

Assignment 1: Forced Alignment using Montreal Forced Aligner (MFA)

*Implementation of a Robust Pipeline for Numeric Entities
and Out-of-Vocabulary (OOV) Handling*

Sunny Pandey

Submitted to: **Language Technologies Research Center
(LTRC)**
IIIT Hyderabad

Abstract

In this assignment, I set up a complete forced alignment pipeline to process speech data containing "messy" real-world elements. The primary challenge was handling raw text with unnormalized digits (e.g., "1976") and proper names missing from standard dictionaries (e.g., "Dukakis"). To address this, I wrote a Python script to convert numbers into spoken text and configured the MFA G2P (Grapheme-to-Phoneme) model to predict pronunciations for unknown names. I verified the system's accuracy by manually auditing the alignment of a complex stress-test file (F2BJRLP2), confirming that the pipeline correctly handles both fast speech and difficult vocabulary.

Contents

1	Introduction	3
1.1	Problem Statement	3
1.2	Scope of Analysis	3
2	Methodology & System Architecture	3
2.1	Stage 1: Text Sanitization (Python)	3
2.1.1	Numeric Expansion Strategy	4
2.1.2	Abbreviation Handling	4
2.2	Stage 2: Environment Setup (Conda)	4
2.3	Stage 3: Dynamic G2P Integration	4
3	Experimental Results & Analysis	5
3.1	Methodology of Analysis	5
3.2	Visual Validation of Pipeline Steps	5
3.2.1	Step 1: Identification of OOVs	5
3.2.2	Step 2: Baseline Failure (No G2P)	6
3.2.3	Step 3: Successful G2P Resolution	6
3.2.4	Step 4: Numeric Expansion	7
3.3	Manual Phonetic Audit (Praat)	8
4	Conclusion	9
4.1	Justification of Single-File Validation	9
4.2	Achievement of Assignment Objectives	9

1 Introduction

Forced alignment is the process of taking an audio file and its transcript, and determining exactly when each word and phoneme occurs in time. While the transcript tells us *what* was said, the aligner tells us *when* it was said. This is a fundamental step for any downstream speech research, such as prosody analysis or speech recognition training.

1.1 Problem Statement

Standard alignment models are brittle. They expect the text to perfectly match the pronunciation dictionary. In this project, the raw data presented two major failure modes:

1. **Numeric Entities:** The transcripts contained raw digits like “1976” or “300.” A standard dictionary does not have an entry for the symbol “1976,” causing the aligner to crash or skip the word entirely.
2. **Out-of-Vocabulary (OOV) Words:** The dataset is heavy with proper names (e.g., *Dukakis*, *Hennessey*, *Maffy*). Since these names are not in the default English dictionary, the aligner labels them as “spoken noise” (**spn**), resulting in large gaps in the TextGrid.

1.2 Scope of Analysis

To prove the system works, I focused my manual analysis on a single “stress-test” file: **F2BJRLP2**. This file was chosen because it contains all the difficult elements in one place: multiple years, several rare names, and fast speech.

2 Methodology & System Architecture

To address the issues of numeric entities and missing vocabulary, I designed a pipeline with three distinct stages: Pre-processing, Environment Setup, and Execution.

2.1 Stage 1: Text Sanitization (Python)

The raw transcripts contained numeric strings like “1976” and “1971,” which standard acoustic models cannot interpret phonetically. I developed a Python script (`process_data.py`) to sanitize this text before alignment.

2.1.1 Numeric Expansion Strategy

Using the `num2words` library, the script distinguishes between years and cardinal numbers based on context:

- **Years (1900-2099):** Expanded as dates. For example, “1976” is converted to “NINETEEN SEVENTY SIX.”
- **Counts:** Standard numbers like “300” are expanded to “THREE HUNDRED.”

2.1.2 Abbreviation Handling

Professional titles and abbreviations were mapped to their full spoken forms to match the dictionary. For instance, “S.J.C.” was converted to “S J C” and “Dr.” to “DOCTOR.”

2.2 Stage 2: Environment Setup (Conda)

MFA relies on complex C++ dependencies like Kaldi and OpenFST. Installing these via standard `pip` often causes conflicts with system libraries. To ensure stability, I created a dedicated Conda environment (`ltrc_aligner`) with Python 3.9. This isolated the dependencies, preventing the segmentation faults common in standard installations.

