



FORCED ALIGNMENT FOR UNDERSTUDIED LANGUAGE VARIETIES

**Lisa M. Johnson, Marianna Di Paolo,
Adrian Bell, Carter Holt**

University of Utah
ICLDC5 | March 2, 2017



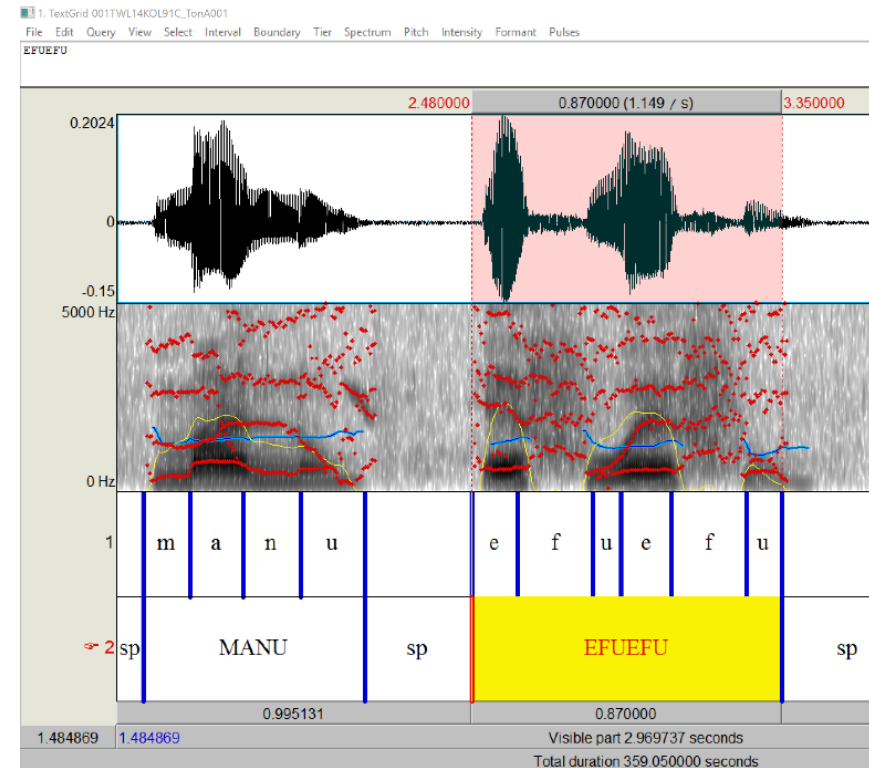
THE PROBLEM

It's necessary in both sociolinguistics, esp. **sociophonetics**, and **language documentation** to efficiently process **large corpora of recorded speech**.

Processing recordings for acoustic analysis is **very time consuming**.

By some estimates, **manual phone-level alignment** may take up to 800x the duration of the audio!

Today we will focus on how these similar processing problems may share similar solutions.

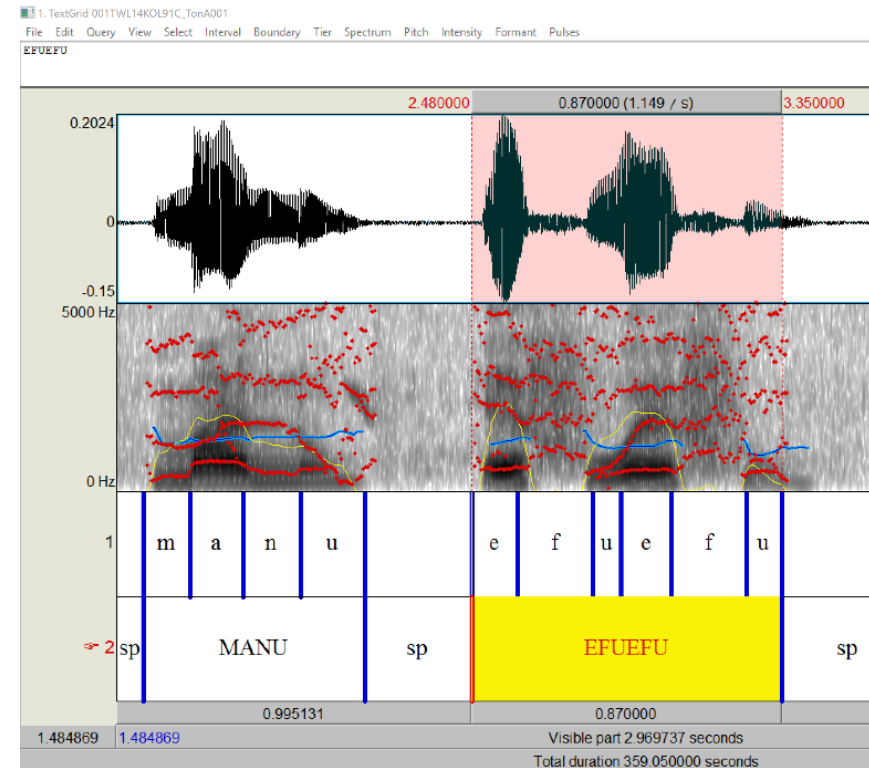


THE PROBLEM

The time and cost associated with processing audio recordings can **limit** the **amount and kinds of data analyzed** and even **the kinds of questions explored**.

