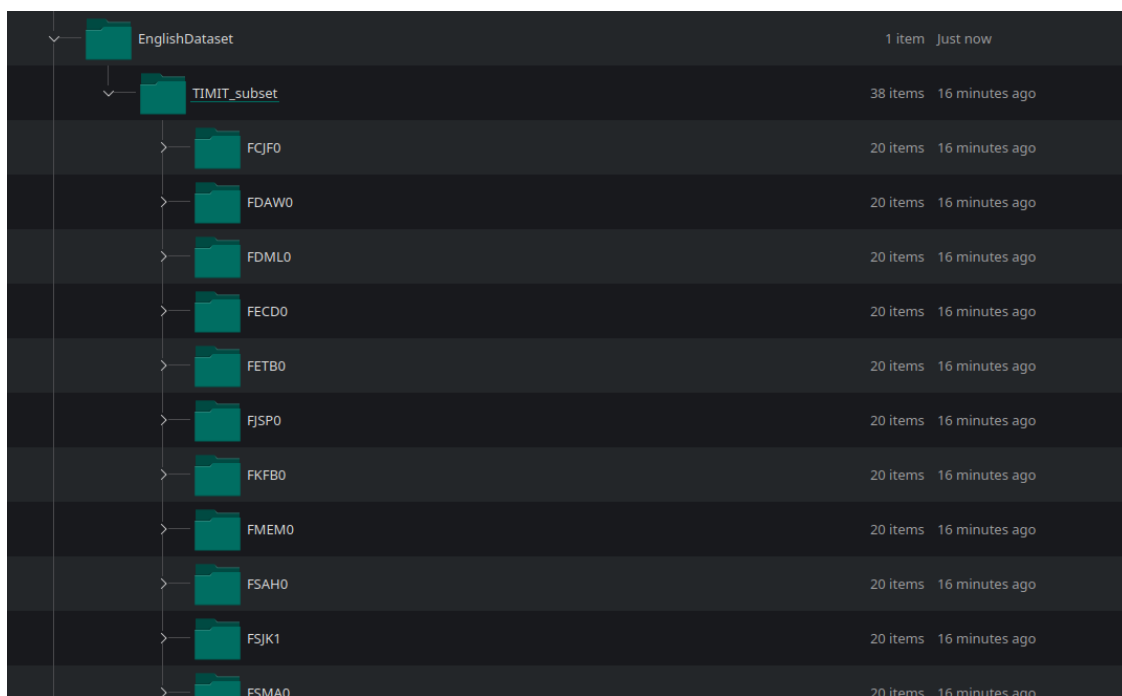# Use Case 3: English alignments

We will align English data in this tutorial. The dataset we use is a subset of the DARPA TIMIT Acoustic-Phonetic Continuous Speech database (Garofolo et al. 1993). The full corpus can be obtained here: https://www.kaggle.com/datasets/mfekadu/darpa-timit-acousticphonetic-continuous-speech .

## 1) Create a corpus folder and extract the data

You can either directly etract the TIMIT subset somewhere on your computer, or create a new folder that will contain the raw data and the alignments. For this use case, I generated a folder called "EnglishDataset" and deployed the data there. The folder structure should look like this, with several subdirectories for each speaker:
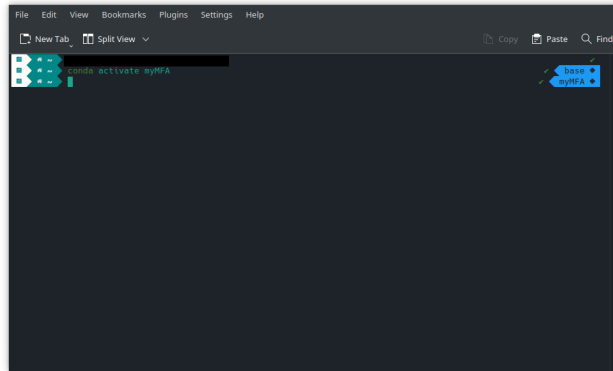


This subset is already prepared for the alignment in the MFA, meaning that the .txt files containing the utterances where already transformed into .lab files. Note that you can use also the respective .txt files in the original corpus, but you have to replace the .txt suffix with .lab and you have to remove the time stemps within the files.