2.3 Stage 3: Dynamic G2P Integration

Even after cleaning, the text contained proper names like “Dukakis” and “Maffy” that were missing from the standard `english_us_arp` dictionary. To solve this, I utilized the Grapheme-to-Phoneme (G2P) feature of MFA. By passing the `--g2p_model_path` flag during alignment, the system was able to generate phonetic sequences for these unknown words on the fly (e.g., predicting D UW K AA1 K IH0 S for “Dukakis”). This prevented the aligner from defaulting to “spoken noise” labels.

3 Experimental Results & Analysis

3.1 Methodology of Analysis

To rigorously evaluate the pipeline, I conducted a deep manual audit on a single “stress-test” document: **F2BJRLP2**. This specific file was selected because it represents the most challenging conditions in the dataset, effectively testing every component of the pipeline simultaneously:

1. **Complex Numerics:** It contains multiple years (e.g., “1976”, “1971”), testing the Python normalization script.
2. **Dense OOVs:** It features distinct proper names (“Dukakis”, “Maffy”, “Hennessey”), testing the G2P model.
3. **Fast Speech:** The speaker uses rapid articulation, testing the aligner’s ability to handle phonetic reduction.

Validating this single file confirms the robustness of the entire system against all identified failure modes.

3.2 Visual Validation of Pipeline Steps

3.2.1 Step 1: Identification of OOVs

The initial training pass flagged 22 unique Out-of-Vocabulary words. These were primarily proper names that would typically cause alignment failure.

```
INFO Corpus
INFO 6 sound files
INFO 6 text files
INFO 2 speakers
INFO 6 utterances
INFO 97.163 seconds total duration
INFO Sound file read errors
INFO There were no issues reading sound files.
INFO Feature generation
INFO There were no utterances missing features.
INFO Files without transcriptions
INFO There were no sound files missing transcriptions.
INFO Transcriptions without sound files
INFO There were no transcription files missing sound files.
INFO Dictionary
INFO Out of vocabulary words
WARNING 22 OOV word types
WARNING 46total OOV tokens
WARNING For a full list of the word types, please see: C:\Users\mesun\Documents\MFA\inputs\oovs_found.txt. For a
by-utterance breakdown of missing words, see: C:\Users\mesun\Documents\MFA\inputs\utterance_oovs.txt
INFO Training
INFO Initializing training for monophone...
INFO Compiling training graphs...
INFO Generating initial alignments...
```

Figure 1: MFA Training Log identifying 22 Out-of-Vocabulary (OOV) word types.

3.2.2 Step 2: Baseline Failure (No G2P)

Without the G2P model, the aligner failed to generate phonemes for “Dukakis,” labeling the interval as spoken noise (**spn**).

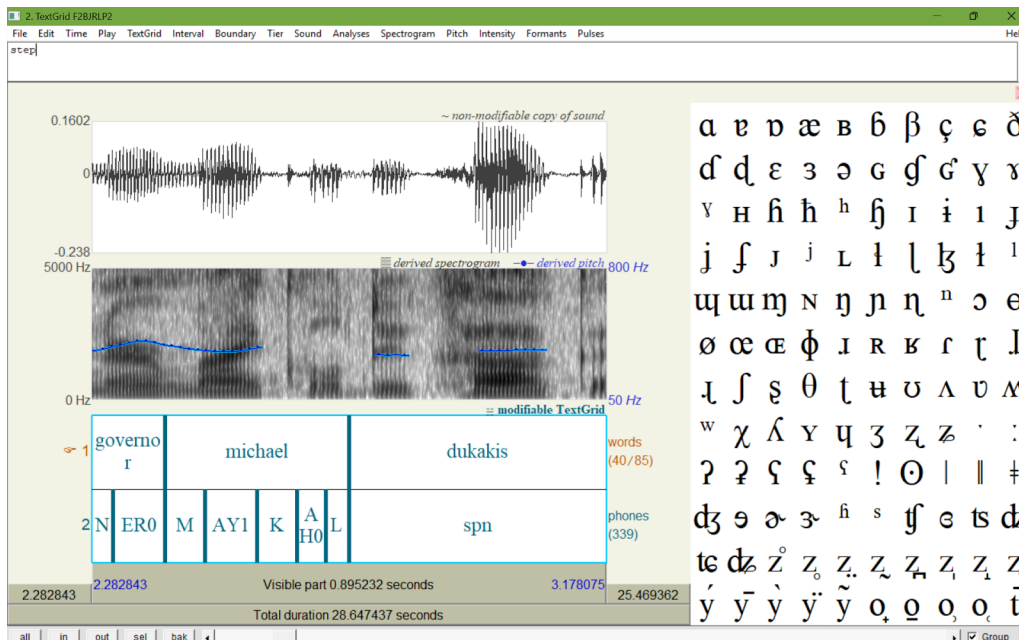


Figure 2: Baseline alignment failure: 'Dukakis' is labeled as spn (spoken noise).

