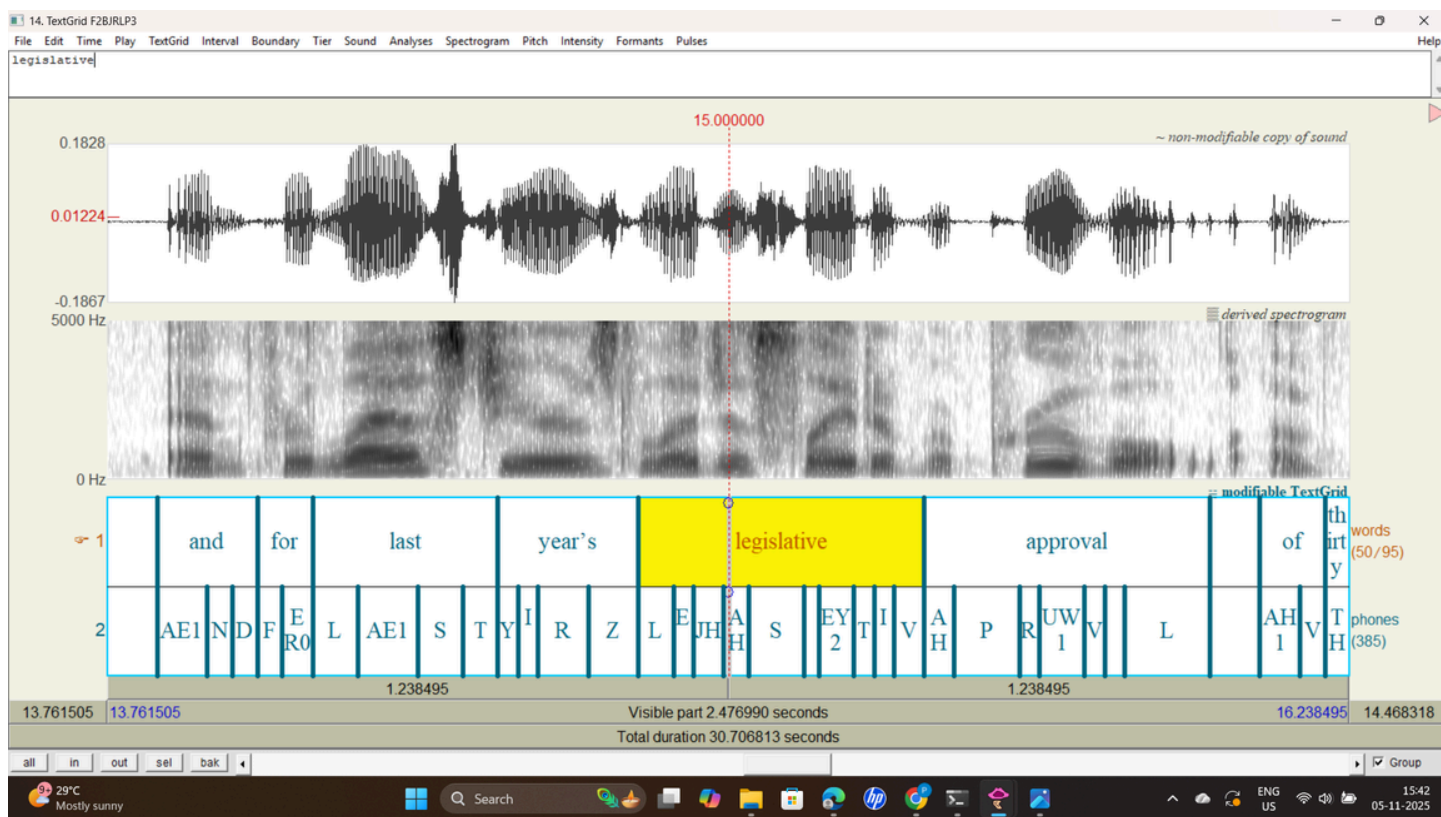
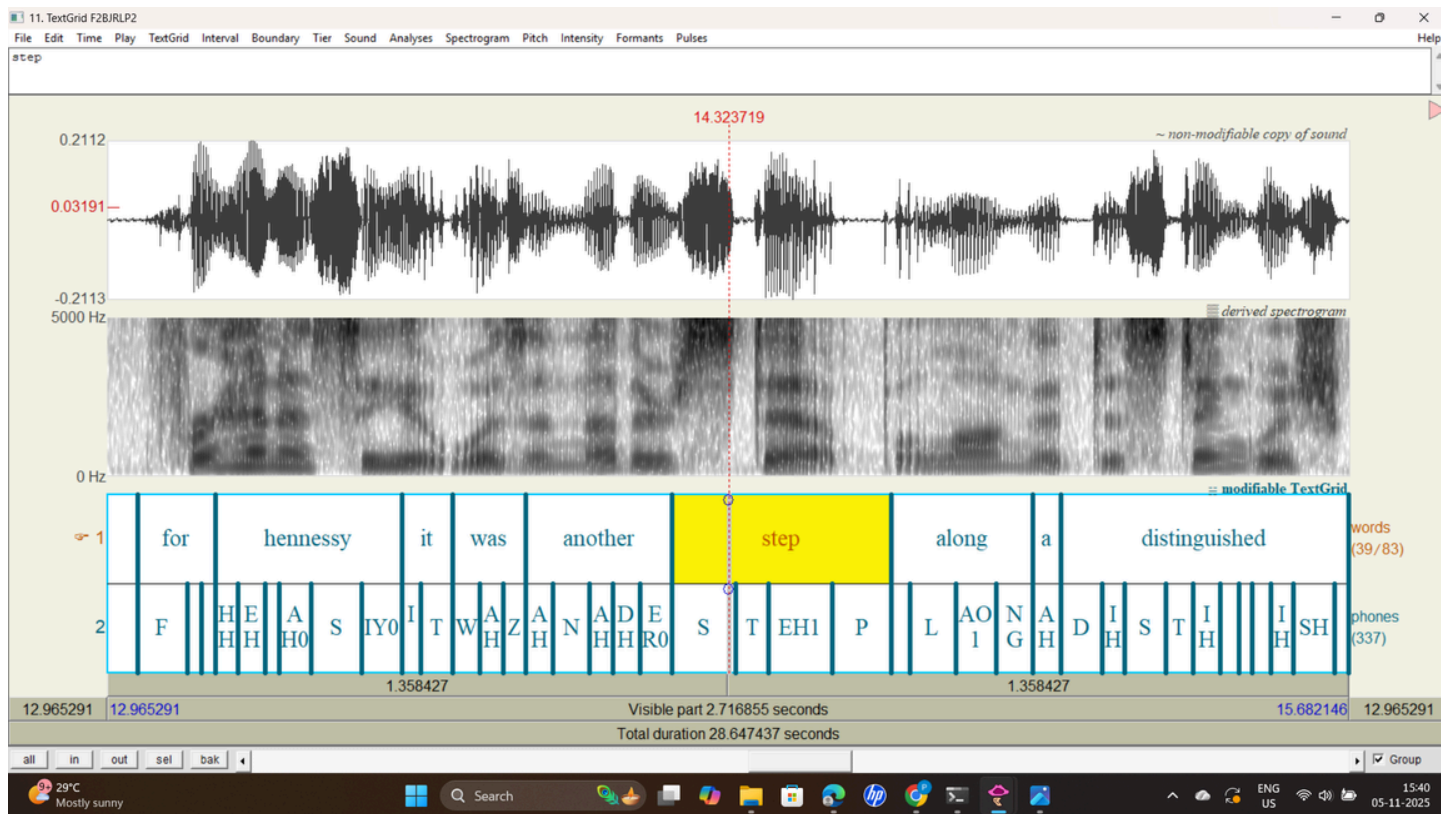
4. Analysis and Observations (Praat Inspection)

