

# Autophon user guide

## British English

Model: Montreal Forced Aligner, English version 1.0

### 1 What is forced alignment?

*Forced alignment* (FA) refers to the automatic process by which speech recordings are phonetically time-stamped with the help of Hidden Markov models or Deep Neural Networks. Autophon uses the latter by means of the *Montreal Forced Aligner* (McAuliffe, Socolof, Mihuc, Wagner, and Sonderegger 2017). The software outputs a time-stamped phonetic annotation, readable in Praat (Boersma and Weenink 2017), that is based on an optimization of two user inputs: (1) the speech recording and (2) a corresponding orthographic transcription. For an FA tool to work for a particular language, an acoustic model must be trained and an accompanying pronunciation lexicon must be built that covers every word in the language. FA is important because it automates something that is resource-intensive when done manually. A typical phonetic annotation can take between 250 and 400 minutes per recorded minute. In a place like Scandinavia – where labor costs are high – this cost has presented a barrier for linguists.

### 2 How to cite

Any study that makes use of Autophon English should cite the following sources:

- Boersma, P., & Weenink, D. (2017). Praat: Doing phonetics by computer [Computer software], Version 6.0.36. <http://www.praat.org/>
- McAuliffe, M., Socolof, M., Mihuc, S., Wagner, M., & Sonderegger, M. (2017). Montreal Forced Aligner: Trainable text-speech alignment using Kaldi. *Proceedings of Interspeech*, 498–502.
- Young, N. J. (2023). Autophon – Automatic phonetic annotation of Nordic languages (web application). [www.autophon.se](http://www.autophon.se)

### 3 Instructions

**Logging in** Creating an account and logging in are relatively intuitive. You must verify your email address before you can use Autophon. Wait at least 1 hour for your verification email.

**Cost** Autophon is free of charge to the public.

**Uploading files** Files can be uploaded individually, as a group, or in a zip file. Every transcription file must have a corresponding audio file by the same name, e.g., *Micke\_spontan.TextGrid* and *Micke\_spontan.wav*. You can upload/align a maximum of five files, totaling 350 MB, at a time. The following formats are supported:

**Transcription formats:** Transcription files can be made in Praat (**.TextGrid**) (soon also ELAN **.eaf**) and must only have *one tier each*. Multiple-tier transcriptions will be rejected. If you have a dialogue, extract each tier as a separate file, duplicate the sound file, and give it corresponding names to each extracted tier. TextGrid encoding needs to be either **UTF-8/Unix** or **UTF-16/Windows CRLF**. If you have older Praat files from the mid to early 2000s, they may produce a yellow error box, in which case you should open them on your desktop in a current version of Praat, resave them, and then try uploading them to Autophon again.

**Audio formats:** Audio files can be **AAC(M4A)**, **AC-3**, **AIFF**, **AIFF/24bit**, **AIFF/32bit**, **ALAC**, **FLAC**, **M4R**, **MP3**, **OGG**, **OPUS**, **WAV/8bit**, **WAV/24bit**, **WAV/32bit**, **WAV/A-law**, **WAV/mu-law**, **WMA**. Autophon will automatically consolidate stereo files to mono, which can compromise quality. Therefore, you may wish to convert your audio to mono yourself.

**Transcription preparation** Transcriptions should be broken down into chunks that are between one and 20 words, and the boundary demarcations should have at least 0.01 seconds of buffer before and after the speech stream. This is illustrated in Figure 1. In natural-speech data, sometimes a researcher will encounter utterances of rapid speech with more than 20 words that have no pauses. In these cases, one must violate one rule in order to appease the other, and it is unclear to us at Autophon what the best approach is. Either force in a break to stick to the 20-word recommendation – violating the 0.01-second buffer – or permit the utterance to go beyond 20 words.

Some users will have transcriptions of single words or phrases that are segmented in Praat down to the exact start and finish time of the speech stream. Autophon will unfortunately perform poorly and not maintain the manual start and finish times. This, however, is something we would be interesting in remedying if a user expresses a need for his/her project.

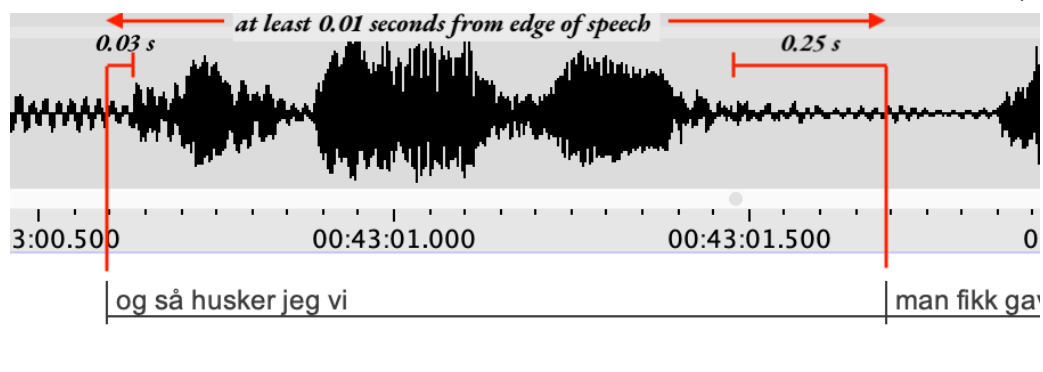


Figure 1: A sample transcription that follows our recommendations: a five-word chunk whereby the start boundary is approximately 0.03 seconds from the speech stream and the end boundary is approximately 0.25 seconds from the speech stream.