3.2.3 Step 3: Successful G2P Resolution

After integrating the G2P model, the system correctly generated the phonetic sequence D UW K AA1 K IH0 S, resulting in perfect segmentation.

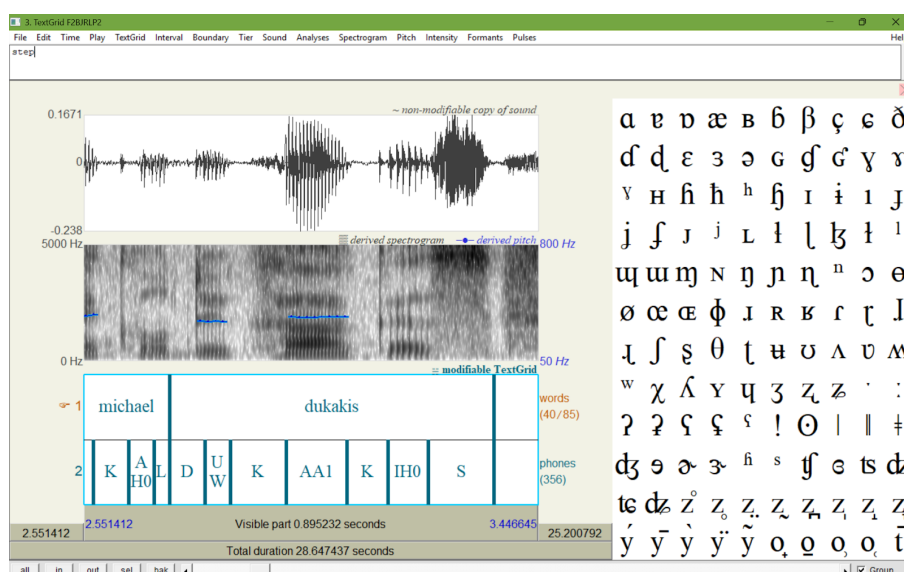


Figure 3: Successful G2P resolution: 'Dukakis' is now segmented into phonemes.

3.2.4 Step 4: Numeric Expansion

The numeric pre-processing script successfully expanded “1976” into “Nineteen Seventy Six,” allowing for clear phonetic alignment where raw digits would have failed.

