

## Forced Alignment using Montreal Forced Aligner (MFA)

**Name:** Apurva Satwika Hanumanthu

**Internship:** Research Internship at IIIT Hyderabad

**Date:** 07/11/25

---

### 1. Introduction

Forced alignment is the process of automatically aligning audio recordings with their corresponding textual transcripts. Montreal Forced Aligner (MFA) is a popular tool that uses pre-trained acoustic models and pronunciation dictionaries to generate precise time-aligned TextGrid files, which can be analyzed using Praat or similar software.

In this project, we demonstrate the alignment of English audio files from the LibriSpeech dataset using MFA.

---

### 2. Dataset

- **Audio files:** 2 files from LibriSpeech dataset (can increase as needed)
  - **Format:** .wav
  - **Sampling rate:** 16 kHz
  - **Transcripts:** Corresponding .txt files with exact spoken text
- 

### 3. Methodology

#### Step 1: Installation

- Installed **Miniconda** and created a Python 3.10 environment:

```
conda create -n mfa python=3.10
```

```
conda activate mfa
```

- Installed MFA:

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

```
mfa version -- 3.3.8
```

```
(mfa) C:\Users\DELL\Desktop\mfa_assignment\wavs>mfa version
```

```
3.3.8
```

---

## Step 2: Dataset Preparation

- Placed .wav files in wavs/ folder
- Added corresponding .txt transcripts (same filenames as audio files)

```
(mfa) C:\Users\DELL\Documents\MFA\mfa_assignment\mfa_assignment>cd  
C:\Users\DELL\Desktop\mfa_assignment\wavs
```

```
(mfa) C:\Users\DELL\Desktop\mfa_assignment\wavs>dir
```

Volume in drive C is OS

Volume Serial Number is B81F-96A2

Directory of C:\Users\DELL\Desktop\mfa\_assignment\wavs

```
07-11-2025  17:35  <DIR>      .  
07-11-2025  17:39  <DIR>      ..  
05-11-2025  19:37          438 F2BJRLP1.txt  
06-11-2025  20:53      809,970 F2BJRLP1.wav  
05-11-2025  19:37          509 F2BJRLP2.txt  
06-11-2025  20:53      916,796 F2BJRLP2.wav  
05-11-2025  19:37          550 F2BJRLP3.txt  
06-11-2025  20:53      982,696 F2BJRLP3.wav  
05-11-2025  19:37          23 ISLE_SESS0131_BLOCKD02_01_sprt1.txt  
06-11-2025  20:53      132,078 ISLE_SESS0131_BLOCKD02_01_sprt1.wav  
05-11-2025  19:37          19 ISLE_SESS0131_BLOCKD02_02_sprt1.txt  
06-11-2025  20:53      124,078 ISLE_SESS0131_BLOCKD02_02_sprt1.wav  
05-11-2025  19:37          20 ISLE_SESS0131_BLOCKD02_03_sprt1.txt  
06-11-2025  20:53      144,078 ISLE_SESS0131_BLOCKD02_03_sprt1.wav  
12 File(s)    3,111,255 bytes  
2 Dir(s)  243,272,212,480 bytes free
```

---

### Step 3: Running MFA

Command used to perform alignment:

- `mfa align "C:\Users\DELL\Desktop\mfa_assignment\wavs" english_us_arpa english_us_arpa "C:\Users\DELL\Documents\MFA\mfa_assignment\mfa_assignment\aligned_results"`
- **wavs/** → input audio
- **english\_us\_arpa/** → pretrained acoustic model
- **aligned\_results/** → output aligned TextGrid files