## 2) Open Anaconda and activate your MFA environment

Open your terminal/console/command line and activate the environment in which the MFA installed:

```
[ ]: conda activate myMFA
```

## 3) Download the acoustic model and dictionary for English

In this step you have to download the acoustic model and dictionary for the English alignments. These can be found here: https://mfa-models.readthedocs.io/en/latest/ . If you search for English, you will find two different acoustic models. You will see that these models differ by the feature extraction and amount of data that were used for training.

# English (US) ARPA acoustic model v2.0.0a

- **Maintainer:** Montreal Forced Aligner
- **Language:** English
- **Dialect:** General American English
- **Phone set:** ARPA
- **Model type:** Acoustic
- **Features:** MFCC
- **Architecture:** gmm-hmm
- **Model version:** v2.0.0a
- **Trained date:** 2022-05-11
- **Compatible MFA version:** v2.0.0
- **License:** CC BY 4.0
- **Citation:**

```
@techreport{mfa_english_us_arpa_acoustic_2022,
        author={McAuliffe, Michael and Sonderegger, Morga
        title={English (US) ARPA acoustic model v2.0.0a},
        address={\url{https://mfa-models.readthedocs.io/a
        year={2022},
        month={May},
}
```

- If you have comments or questions about this model, you can check previous MFA model discussion posts or create a new one.

🔔 **Training corpora**

- LibriSpeech English:
    - **Hours:** 982.10
    - **Speakers:** 2,484
    - **Utterances:** 292,367

🔤 **Pronunciation dictionaries**

- English (US) ARPA dictionary v2.0.0a

## Installation

Install from the MFA command line:

```
mfa model download acoustic english_us_arpa
```

## English MFA acoustic model v2.0.0a

- **Maintainer:** Montreal Forced Aligner
- **Language:** English
- **Dialect:** N/A
- **Phone set:** MFA
- **Model type:** Acoustic
- **Features:** MFCC + pitch
- **Architecture:** gmm-hmm
- **Model version:** v2.0.0a
- **Trained date:** 2022-05-14
- **Compatible MFA version:** v2.0.0
- **License:** CC BY 4.0
- **Citation:**

```
@techreport{mfa_english_mfa_acoustic_2022,
        author={McAuliffe, Michael and Sonderegger, Morga
        title={English MFA acoustic model v2.0.0a},
        address={\url{https://mfa-models.readthedocs.io/a
        year={2022},
        month={May},
}
```
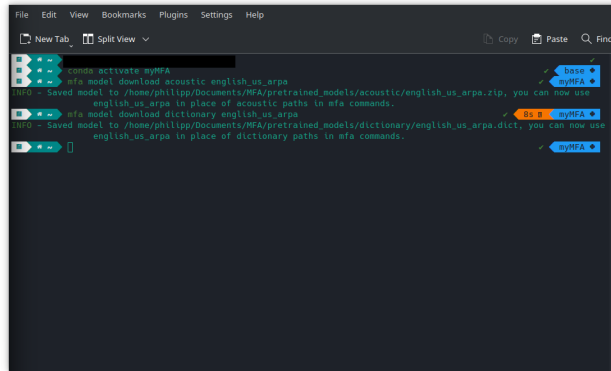
- If you have comments or questions about this model, you can check previous MFA model discussion posts or create a new one.

🔔 **Training corpora**

- Common Voice English v8.0:
  - **Hours:** 2479.95
  - **Speakers:** 74,811
  - **Utterances:** 1,781,717
- LibriSpeech English:
  - **Hours:** 982.10
  - **Speakers:** 2,484
  - **Utterances:** 292,367
- Corpus of Regional African American Language v2021.07:
  - **Hours:** 124.31
  - **Speakers:** 193
  - **Utterances:** 236,792
- Google Nigerian English:
  - **Hours:** 5.77