Such limitations are especially problematic when they **inhibit work on underdocumented languages**.

Is there a way to expedite this process?



FORCED ALIGNMENT TOOLS

In recent years, new tools have been developed to time-align orthographic transcriptions to recorded speech at the word and phone level.

These forced alignment tools use speech recognition technology to **create a statistical model associating phonetic symbols to speech signals.**

Sociophonetics has benefited greatly from the use of forced alignment technology

- Developed primarily for majority languages like English, with large extant corpora available

Examples: Forced Alignment and Vowel Extraction (FAVE) (Rosenfelder 2013, Rosenfelder et al. 2011); EasyAlign (Goldman 2011); MAUS/WebMAUS (Kisler, Schiel, and Sloetjes 2012); Prosodylab-Aligner (Gorman, Howell, and Wagner 2011); and the Dartmouth Linguistic Automation suite (DARLA) (Reddy and Stanford 2015)

FORCED ALIGNMENT

Two digital tools developed for forced alignment of underdocumented languages:

- Prosodylab-Aligner (PL-A)
- Montreal Forced Aligner (MFA)

(Both developed at McGill University Prosody Lab)

Key features:

- Don't require a pretrained model or a large corpus
- Allow model training and alignment using the same dataset

TONGAN ETHNOLINGUISTIC STUDY

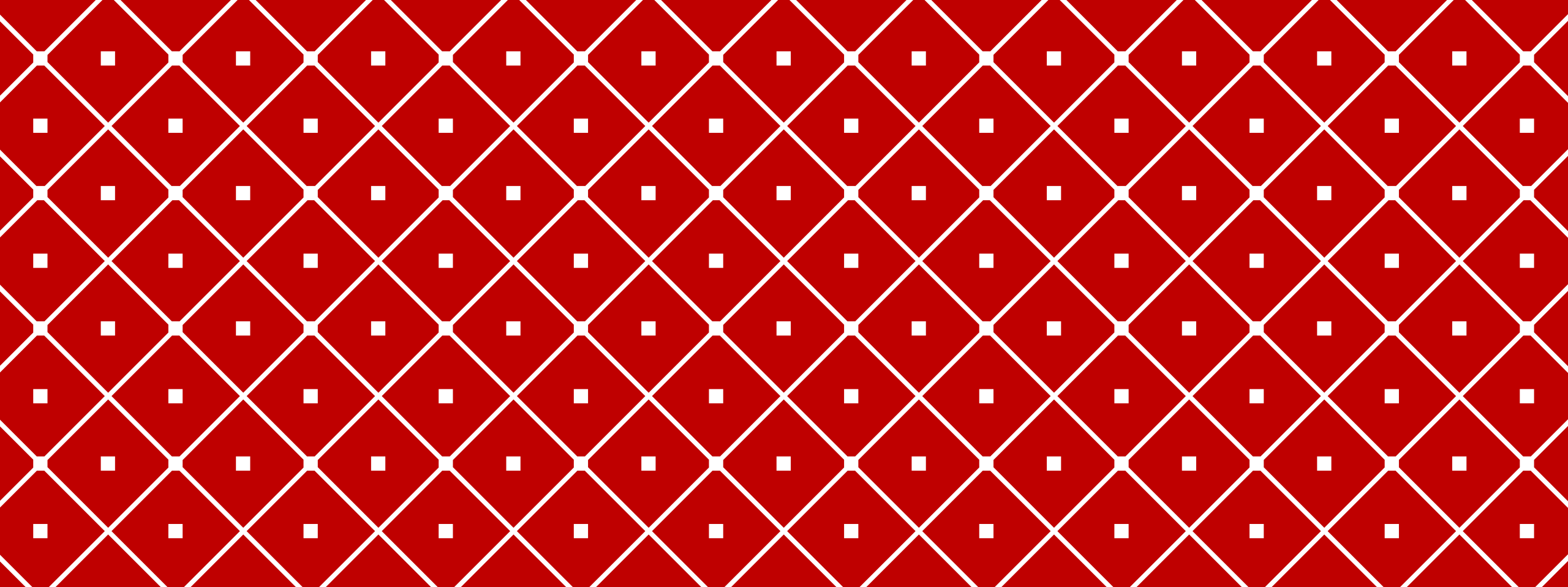
Large-scale ethnographic and linguistic study of post-migration Tongans/Tongan Americans in the U.S. (Adrian Bell, PI)