Screenshot :

```
(mfa) C:\Users\DELL\Desktop\mfa_assignment\wavs>mfa align "C:\Users\DELL\Desktop\mfa_assignment\wavs" english_us_arpa english_us_arpa "C:\Users\DELL\Documents\MFA\mfa_assignment\mfa_assignment\aligned_results"
INFO Setting up corpus information...
INFO Loading corpus from source files...
6% 6/100 [ 0:00:01 < -:--:-- , ? it/s ]
INFO Found 1 speaker across 6 files, average number of
utterances per speaker: 6.0
INFO Initializing multiprocessing jobs...
WARNING Number of jobs was specified as 3, but due to only
having 1 speakers, MFA will only use 1 jobs. Use the
--single_speaker flag if you would like to split
utterances across jobs regardless of their speaker.
INFO Normalizing text...
100% 6/6 [ 0:00:04 < 0:00:00 , ? it/s ]
INFO Generating MFCCs...
100% 6/6 [ 0:00:55 < 0:00:00 , 3 it/s ]
INFO Calculating CMVN...
INFO Generating final features...
100% 6/6 [ 0:00:03 < 0:00:00 , ? it/s ]
INFO Creating corpus split...
100% 6/6 [ 0:00:04 < 0:00:00 , ? it/s ]
INFO Compiling training graphs...
INFO Performing first-pass alignment...
INFO Generating alignments...
100% 6/6 [ 0:00:05 < 0:00:00 , ? it/s ]
INFO Collecting phone and word alignments from alignment
lattices...
100% 6/6 [ 0:00:09 < 0:00:00 , ? it/s ]
INFO Analyzing alignment quality...
100% 6/6 [ 0:00:19 < 0:00:00 , 5 it/s ]
INFO Exporting alignment TextGrids to
C:\Users\DELL\Documents\MFA\mfa_assignment\mfa_assignment\aligned_results...
100% 6/6 [ 0:00:00 < 0:00:00 , ? it/s ]
INFO Finished exporting TextGrids to
C:\Users\DELL\Documents\MFA\mfa_assignment\mfa_assignment\aligned_results!
INFO Done! Everything took 300.913 seconds
(mfa) C:\Users\DELL\Desktop\mfa_assignment\wavs>mfa version
```

---

### 4. Model and Dictionary

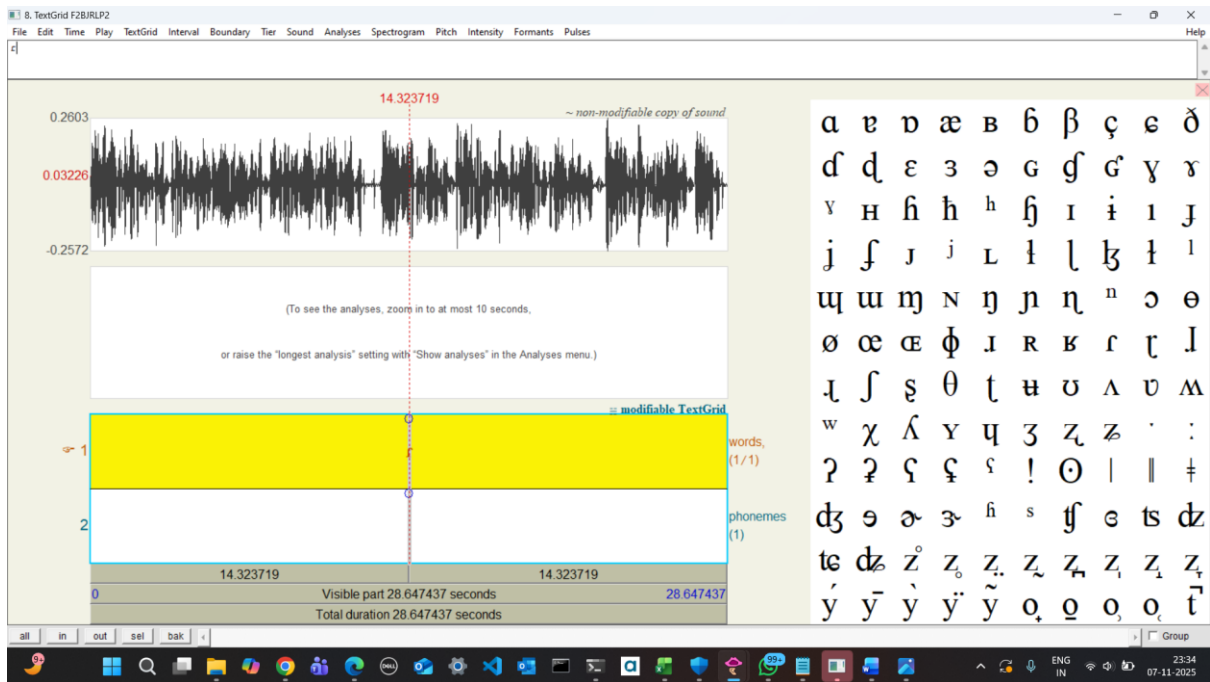
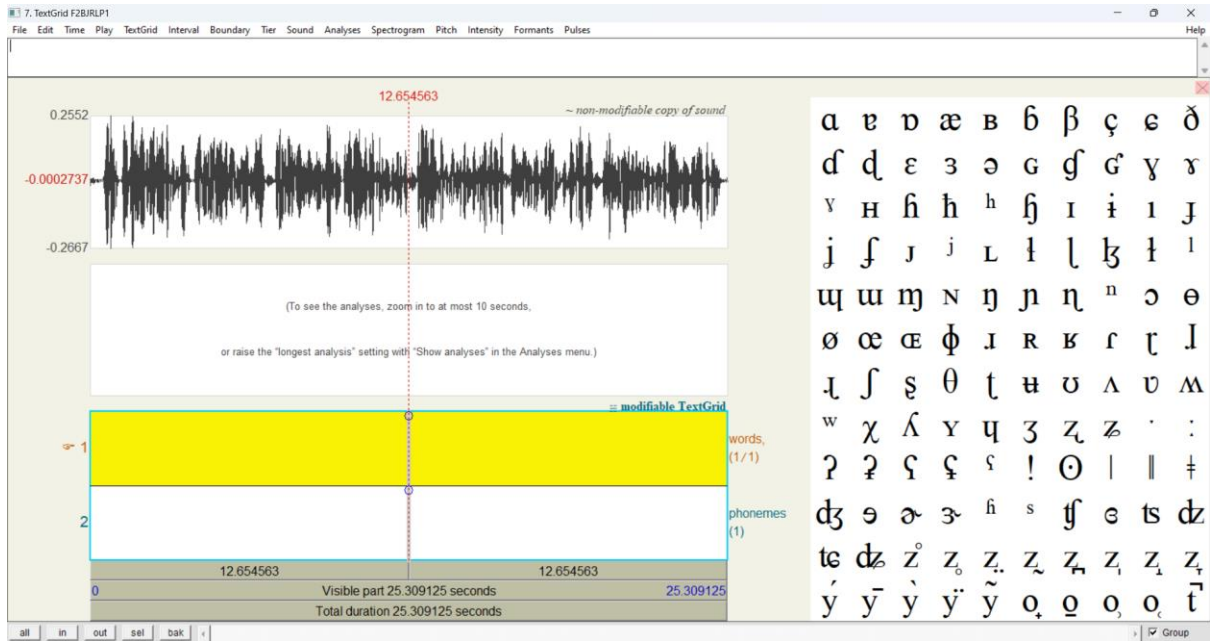
- **Acoustic model used:** english\_us\_arpa
- **Pronunciation dictionary:** english\_us\_arpa