**Select a language** Once the files are uploaded, your files will appear in the topmost upload list where you can proof them before alignment. Autophon will automatically suggest a language for you. You can, of course, change this selection.

**Proofing files** In the upload list, metrics are provided for your files, including word count, file size, language, and missing words. It is here you can either delete the upload and start over, change the language, add words, or proceed with alignment by clicking “Align!”.

**Missing words** Alignment is executed with a language-specific model and a language-specific pronunciation dictionary. These pronunciation dictionaries are finite and do not contain every word in the language. Therefore, the “missing words” column outputs any word in your transcription that is not found in Autophon’s dictionary. Specifically, the following process happens: (1) Your transcription is reconciled against Autophon’s dictionary. (2) Missing words are extracted. (3) A grapheme-to-phoneme script that was trained on the original dictionary is run on each missing word to offer suggested pronunciations. (4) These are reported in the “missing words” column, available for download. You then have the option to either accept these suggestions or change them and upload your own changes into the user dictionary (see next section).

**Your custom pronunciations** You may decide that you either (a) do not agree with Autophon’s suggested pronunciations in the “missing words” column or (b) do not agree with the pronunciations in the back-end dictionary. In this box you can enter your own pronunciations that will override *both*. Entries must be done in the ASCII-bet specific to the language at hand (DanFAbet – Danish, ARPAbet – English, NoFAbet – Norwegian, etc.). The phoneme key can be accessed in Table 1.

Entries can either be made directly into the dictionary box or they can be uploaded from a txt file. You are limited to 1 million characters. Format entries as LOOKUP-SPACE/TAB-PHONEME-SPACE-PHONEME. Note that each phoneme must be separated by a space and that *every vowel must have a numerical stress attached to it (unstressed vowels take zero)*. Note also that the lookup cannot be two or more words because that will confuse Autophon and make it treat the second word as a phoneme.

If you would like Autophon to optimize and select from more than one pronunciation option, simply enter in multiple pronunciations for the same lookup.

**Correct:**

```
jabberwocky D J AEEH1 B EH0 W AO2 K IH0
jabberwocky J AEEH1 B EH0 W AO2 K IH0
```

**Incorrect:**

```
jabberwocky D J AEEH1 B EH W AO2 K IH
jabberwocky D J AEEH BEH WAO2 K IH
jabber wocky J AEEH1 B EH0 W AO2 K IH0
```

**Aligning files** Click “Align!” to the far right of the upload list to initiate alignment. This can take between a few minutes and an hour, depending on the amount of users accessing the aligner at that moment. You may have to log out and then log back in to see the outputted files.

**Downloading the annotations** Once alignment is finished, your annotations will be available in the bottom-side download list. You may have to log out and then log back in to see the outputted files. Click on the item to download the folder. The annotations will be in **TextGrid** format and are only readable in Praat (Boersma and Weenink 2017).

## 4 Phoneme key

Autophon will output two versions of the same TextGrid for every file you align: (1) in ARPAbet and (2) in IPA. ARPAbet is the ASCII-based phoneme coding that is specific to the English language<sup>1</sup>. The key is located in Table 1.

ARPAbet	IPA	ex.	ARPAbet	IPA	ex.	ARPAbet	IPA	ex.	ARPAbet	IPA	ex.
Vowels			AY	aɪ	<i>bite</i>	L	l	<i>lie</i>	ZH	ʒ	<i>pleasure</i>
AA	ɑ	<i>father</i>	EY	eɪ	<i>bait</i>	M	m	<i>my</i>	Syllabic consonants		
AE	æ	<i>bat</i>	OY	ɔɪ	<i>boy</i>	N	n	<i>nigh</i>	ER	ɚ	<i>bird, foreword</i>
AH	ʌ	<i>butt</i>	Consonants			NG	ŋ	<i>sing</i>	Lexical stress		
AO	ɔ	<i>caught</i>	B	b	<i>buy</i>	P	p	<i>pie</i>	AA0	ɑ	<i>ban<u>ana</u></i>
EH	ɛ	<i>bet</i>	CH	tʃ	<i>China</i>	R	r	<i>rye</i>	AA1	'ɑ	<i>ban<u>ana</u></i>
IH	ɪ	<i>bit</i>	D	d	<i>die</i>	S	s	<i>sigh</i>	AA2	,ɑ	<i>bar<u>nyard</u></i>
IY	i	<i>beat</i>	DH	ð	<i>thy</i>	SH	ʃ	<i>shy</i>			
OW	oo	<i>boat</i>	F	f	<i>fight</i>	T	t	<i>tie</i>			
UH	ʊ	<i>book</i>	G	g	<i>guy</i>	TH	θ	<i>thigh</i>			
UW	u	<i>boot</i>	HH	h	<i>high</i>	V	v	<i>vie</i>			
Diphthongs			JH	dʒ	<i>jive</i>	W	w	<i>wise</i>			
AW	aʊ	<i>bout</i>	K	k	<i>kite</i>	Y	j	<i>yacht</i>			
						Z	z	<i>zoo</i>			

Table 1: Phoneme key: ARPAbet, IPA, and lexical examples

Every ARPAbet vowel is followed by a numerical code that denotes suprasegmental information. XX0 refers to lexically unstressed vowels; XX1 – primary lexical stress; XX2 – secondary lexical stress.