- Formation of new post-migration ethnolinguistic identities
- Longitudinal and cross-sectional data
- Includes data collection by crowd-sourcing
 - Leads to huge linguistic data set
 - Must expedite the linguistic analysis

To identify potentially important linguistic variables in these newly formed U.S. Tongan American communities, exploring

- **linguistic variation in Tonga**
 - potential Tongan sociolinguistic variables
 - varieties of English used in Tonga
- **linguistic variation** the U.S. **English** contact varieties in **Salt Lake Valley**

Best to use same digital tools for both Tongan & ambient English



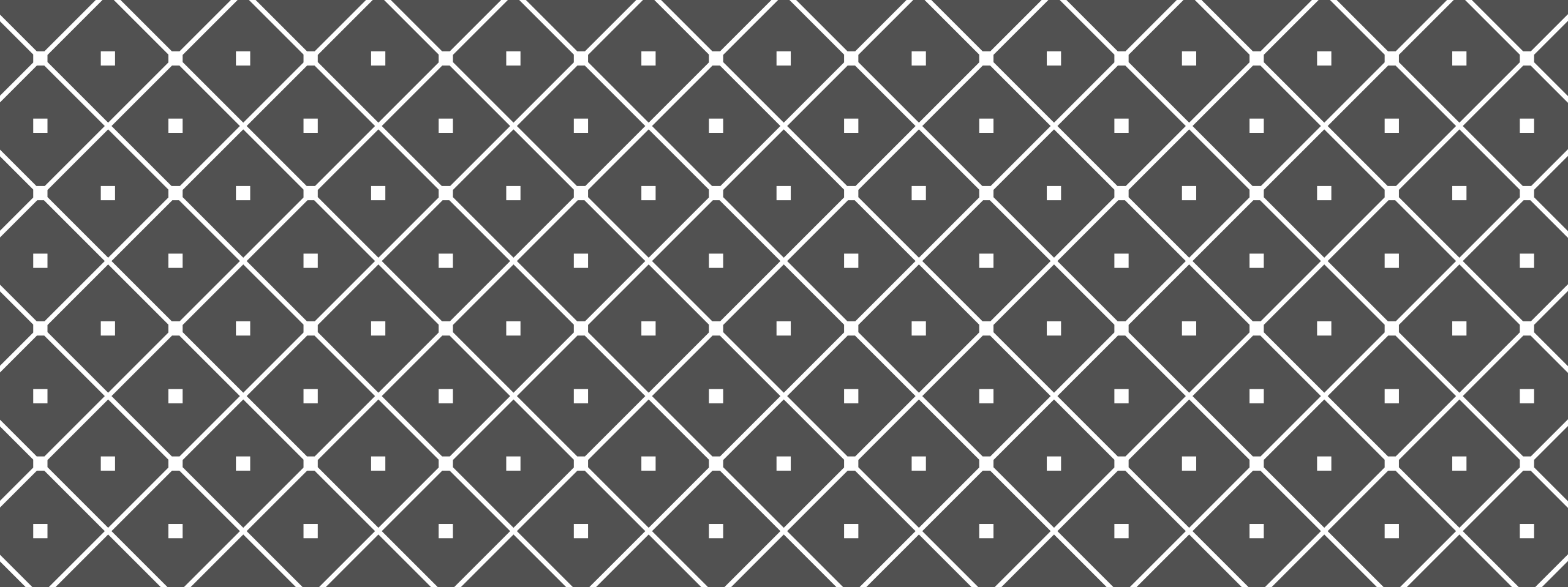
PROSODYLAB-ALIGNER

For Training and Alignment
(Understudied Languages)

WHAT IS PROSODYLAB-ALIGNER?

- A set of scripts that use HTK (Hidden Markov Toolkit) speech recognition software to create time-aligned TextGrid transcriptions
- Designed with laboratory data in mind, best with short audio files
- Includes a pre-trained North American English model
- Supports model training on user-supplied data
- Does not require and time-aligned training data (uses simple text transcriptions)
- Has been used for a variety of majority and minority languages
English (U.S.A., Canadian, British, Aviation, South African), French, Arabic (Gulf), Irish, Cantonese, German, Polish, Mandarin, Tagalog, Spanish, Cho'ol, Mi'gmaq, and Kaqchikel. (Gorman, p.c.)