These two models are also accompanied by different dictionaries. We fill first download the ARPA acoustic model and dictionary:

```
[ ]: mfa model download acoustic english_us_arpa
     mfa model download dictionary english_us_arpa
```
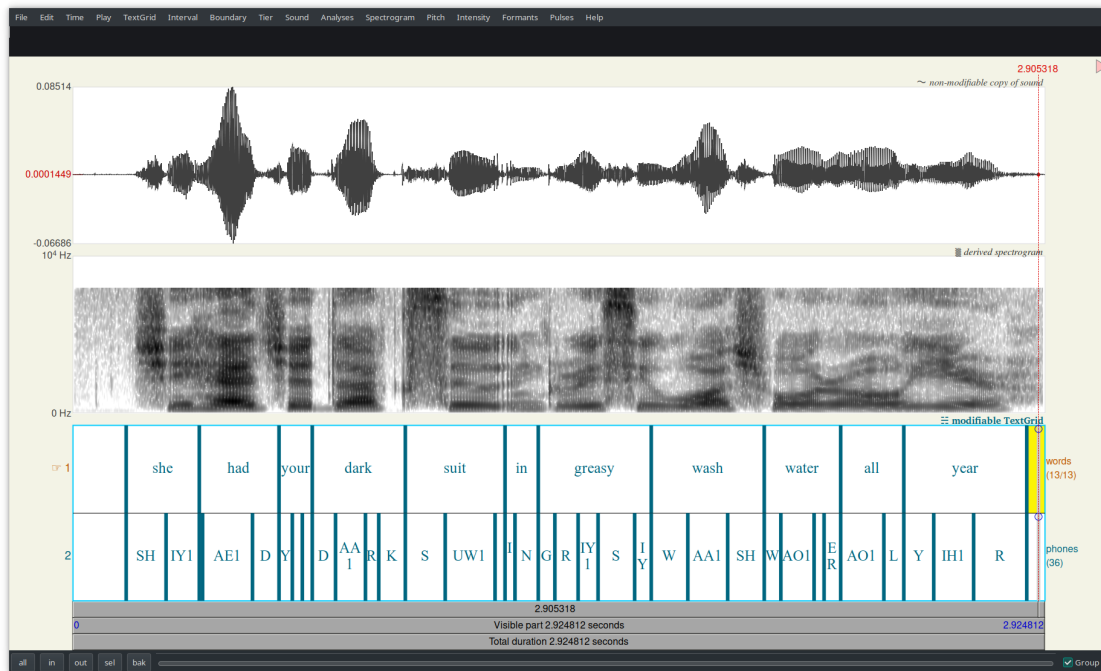
## 4) Align the English data

Since we are using the ARPA model, enter the following command to generate the alignments. Adapt the paths, if necessary.

```
[ ]: mfa align --clean ~/Documents/EnglishDataset/TIMIT_subset english_us_arpa
     ↪english_us_arpa ~/Documents/EnglishDataset/ARPA_alignments/
```

After a successfull run, your should see the alignments in the ARPA_alignments folder. Inspect the TextGrids in Praat.
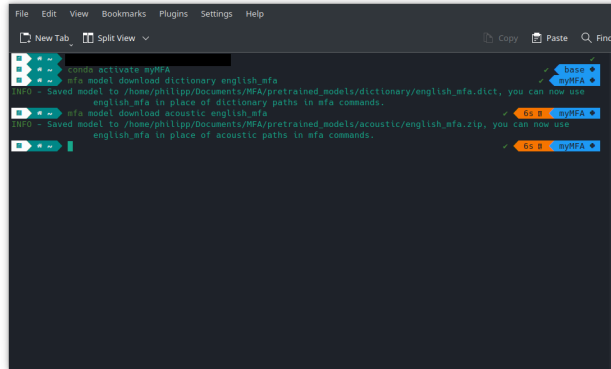


You will see that the phone labels are not in IPA as before, but in the ARPABET ( https://en.wikipedia.org/wiki/ARPABET ). Lets compare the labels and alignments with the second dictionary from the MFA.