## 5 Performance and metrics

The model used here is the original English model developed for version 1.0 of the Montreal Forced Aligner. Metrics for this specific model can be accessed in (McAuliffe, Socolof, Mihuc, Wagner, and Sonderegger 2017) (attached as an Appendix here). The pronunciation dictionary used is the CMU dictionary, altered to reflect Southern British pronunciation.

## 6 Data security

Everything you upload is encrypted and sent to a server in Frankfurt, Germany, that is run by *Digital Ocean*. Transcriptions and sound files are deleted immediately after alignment, which significantly reduces the chance of a data breach. On the other hand, finished TextGrids are stored in your account for as long as you like. Once, however, you delete them, they will be removed from our server permanently.

## 7 Features and limitations

### What Autophon is

Autophon is a frontend web application for the Nordic languages that uses the Montreal Forced Aligner (MFA) (McAuliffe, Socolof, Mihuc, Wagner, and Sonderegger 2017) as a core component of its backend. The language-specific models and pronunciation dictionaries were constructed by Nate Young. The language-specific models are trained on various corpora, and the pronunciation dictionaries are usually adaptations of existing dictionaries available online.

The main advantages of using Autophon are:

1. Autophon is a web app, which means it is OS-agnostic.
2. As a web app, Autophon also requires less computational knowledge, which expands access to more researchers and students.
3. Autophon takes nearly all types of transcription and sound formats. It runs your transcriptions and sound files through a converter to ensure that they can be aligned with the MFA backend.

<sup>1</sup><http://www.speech.cs.cmu.edu/cgi-bin/cmudict>

4. Autophon offers an easy way to add the pronunciation of missing words by integrating grapheme-to-phoneme algorithms into the user interface.
5. Autophon allows the user to have access to the latest language training models without needing to constantly check for updates.

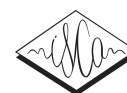
### What Autophon is *not*

Important limitations are:

1. Autophon is no magic bullet. You need a highly accurate orthographic transcription for it to work. And even then, you may not be satisfied with the results.
2. Autophon varies in its accuracy, and this accuracy depends on the language. Accuracy metrics are provided above.
3. Autophon will be slower to implement core MFA updates because it consists of layers and layers of code packed around MFA.

### References

- Boersma, P., & Weenink, D. (2017). Praat: Doing phonetics by computer [Computer software], Version 6.0.36. <http://www.praat.org/>
- McAuliffe, M., Socolof, M., Mihuc, S., Wagner, M., & Sonderegger, M. (2017). Montreal Forced Aligner: Trainable text-speech alignment using Kaldi. *Proceedings of Interspeech*, 498–502.
- Young, N. J. (2023). Autophon – Automatic phonetic annotation of Nordic languages (web application). [www.autophon.se](http://www.autophon.se)



## Montreal Forced Aligner: trainable text-speech alignment using Kaldi

Michael McAuliffe<sup>1</sup>, Michaela Socolof<sup>2</sup>, Sarah Mihuc<sup>1</sup>, Michael Wagner<sup>1,3</sup>, Morgan Sonderegger<sup>1,3</sup>

<sup>1</sup>Department of Linguistics, McGill University, Canada

<sup>2</sup>Department of Linguistics, University of Maryland, USA

<sup>3</sup>Centre for Research on Brain, Language, and Music, McGill University, Canada

michael.mcauliffe@mail.mcgill.ca, msocolof@umd.edu, sarah.mihuc@mail.mcgill.ca,  
chael@mcgill.ca, morgan.sonderegger@mcgill.ca

### Abstract

We present the Montreal Forced Aligner (MFA), a new open-source system for speech-text alignment. MFA is an update to the Prosodylab-Aligner, and maintains its key functionality of trainability on new data, as well as incorporating improved architecture (triphone acoustic models and speaker adaptation), and other features. MFA uses Kaldi instead of HTK, allowing MFA to be distributed as a stand-alone package, and to exploit parallel processing for computationally-intensive training and scaling to larger datasets. We evaluate MFA's performance on aligning word and phone boundaries in English conversational and laboratory speech, relative to human-annotated boundaries, focusing on the effects of aligner architecture and training on the data to be aligned. MFA performs well relative to two existing open-source aligners with simpler architecture (Prosodylab-Aligner and FAVE), and both its improved architecture and training on data to be aligned generally result in more accurate boundaries.

**Index Terms:** forced alignment, automatic segmentation, acoustic analysis

### 1. Introduction

In *forced alignment*, speech and its corresponding orthographic transcription are automatically aligned at the word and phone level, given a way to map graphemes to phonemes (typically a pronunciation lexicon) and a statistical model of how phones are realized. Forced alignment has become widely used in scientific research on language over the past ~10 years, including in sociolinguistics, phonetics, language documentation, and psycholinguistics (e.g. [1, 2, 3, 4, 5]). This use has been driven by the availability of accurate, pre-built, and easily usable aligners, such as FAVE/P2FA, (Web)MAUS, and Prosodylab-Aligner [6, 7, 8]. We focus on this broad use case: forced alignment for language sciences using publicly-available software, when at least an orthographic transcript is available.<sup>1</sup>

Many such forced aligners now exist (e.g. [6, 7, 8, 12, 13, 14, 15, 16, 17]), which differ in two key ways. First, in *architecture*, including the acoustic model used to model the realization of phones, and whether the acoustic features are transformed to account for speaker variability. Second, in *trainability*: most aligners ship with pre-trained acoustic models only, while others can be retrained on new data [8, 17].