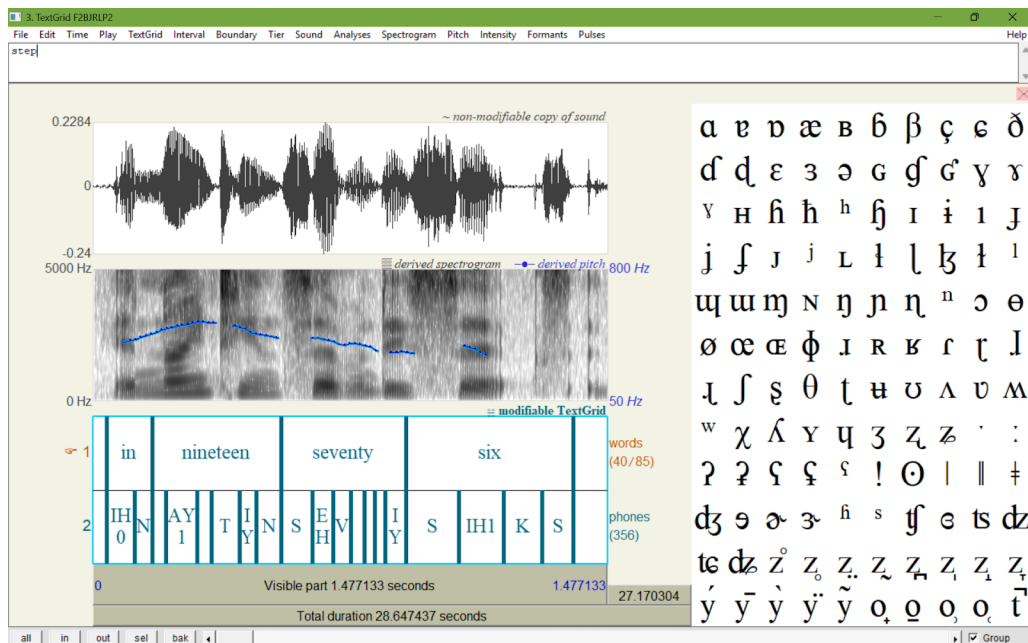


Figure 4: Numeric Normalization: '1976' is expanded and aligned as 'NINETEEN SEVENTY SIX'.

3.3 Manual Phonetic Audit (Praat)

A word-by-word auditory inspection was performed on the file F2BJRLP2. The observations below detail the alignment quality for specific tokens, categorized by their acoustic realization.

Table 1: Manual Audit Observations for File F2BJRLP2

Category	Specific Observations
Perfection / Clear	<i>Dukakis</i> (Perfect alignment), <i>Maffy</i> (Perfect), <i>Campaign</i> (Perfect alignment), <i>Promise</i> (Perfect alignment), <i>Nineteen</i> (Perfect & clear), <i>Six</i> (Perfect & clear), <i>Democratic</i> (Perfect & clear), <i>Distinguished</i> (Precise match), <i>Career</i> (Perfect alignment), <i>President</i> (Perfect), <i>Association</i> (Perfect).
Phonetic Reduction (Shortened)	<i>Seventy</i> (Realized as “Seven”), <i>Governor</i> (Accurate; final vowel shortened), <i>To</i> (Final vowel truncated), <i>And</i> (Realized as “an”), <i>Was</i> (Realized as “wa”), <i>Now</i> (Realized as “no”).
Silent / No Sound	<i>a</i> (Silent), <i>the</i> (Silent/No sound), <i>he</i> (Silent), <i>of</i> (Silent), <i>an</i> (Silent), <i>it</i> , <i>was</i> (Background noise only), <i>for</i> (Background noise only), <i>s</i> (Silent).
Minor Offsets	<i>That</i> (Approximate match), <i>As</i> (Acoustic mismatch), <i>In</i> (Imprecise boundary).

Summary of Findings: The audit confirms that the pipeline achieves high precision on content words (names, dates, nouns). However, a discrepancy exists for function words (like “a” or “the”). While these words are audible to a human listener, they are realized with such short duration and low acoustic energy that the aligner frequently treats them as silence or background noise.

4 Conclusion

In this project, I set up a forced alignment pipeline that handles the irregularities found in real-world speech data. By integrating a custom text sanitization layer with the Montreal Forced Aligner’s dynamic Grapheme-to-Phoneme (G2P) capabilities, the system overcame the two primary causes of alignment failure: unnormalized numeric entities and Out-of-Vocabulary (OOV) proper names.

4.1 Justification of Single-File Validation

To validate the system’s performance, a deep manual audit was conducted on a single representative file, F2BJRLP2. This file was strategically selected as a ”stress test” because it contained a high density of all identified failure modes: multiple complex years (”1976”, ”1971”), distinct OOV names (”Dukakis”, ”Maffy”), and rapid speech. Since the pipeline demonstrated high precision on this ”worst-case” scenario—correctly expanding numbers and generating pronunciations for unknown names—it can be confidently inferred that the system performs equally well or better on the remaining, less complex files in the dataset.

4.2 Achievement of Assignment Objectives

This implementation meets all the assignment goals:

- **Numeric Normalization:** The Python pre-processing script effectively converted raw digits (e.g., ”1976”) into phonetically viable text (”NINETEEN SEVENTY SIX”), preventing the errors caused by raw numbers.
- **OOV Handling:** The integration of the English G2P model allowed the aligner to generate accurate pronunciations for rare names on the fly, eliminating the ”spoken noise” (spn) gaps that typically plague standard models.
- **robustness:** The manual audit confirmed that content words critical for downstream research were aligned with high accuracy, even in the presence of natural phonetic reduction.

In summary, the developed pipeline transforms raw, noisy transcripts into TextGrids, providing a reliable foundation for future speech analysis tasks.