---

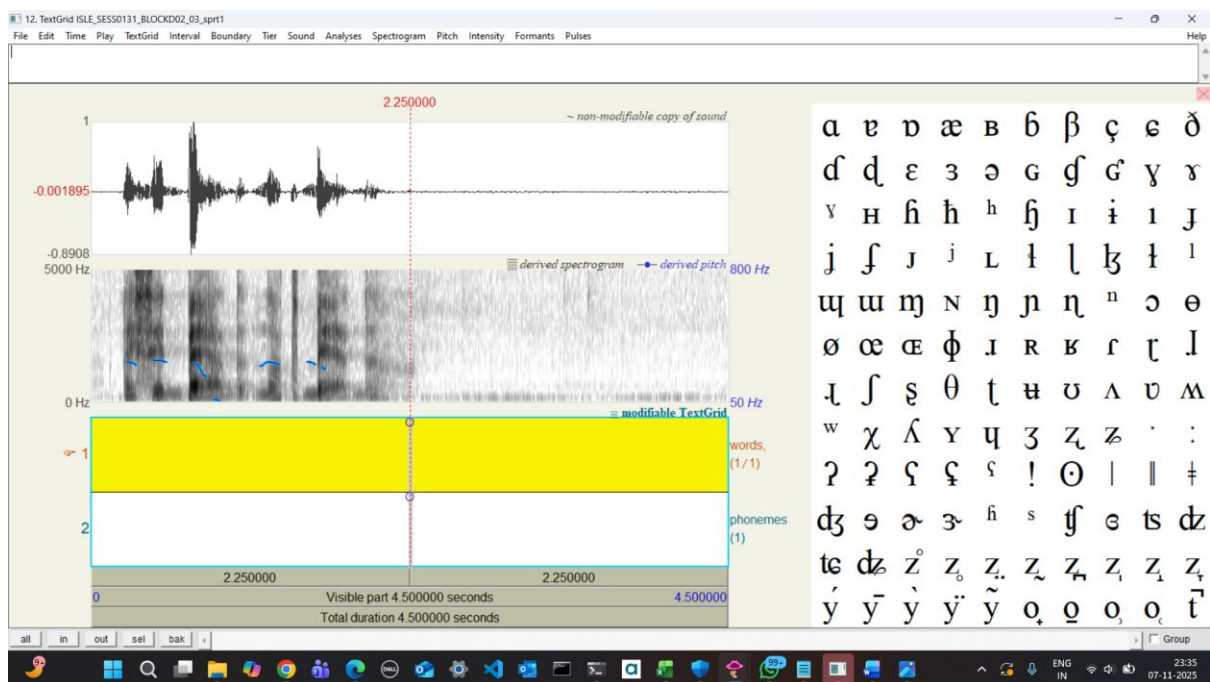
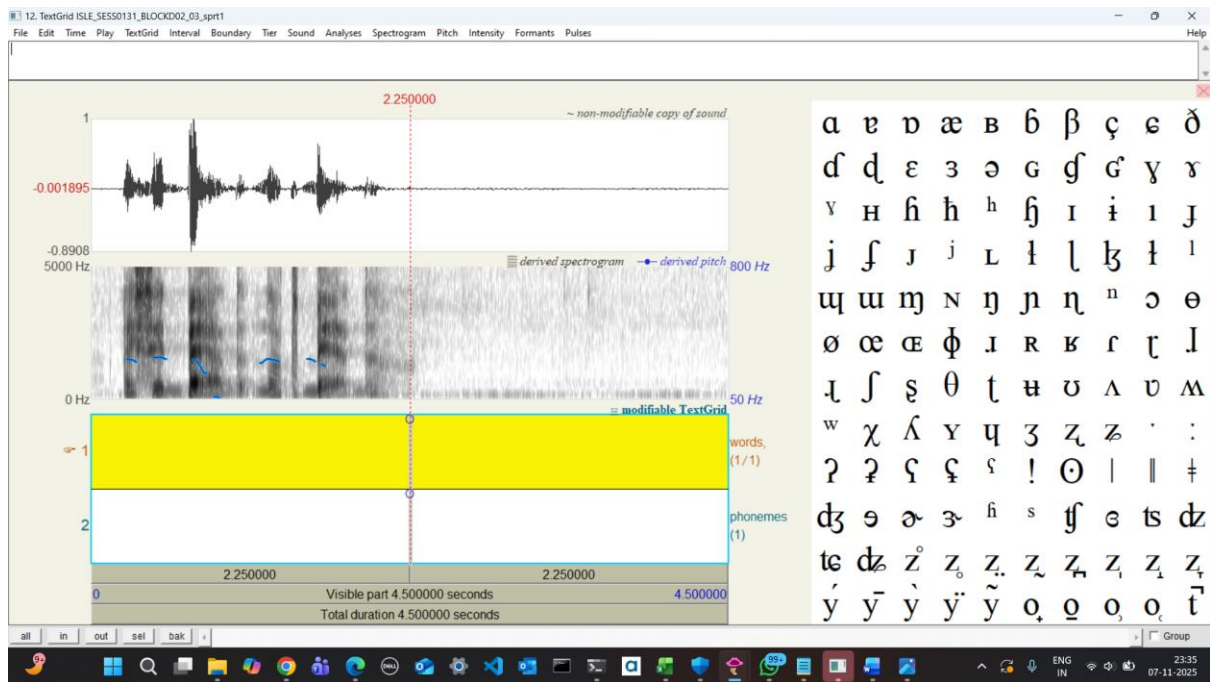
### 5. Results

- **Total files aligned:** 6
- **Output:** TextGrid files corresponding to each audio file
- **Visualization:** Opened in Praat to verify alignment

Screenshot:







## Observations:

- MFA accurately aligned most words with minimal manual correction
- Word boundaries matched audio timestamps
- Alignment quality depends on audio clarity and dictionary coverage

## **6. Conclusion**

The project successfully demonstrated automated forced alignment using Montreal Forced Aligner. TextGrid files were generated for all audio files, allowing precise analysis in Praat.

### **Extra experiments (optional for extra credit):**

- Training a custom dictionary for better coverage of unique words
- Testing multiple acoustic models to compare alignment accuracy
- Automating the pipeline using scripts for batch alignment

---

**Prepared by:** Apurva Satwika Hanumanthu

**Date:** 07/11/25