We describe the Montreal Forced Aligner (MFA), new open-source forced alignment software which is a successor to the Prosodylab-Aligner. MFA maintains Prosodylab-Aligner's

trainability and updates its architecture. MFA uses triphone acoustic models to capture contextual variability in phone realization, in contrast to monophone acoustic models used in Prosodylab-Aligner and other current aligners (e.g. FAVE). MFA also includes speaker adaptation of acoustic features to model interspeaker differences. MFA uses the Kaldi speech recognition toolkit [18], which offers advantages over the HTK toolkit underlying most existing aligners.

We evaluate MFA's performance on detecting word and phone boundaries in laboratory and conversational speech. Our experiments test whether the more complex architecture and trainability of MFA affect performance, by comparing to two existing monophone acoustic model aligners and varying the training data.

### 2. Montreal Forced Aligner

MFA is an open-source command line utility, with prebuilt executables for Windows and Mac OSX, and online documentation.<sup>2</sup> MFA is built on top of Kaldi, an actively maintained, open-source automatic speech recognition toolkit [18], and has three key usability features: it builds on the *trainability* of Prosodylab-Aligner, and improves *portability* and *scalability*. The use of Kaldi as the ASR toolkit rather than HTK allows for easier distribution due to Kaldi's more permissive license, so no compilation from source is required by the user. MFA's use of Kaldi is highly parallel, which mitigates run time when using larger corpora and more computationally-intensive training.

The ASR pipeline that MFA implements uses a standard GMM/HMM architecture, adapted from existing Kaldi recipes. To train a model, monophone GMMs are first iteratively trained and used to generate a basic alignment. Triphone GMMs are then trained to take surrounding phonetic context into account, along with clustering of triphones to combat sparsity. The triphone models are used to generate alignments, which are then used for learning acoustic feature transforms on a per-speaker basis, in order to make the models more applicable to speakers in other datasets [19]. MFA has been successfully applied to 29 languages from GlobalPhone [20], the NCHLT corpora of South African languages [21], and other corpora.

MFA uses mel-frequency cepstral coefficients (MFCCs) as acoustic features. Thirteen MFCCs are calculated with a 25 ms window size and 10 ms frame shift. The feature calculation has a frequency ceiling of 8 kHz, allowing for acoustic models to be built and used regardless of sampling rate (i.e., models trained on 16 kHz sampled files can be applied to 44.1 kHz sampled files without manual resampling). Delta and delta-delta features from surrounding MFCC frames are also included, giving

<sup>1</sup>We do not address related work, such as on linguistic analysis of untranscribed speech [9], or phoneme boundary detection [10], or text-speech alignment for TTS [11].

<sup>2</sup><https://montrealcorpusools.github.io/Montreal-Forced-Aligner/>

39 features per frame. Following MFCC generation, CMVN is applied to the features on a per-speaker basis to increase robustness to speaker variability. In the final round of training, feature transforms for each speaker are estimated using feature space Maximum Likelihood Linear Regression (fMLLR) [19]. Speaker adaptation is also done when aligning using pre-trained models, but can be disabled for faster alignment.

During training, MFA does 40 iterations of monophone GMM training, with realignment done during 20 of the iterations. Following monophone training, 35 iterations of triphone training are done, with 15 iterations that perform realignment. Speaker-adapted triphone training includes another 35 iterations with 15 realignment iterations, as well as 5 iterations that include fMLLR estimation. Multiprocessing is used extensively during feature calculation and training, allowing MFA to handle training and alignment of large corpora. For instance, the 1000-hour LibriSpeech corpus was aligned in 80 hours (on a desktop using 12 3.4-GHz processors, 32 GB memory), and training from scratch on the 20-hour Buckeye corpus (Sec. 3) took 2 hours (on a laptop using 4 2.5-GHz processors, 8GB memory).

MFA ships with a pre-trained model for English that has been trained on the LibriSpeech corpus [22] (~1000 hours of audiobooks), and pre-trained acoustic models (mostly from GlobalPhone corpora [20]) and grapheme-phoneme models for generating pronunciation dictionaries are publicly available in the online documentation for 20+ languages. A key feature of MFA is trainability of acoustic models on new data, as in the Prosodylab-Aligner [8]. Thus, a user can align their dataset either using pre-trained models, or by training from scratch on the dataset. Alignment can be significantly better when using acoustic models trained from scratch—especially when the dataset to be aligned is sufficiently large and varied. We recommend experimenting with pre-trained models and retraining, as it is an empirical question which method gives better alignments.<sup>3</sup> The experiments in Section 3 address this question.

There are two primary transcription formats used in current forced aligners, exemplified by Prosodylab-Aligner and FAVE. Prosodylab-Aligner aligns short wav files, each with an associated text file specifying the transcription. This format is common to lab speech where individual trials keep speech segments naturally short. FAVE aligns long files containing time-aligned periods of transcribed speech, a format more common to sociolinguistic data and spontaneous speech. MFA supports both formats, building on the Prosodylab-Aligner format and adding support for Praat [23] TextGrids as a way to specify transcriptions in longer sound files. The TextGrid format allows for the user to specify transcriptions for multiple speakers in the same file. The output of alignment is then a TextGrid for each input file, with separate word and phone tiers for each speaker.