## 5) Download the English MFA acoustic model and dictionary

Use the following commands to download the acoustic model and dictionary of the MFA:

```
[ ]: mfa model download acoustic english_mfa
     mfa model download dictionary english_mfa
```
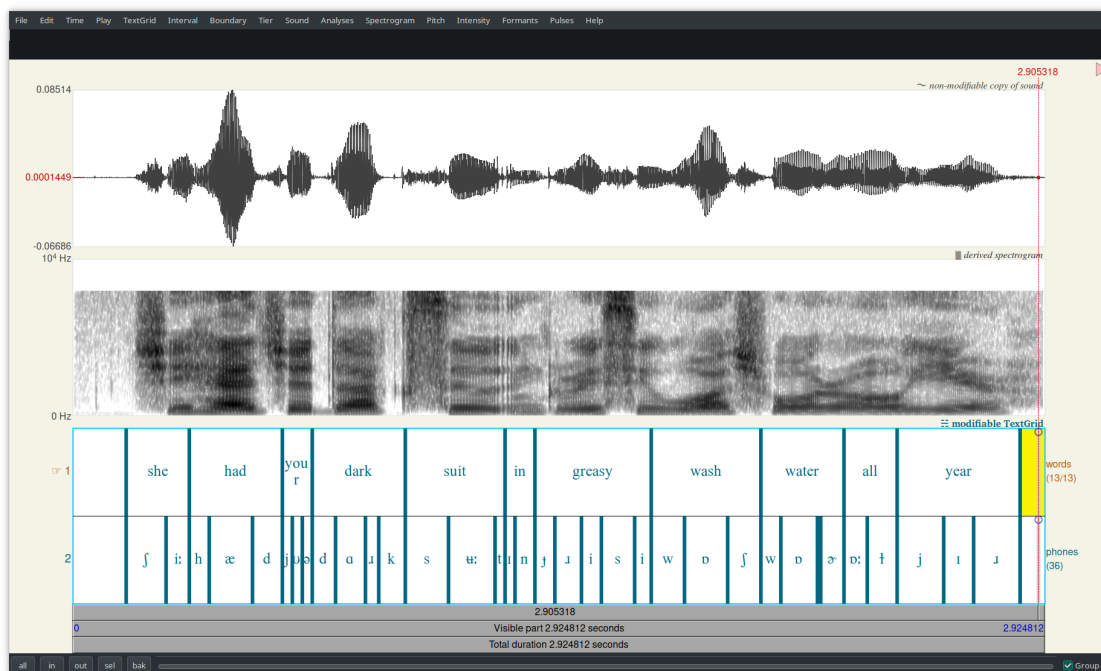


By inspecting the dictionaries under ~/Documents/MFA/pretrained_models/dictionary/ (Windows: C:/Users/username/Documents/MFA/pretrained_models/dictionary/), it can be seen that they differ not only in the transcription standard (ARPA: ARPABET, MFA: IPA) but also to the number of entries (ARPA: 206.366, MFA: 200.695).
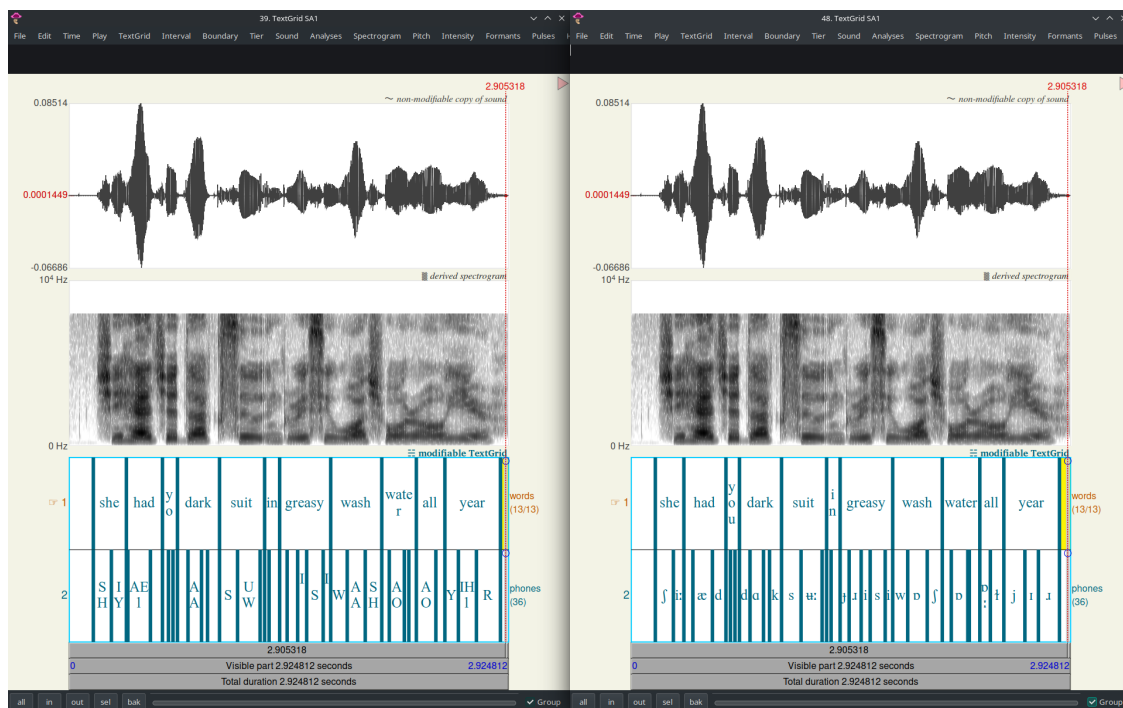
## 5) Align the English data using the MFA model