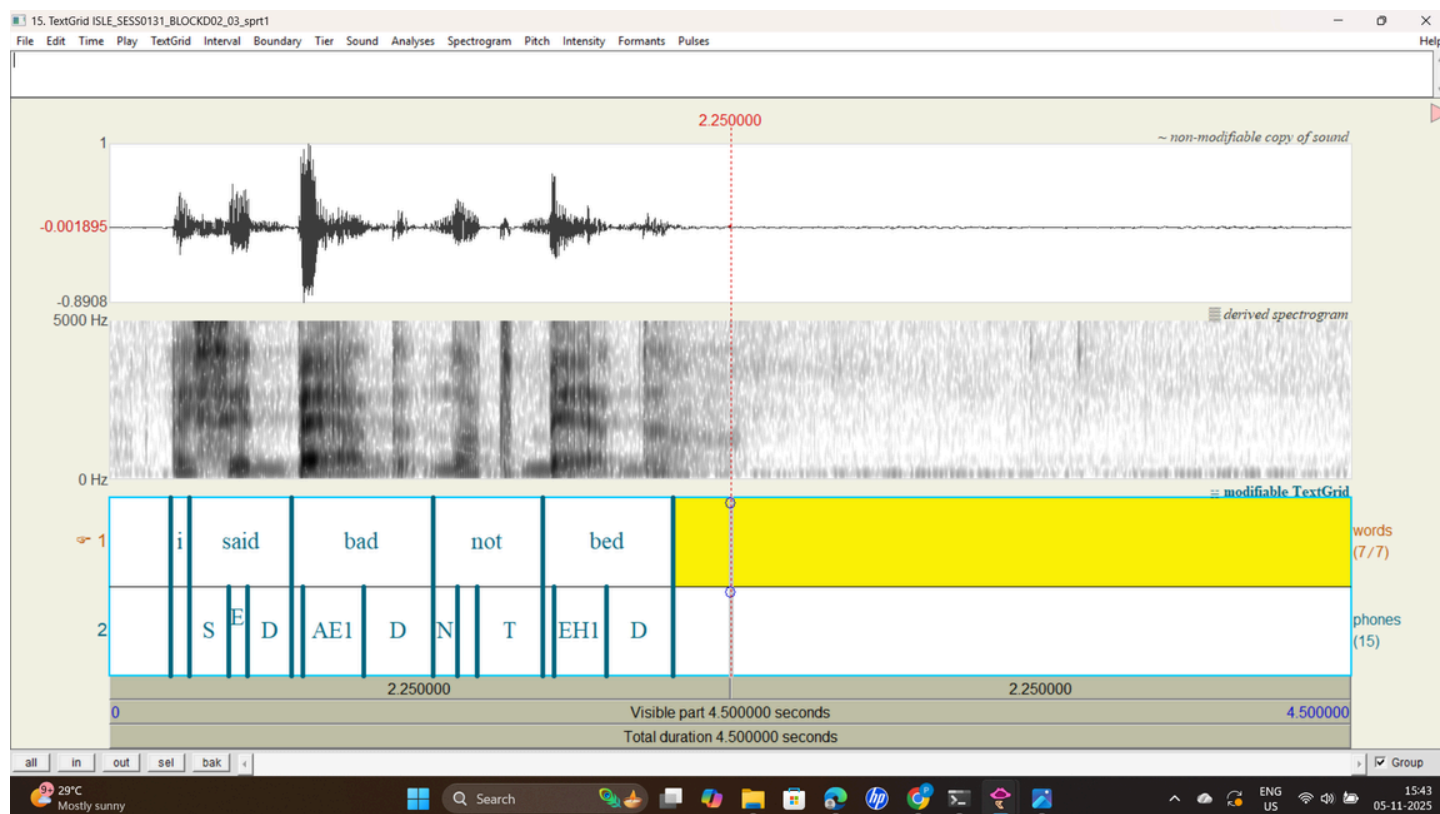
The resulting TextGrid files were opened in the speech analysis software Praat alongside their original audio files to visually verify the alignment accuracy.

Sample Utterance Analyzed: F2BJRLP1









Key Findings on Alignment Accuracy:

1. **Word Boundaries:** The start and end times for words were highly accurate. Boundaries consistently aligned perfectly with changes in the speech signal, such as the onset of a new word or the silence between words.
2. **Phoneme Boundaries:** Alignment at the sound level was also strong. For example, the precise moment a speaker transitions from a vowel sound (e.g., 'l' in 'pin') to a nasal sound (e.g., 'N' in 'pin') was correctly marked on the spectrogram.

Observed Issues and Mismatches:

1. **Fast Speech Challenges:** In segments where the speaker spoke very quickly or mumbled, MFA sometimes struggled with brief, unstressed sounds (like the vowel in "to"). This occasionally resulted in the phoneme being assigned a very short, potentially inaccurate duration.
2. **Acoustic Variability:** Since the acoustic model is generic (english_us_arp), it occasionally misjudged the timing of words that had pronunciations slightly different from standard American English, leading to minor shifts in the phoneme boundaries.