MFA contains other upgrades to the Prosodylab-Aligner. Instead of requiring every word in the transcripts to be in the pronunciation dictionary, MFA includes an explicit model for unknown words as having a unique phone, which allows them to be modeled while maintaining alignment of surrounding words. The unknown word’s phone is constructed similarly to the silence phone, and can match any amount of vocal noise or speech (e.g. words of different lengths). Before performing alignment, MFA prompts the user if unknown words are found, including their location, to deal with simple typos for existing words. Anecdotally, MFA’s alignment quality remains very good when up to 5–10% of word types are unknown.

<sup>3</sup>Similarly, disabling speaker adaptation may lead to better alignments if there is little enough data per speaker.

A common source of alignment errors in read speech like audio books or laboratory experiments is deviations from the prompt, such as filled pauses, restarts, or speech errors. Transcriptions of spontaneous speech often contains analogous transcription errors, since listeners are prone to filtering out such deviations. Rather than manual inspection of each audio file for deviations from the transcription, MFA offers a feature from Kaldi to facilitate finding and correcting them. A limited lexicon per utterance is generated, supplemented with frequent words, and a simple speech recognition pass is run on the file to generate a transcript. This generated transcript is compared to the original transcript and deviations are saved to facilitate manual inspection.

### 3. Evaluation

Our evaluation of MFA addresses three questions: (1) how good is the aligner’s performance relative to manual annotation, and what is the effect on performance of the two key aspects of MFA: (2) architecture (acoustic model and speaker adaptation) and (3) trainability? We evaluate MFA’s performance by examining its accuracy on detecting phone and word boundaries in two datasets, representing types of speech commonly used in language research: isolated-word lab speech and conversational interview speech. We compare MFA to two existing widely-used aligners with simpler architectures—FAVE and Prosodylab-Aligner—and vary the training data for aligners where possible.

#### 3.1. Datasets

The first dataset used in our evaluations was the Buckeye Corpus [24], which contains 20.7 hours of conversational speech from 40 speakers. Buckeye comes with manual transcription and boundaries at the phone and word level, which were produced by forced alignment followed by manual correction. The Buckeye phone set represents more subphonemic detail (e.g. flapping) than needed for our evaluations; we thus mapped it to the phone set used in our pronunciation dictionary (see below).

HTK-based aligners, such as FAVE and Prosodylab-Aligner, require relatively short speech chunks. We thus broke up Buckeye into chunks bounded by non-speech (pauses, noise, interviewer speech) of >150 msec marked in the transcription files, using PolyglotDB.<sup>4</sup> Each of these chunks consists of an orthographic transcription and speech, as well as corresponding word and phone-level manual alignments. In our evaluation, the transcription and speech are force-aligned, and the manual alignments used as the gold standard.

Utterances were excluded if they contained words not in the pronunciation dictionary used in evaluation, for comparability between FAVE/Prosodylab-Aligner (which require all words to be in the dictionary) and MFA (which does not).

The second dataset, Phonsay, consists of 48 minutes of lab speech from 45 participants from two experiments. Participants said words in the frame “Please say \_\_\_\_ again”. The target words all contained vowels followed by a consonant: a voiced obstruent, unvoiced obstruent, or sonorant (e.g. *buzz*, *bus*, *bun*). The boundaries of the vowel and the following consonant were hand-annotated, and these manual annotations are the gold standard in our evaluation.

In the evaluation, we examine two kinds of boundaries. First, left and right *word boundaries*, across all words, for Buckeye only. (Most word boundaries in Phonsay were not anno-

<sup>4</sup><https://github.com/MontrealCorpusTools/PolyglotDB>

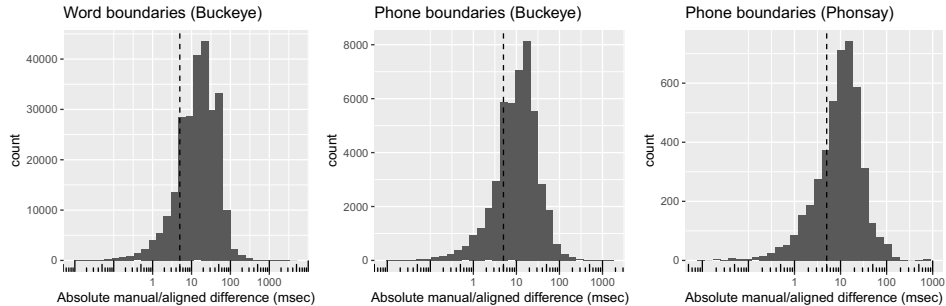


Figure 1: Histograms of absolute differences (on log scale) between force-aligned word and phone boundaries using MFA-LS aligner and gold-standard annotations. Dashed line is at 1/2 frame rate (5 msec), which is a lower bound on average absolute difference.

tated.) Second, *phone boundaries*, for each phone boundary of CVC words in either dataset, that corresponds to a manually-annotated boundary. For Buckeye, this is all four boundaries (denoted .CVC, C.VC, CV.C, CVC.). The CVC words in Buckeye were those from the list of [25], with the additional criterion of having all three segments realized in some way according to the manual transcription. For Phonsay, the boundaries were C.VC, CV.C, and CVC. for the target word in every sentence.

### 3.2. Aligners and training

Our evaluation uses MFA and two HTK-based aligners which are currently used in language research: FAVE, the most widely-used aligner in recent work, and Prosodylab-Aligner (PLA). PLA and FAVE are used as representative of aligners using GMM-HMM monophone acoustic models<sup>5</sup> without speaker adaptation, which are and are not trainable, respectively. Many existing aligners fall into these two categories (e.g. [6, 7, 8, 15]).