In order to align the data with the MFA model, use the code above, but replace "english_us_arpa" with "english_mfa" for both the dictionary and the acoustic model:

```
[ ]: mfa align --clean ~/Documents/EnglishDataset/TIMIT_subset english_mfa␣
     ↪english_mfa ~/Documents/EnglishDataset/MFA_alignments/
```

Under MFA_alignments you should find the TextGrids to use them in Praat:

You can now also compare the alignments of the two models:



There should not be too much difference between the segmentation. However, you may prefer one acoustic model over the other.
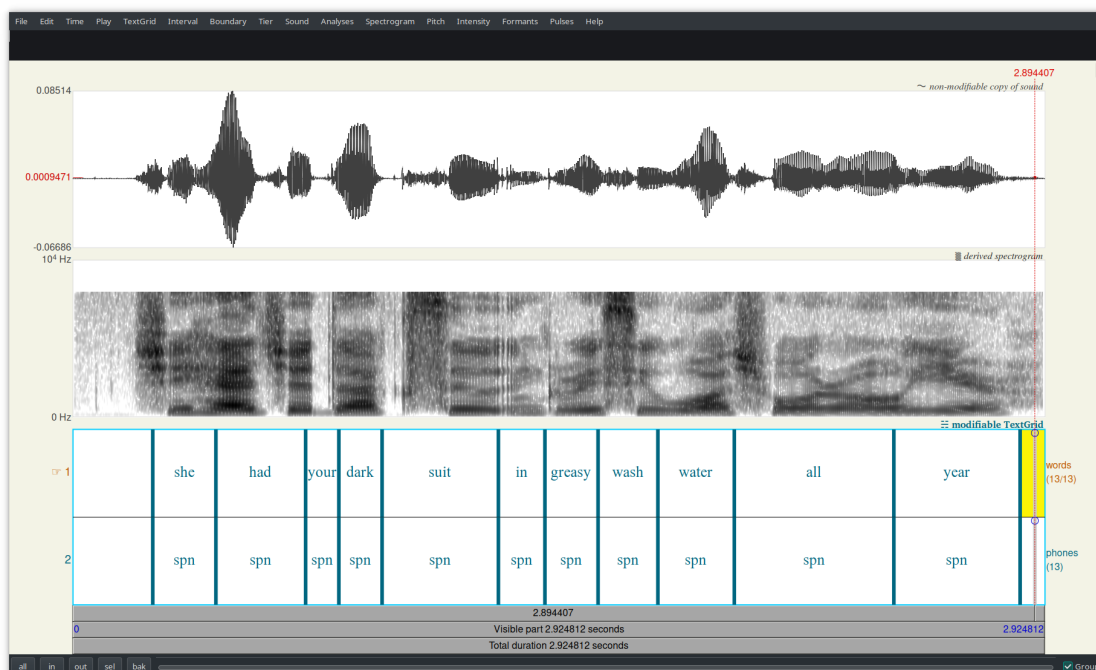
## Final remark

We saw that for a language like English we can have multiple acoustic models and dictionaries. Two points are noteworthy: First, if you plan to extend a dictionary manually, ensure to use the transcription standard for the phones. Otherwise the respective word will note be recognized. Second, do not use an acoustic model and another dictionary that was not used for training. For example, if we want to align the English data with the English MFA dictionary but with the ARPA acoustic model:

```
[ ]: mfa align --clean ~/Documents/EnglishDataset/TIMIT_subset english_mfa␣
     ↪english_us_arpa ~/Documents/EnglishDataset/MFA_alignments/
```

you will see this error message:



The reason is, that the phoneset of the MFA dictionary is different from the one that was used for the ARPA acoustic model. The MFA is thus not able to recognize the specific phones of the transcriptions of the words. This leads to TextGrids containing mis-alignments in the word tier and "spn" labels (default label to model unknown words) in the phone tier.

# References

Garofolo, John S., et al. TIMIT Acoustic-Phonetic Continuous Speech Corpus LDC93S1. Web Download. Philadelphia: Linguistic Data Consortium, 1993