In order to minimize out-of-vocabulary words for PLA and FAVE, the pronunciation dictionaries which ship with each of the three aligners were combined into one Arpabet-based dictionary, which was used across all three aligners for training (MFA, PLA) and alignment (MFA, PLA, FAVE).

Both MFA and PLA were trained in two ways: on the LibriSpeech corpus, and on the corpus to be aligned: Buckeye (the subset without unknown words) or Phonsay. For training on LibriSpeech, MFA was trained on the full corpus (~1000 hours), while PLA was trained on the ‘clean’ subset (~450 hours), due to technical difficulties in HTK training on large datasets. For training on Buckeye, we treated the corpus as if only utterance boundaries and the orthographic transcription were known, to simulate the most common case in aligning speech in linguistic research. We refer to the resulting trained aligners as *MFA-LS*, *MFA-Retrained*, *PLA-LS*, and *PLA-Retrained*, where the “retrained” aligners refer to the version trained on Buckeye or the version trained on Phonsay, when discussing each corpus. We also used the existing version of *FAVE*, which uses acoustic models trained on the SCOTUS corpus (25 hours) [26]. Thus, our experiments compare five types of aligner (*MFA-LS*, *MFA-Retrained*), *PLA-LS*, *PLA-Retrained*), *FAVE*).

Each type of aligner was applied to align the Buckeye and Phonsay datasets, resulting in predicted word and phone boundaries. Note that we did not split the datasets into training and

<sup>5</sup>While it is possible to use triphone models in HTK, all distributed software packages for alignment use monophone models.

Table 1: Accuracies at different tolerances (percentage below a cutoff) for absolute differences between force-aligned boundaries using MFA-LS aligner, and gold-standard annotations.

	Tolerance (ms)			
	<10	<25	<50	<100
Word boundaries (Buckeye)	0.33	0.68	0.88	0.97
Phone boundaries (Buckeye)	0.41	0.77	0.93	0.98
Phone boundaries (Phonsay)	0.36	0.72	0.88	0.95

test sets, as the common use case for a trainable aligner is to simultaneously train on and align the entire dataset of interest.

Our evaluation considers two subsets of the predicted boundaries, described above: word boundaries (Buckeye only), and phone boundaries (Buckeye and Phonsay). The metric we use for accuracy of a force-aligned boundary is the absolute difference (in msec) from the manually-annotated boundary.

### 3.3. Results

Our results address questions (1)–(3): how good are MFA’s alignments ‘out of the box’ compared to hand annotation, and do the more complex architecture and trainability of MFA lead to more accurate alignments?

#### 3.3.1. Alignment quality

We first consider the performance of MFA-LS, which is the version distributed with the current version of MFA. Performance on the two datasets approximates the performance a user can expect if MFA-LS is applied to lab (Phonsay) or conversational (Buckeye) English data, without retraining.

Figure 1 and Table 1 show the distribution of manual/force-aligned differences, for each kind of boundary, for the two datasets. The distributions of differences are highly right-skewed, as for other forced aligners [8, 26]: 2–5% of tokens have differences of at least 100 msec, while about 90% have differences of less than 50 msec. Table 2 (row 1) gives the mean and median of manual/aligned boundary differences for each case. These measures can be compared for the Buckeye corpus to differences between human transcribers reported by [27]—bearing in mind that the set of word and phone boundaries used there differs from the set used in our evaluation.

For word boundaries, the mean manual/aligned difference is 24 msec, which is comparable to 26 msec intertranscriber

Table 2: Comparison of aligners in detecting word boundaries (Buckeye only) and phone boundaries (Buckeye and Phonsay). Means and medians are over differences between aligned and gold-standard boundaries.

Aligner	Word bound.		Phone boundaries			
	Buckeye		Buckeye		Phonsay	
	mean (ms)	med (ms)	mean (ms)	med (ms)	mean (ms)	med (ms)
MFA-LS	24.1	15.8	<b>17.0</b>	<b>11.2</b>	25.2	11.3
MFA-Retrained	<b>22.6</b>	<b>15.0</b>	17.3	11.8	<b>16.6</b>	<b>10.8</b>
PLA-LS	30.5	15.6	24.0	13.9	40.1	21.5
PLA-Retrained	27.2	15.6	24.7	15.8	25.9	16.5
FAVE	24.7	16.6	19.3	12.0	21.8	13.0

reliability [27]. 68% of manual/aligned differences are under 25 msec, which is significantly lower than the 90% intertranscriber agreement reported at 26 msec tolerance.

For phone boundaries, the mean difference is 17 msec for Buckeye and 25 msec for Phonsay. For Buckeye, an identical figure (17 msec) is reported for intertranscriber agreement [27]. The median difference is comparable (11 msec) for Phonsay and Buckeye, suggesting that the main difference between them is more gross misalignments for Phonsay (visible in Fig. 1 right).

In sum, MFA performs well across both datasets and boundary types. While phone and word-level alignment is comparable to human annotators on average, the force-aligned boundaries do contain more medium-to-large alignment errors (>25 msec).

### 3.3.2. Architecture

To examine the effect of MFA’s more complex architecture—triphone acoustic models and speaker-adapted features, compared to monophone acoustic models without speaker adaptation—we compare MFA-LS to PLA-LS and FAVE. The comparison with PLA-LS is most important, since MFA is essentially the same as PLA except for this modified architecture.

Rows 1, 3, 5 of Table 2 show, for these three aligners, the mean and median differences between manual and force-aligned boundaries for each condition. In most cases (columns of Table 2), the ordering is MFA-LS < FAVE < PLA-LS. However, MFA-LS and PLA-LS have roughly the same median for word boundaries for Buckeye (below FAVE), and FAVE has the lowest mean for phone boundaries for Phonsay.<sup>6</sup> Still, MFA-LS has the best overall performance of the three aligners. The difference between MFA-LS and PLA-LS suggests that MFA’s different architecture led to better alignments.

To what extent is MFA’s performance in this comparison due to the updated acoustic model versus speaker adaptation? Experiments with a version of MFA with speaker adaptation disabled suggest that it is the triphone acoustic model that primarily accounts for MFA’s performance relative to PLA, with 88%/95% of the performance difference for word/phone boundaries (as measured by mean absolute manual/aligned difference) between PLA-LS and MFA-LS on Buckeye coming from just changing the acoustic model.<sup>7</sup>

### 3.3.3. Experiment 3: Training

To examine the effect of retraining on the dataset to be aligned, we compare MFA-Retrained to MFA-LS and PLA-Retrained to

<sup>6</sup>All comparisons are significant (paired Wilcoxon rank-sum test).

<sup>7</sup>Disabling speaker adaptation gives *better* performance as measured by the median, suggesting that enabling speaker adaptation may induce more gross errors, while increasing mean alignment accuracy.

PLA-LS. This comparison represents a common use case: a researcher has a medium-to-large dataset (say 5–50 hours) of speech from speakers of a given type (e.g. Buckeye: Columbus-dialect adults). She can either re-train the aligner’s acoustic models on this data, or use acoustic models which have been pre-trained on a much larger dataset that contains significant interspeaker variation (e.g. LibriSpeech: 1000 hours). Will training on a smaller amount of more similar data or a larger amount of more variable data give better alignments?

The effect of retraining can be evaluated by comparing rows 1 and 2 of Table 2 for MFA, and rows 3 and 4 for PLA, again examining the mean and median of absolute differences between manual and aligned boundaries. In five cases (Buckeye word boundary mean for MFA/PLA, Phonsay phone boundary mean for MFA and mean/median for PLA), retraining leads to better performance, decreasing the mean or median difference by at least 1 msec. In six of the remaining seven cases, retraining makes little difference (< 1 msec mean or median). In only one case (Buckeye phone boundary median for PLA) does retraining lead to clearly worse performance (> 1 msec difference).

On balance, retraining on the dataset to be aligned often improves alignment accuracy relative to using acoustic models pretrained on a larger dataset—and rarely hurts. However, the discrepancy between mean and median values in some conditions suggests that a more thorough evaluation should examine the effect of retraining on gross alignment errors.

## 4. Conclusion

We have presented a new open-source trainable forced aligner for language research, the Montreal Forced Aligner, which updates the Prosodylab-Aligner. MFA uses more complex acoustic models (triphones), and is built using the Kaldi toolkit instead of HTK. MFA showed good performance in aligning word and phone boundaries in one lab speech dataset and one spontaneous speech dataset. Notably, in each test case (columns of Table 2), it is one of the MFA aligners which gives the most accurate alignment relative to the gold standard.

Our evaluation suggests that both MFA’s more complex architecture and the ability to retrain on new data generally improve performance. Using triphone acoustic models in particular seems to improve accuracy, compared to the monophone models commonly used in current aligners. More complex architectures, such as using artificial neural network models implemented in Kaldi (as in [14]), could improve accuracy further and are planned in future development. Retraining on the data to be aligned generally improved alignment accuracy, though it often had little effect—perhaps reflecting the similarity of training data for all aligners tested (North American English).

The mixed results of our evaluations point to the need for more thorough evaluations of forced aligners, to establish best practices for deploying forced alignment in language research [2, 28, 29]. Future work could examine the conditions under which adding speaker adaptation, or adapting an existing forced aligner versus retraining, improves alignment [30].

## 5. Acknowledgements

We acknowledge funding from SSHRC #430-2014-00018, FRQSC #183356 and CFI #32451 to MS, and SSHRC #435-2014-1504 and the SSHRC CRC program to MW.



## 6. References

- [1] M. Adda-Decker and N. D. Snoeren, "Quantifying temporal speech reduction in French using forced speech alignment," *Journal of Phonetics*, vol. 39, no. 3, pp. 261–270, 2011.
- [2] C. DiCanio, H. Nam, D. H. Whalen, H. Timothy Bunnell, J. D. Amith, and R. C. García, "Using automatic alignment to analyze endangered language data: Testing the viability of untrained alignment," *The Journal of the Acoustical Society of America*, vol. 134, no. 3, pp. 2235–2246, 2013.
- [3] W. Labov, I. Rosenfelder, and J. Fruehwald, "One hundred years of sound change in Philadelphia: Linear incrementation, reversal, and reanalysis," *Language*, vol. 89, no. 1, pp. 30–65, 2013.
- [4] B. Schuppler, M. Ernestus, O. Scharenborg, and L. Boves, "Acoustic reduction in conversational Dutch: A quantitative analysis based on automatically generated segmental transcriptions," *Journal of Phonetics*, vol. 39, no. 1, pp. 96–109, 2011.
- [5] J. Yuan, M. Liberman, and C. Cieri, "Towards an integrated understanding of speaking rate in conversation," in *Proceedings of Interspeech*, 2006, pp. 541–544.
- [6] I. Rosenfelder, J. Fruehwald, K. Evanini, and J. Yuan, "FAVE (Forced Alignment and Vowel Extraction) Program Suite [Computer program]," 2011, available at <http://fave.ling.upenn.edu>.
- [7] T. Kislser, F. Schiel, and H. Sloetjes, "Signal processing via web services: the use case WebMAUS," in *Digital Humanities Conference 2012*, 2012.
- [8] K. Gorman, J. Howell, and M. Wagner, "Prosodylab-aligner: A tool for forced alignment of laboratory speech," *Canadian Acoustics*, vol. 39, no. 3, pp. 192–193, 2011.
- [9] S. Reddy and J. Stanford, "A web application for automated dialect analysis," in *Proceedings of HLT-NAACL*, 2015, pp. 71–75.
- [10] J. Yuan, N. Ryant, M. Liberman, A. Stolcke, V. Mitra, and W. Wang, "Automatic phonetic segmentation using boundary models," in *Proceedings of Interspeech*, 2013, pp. 2306–2310.
- [11] A. Stan, P. Bell, and S. King, "A grapheme-based method for automatic alignment of speech and text data," in *Spoken Language Technology Workshop (SLT), 2012 IEEE*. IEEE, 2012, pp. 286–290.
- [12] A. Pettarin, "Aeneas [computer program]," 2017, available at <https://www.readbeyond.it/aeneas/>.
- [13] R. Fromont and J. Hay, "LaBB-CAT: An annotation store," in *Australasian Language Technology Association Workshop 2012*, vol. 113, 2012, pp. 113–117.
- [14] R. M. Ochshorn and M. Hawkins, "Gentle forced aligner [computer program]," 2017, available at <https://github.com/lowerquality/gentle>.
- [15] J.-P. Goldman, "EasyAlign: an automatic phonetic alignment tool under Praat," in *Proceedings of Interspeech*, 2011, pp. 3233–3236.
- [16] A. Katsamanis, M. Black, P. G. Georgiou, L. Goldstein, and S. Narayanan, "SailAlign: Robust long speech-text alignment," in *Proceedings of Workshop on New Tools and Methods for Very-Large Scale Phonetics Research*, 2011.
- [17] B. Bigi, "SPPAS: a tool for the phonetic segmentations of speech," in *Proceedings of LREC*, 2012, pp. 1748–1755.
- [18] D. Povey, A. Ghoshal, G. Boulianne, L. Burget, O. Glembek, N. Goel, M. Hannemann, P. Motlicek, Y. Qian, P. Schwarz *et al.*, "The Kaldi speech recognition toolkit," in *IEEE Workshop on Automatic Speech Recognition and Understanding*, 2011, pp. 1–4.
- [19] D. Povey and G. Saon, "Feature and model space speaker adaptation with full covariance Gaussians," in *Proceedings of Interspeech*, 2006, pp. 1145–1148.
- [20] T. Schultz, N. T. Vu, and T. Schlippe, "GlobalPhone: A multilingual text & speech database in 20 languages," in *Proceedings of ICASSP*, 2013, pp. 8126–8130.
- [21] E. Barnard, M. H. Davel, C. J. van Heerden, F. De Wet, and J. Badenhorst, "The NCHLT speech corpus of the South African languages," in *Proceedings of SLTU*, 2014, pp. 194–200.
- [22] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: an ASR corpus based on public domain audio books," in *Proceedings of ICASSP 2015*, 2015, pp. 5206–5210.
- [23] P. P. G. Boersma *et al.*, "Praat, a system for doing phonetics by computer," *Glott international*, vol. 5, 2002.
- [24] M. Pitt, L. Dilley, K. Johnson, S. Kiesling, W. Raymond, E. Hume, and E. Fosler-Lussier, *Buckeye Corpus of Conversational Speech (2nd release)*. Department of Psychology, Ohio State University, 2007.
- [25] S. Gahl, Y. Yao, and K. Johnson, "Why reduce? phonological neighborhood density and phonetic reduction in spontaneous speech," *Journal of Memory and Language*, vol. 66, no. 4, pp. 789–806, 2012.
- [26] J. Yuan and M. Liberman, "Speaker identification on the SCOTUS corpus," in *Proceedings of Acoustics '08*, 2008, pp. 5687–5790.
- [27] W. D. Raymond, M. A. Pitt, K. Johnson, E. Hume, M. J. Makashay, R. Dautricourt, and C. Hiltz, "An analysis of transcription consistency in spontaneous speech from the Buckeye corpus," in *Proceedings of Interspeech*, 2002.
- [28] P. Milne, "The variable pronunciations of word-final consonant clusters in a force aligned corpus of spoken French," Ph.D. dissertation, Université d'Ottawa/University of Ottawa, 2014.
- [29] T. Knowles, M. Clayards, M. Sonderegger, M. Wagner, A. Nadig, and K. Onishi, "Automatic forced alignment on child speech: Directions for improvement," *Proceedings of Meetings on Acoustics*, vol. 25, p. 060001, 2015.
- [30] L. MacKenzie and D. Turton, "Crossing the pond: Extending automatic alignment techniques to British English dialect data," 2013, talk given at *New Ways of Analyzing Variation 42*.