<http://prosodylab.org/tools/>



GETTING STARTED

Instructions, Issues, and Solutions

WHAT YOU NEED

Requirements/Recommendations

Hardware

- Instructions provided for Mac, Linux, can also be used with Windows

Software Downloads

- Prosodylab-Aligner—GitHub
- Xcode (compiler)—Mac App Store
- HTK (Hidden Markov Toolkit)—HTK website
- Homebrew
- Python
- SoX utilities

What we used

Hardware: Microsoft Surface Pro ³/₄

- Intel Core i5-4300U / 6300U
- 8GB LPDDR3 RAM
- 256 GB SSD (data files on 200GB micro SD)

Software

- Compiled HTK (x64) on Windows using nmake
- Installed Python environment and required packages

WHAT WE LEARNED

The software can be a bit tricky and buggy, but we got it to work.

The Aligner's developer, Kyle Gorman, was very accessible and helpful.

It's good to have one of these on hand



Craig Johnson, programmer

WHAT YOU NEED

Requirements

Audio files (.wav)

- default at 16 kHz (automatically resamples, but you can override)

Example

WHAT YOU NEED

Requirements

Audio files (.wav)

Transcription files (.lab)

- Plain text, UTF-8
- Prescribed format: All caps, single spaces between words, no carriage returns or punctuation, regular spelling conventions (with Unicode characters)

Example

BARACK OBAMA WAS TALKING ABOUT HOW THERE'S A MISUNDERSTANDING THAT ONE MINORITY GROUP CAN'T GET ALONG WITH ANOTHER SUCH AS AFRICAN AMERICANS AND LATINOS AND HE'S SAID THAT HE HIMSELF HAS SEEN IT HAPPEN THAT THEY CAN AND HE'S BEEN INVOLVED WITH GROUPS OF OTHER MINORITIES

WHAT YOU NEED

Requirements

Audio files (.wav)

Transcription files (.lab)

Configuration file (.yaml)

- Not mentioned in tutorials for older versions
- Contains settings and a “list of phones”
- English example included in download

Example

```
1 # for human reading only
2 authors: Kyle Gorman
3 language: English
4 citation: "K. Gorman, J. Howell, and M. Wagner. 2011. Prosodylab-Aligner: A tool for forced alignment of laboratory speech. Canadian Acoustics, 39(3), 192-193."
5 URL: http://prosodylab.org/tools/aligner/
6
7 # basic features
8 samplerate: 16000 # in Hz
9 phoneset: [AA0, AA1, AA2, AB0, AB1, AB2, AB3, AB4, AB5, AB6, AB7, AB8, AB9, AB0, AB1, AB2,
10           AB3, AB4, AB5, AB6, AB7, AB8, AB9, AB0, AB1, AB2, AB3, AB4, AB5, AB6, AB7, AB8, AB9,
11           EY0, EY1, EY2, IH0, IH1, IH2, IY0, IY1, IY2, OW0, OW1, OW2,
12           OY0, OY1, OY2, UH0, UH1, UH2, UW0, UW1, UW2,
13           B, C, D, DG, F, G, HH, TH, K, L, M, N, NG, P, R,
14           S, SH, T, TH, V, W, Y, Z, ZH]
15
16 # specs for feature extractor: change at your own risk
17 HCopy:
18   SOURCEKIND: WAVEFORM
19   SOURCEFORMAT: RAW
20   TARGETRATE: 16000.0
21   TARGETKIND: MFCC_D_A_0
22   WINDOWSIZE: 25600.0
23   PREEMPH: 0.95
24   USERAMING: T
25   ENORMALIZE: T
26   CEPLIFTER: 12
27   NUMKANS: 50
28   NUMKEPS: 12
29
30 # pruning parameters, to use globally: change at your own risk
31 pruning: [256, 160, 5000]
32
33 # specs for flat start: change at your own risk
34 HCompV:
35   F: .01
36
37 # specs for estimation: change at your own risk
38 HEst:
39   TARGETRATE: 16000.0
40   TARGETKIND: MFCC_D_A_0
41   WINDOWSIZE: 25600.0
42   PREEMPH: 0.95
43   USERAMING: T
44   ENORMALIZE: T
45   CEPLIFTER: 12
46   NUMKANS: 50
47   NUMKEPS: 12
48
49 # specs for the decoder: change at your own risk
50 HVite:
51   SFAC: 5
52
```

```
1 # for human reading only
2 authors: Kyle Gorman
3 language: English
4 citation: "K. Gorman, J. Howell, and M. Wagner. 2011. Prosodylab-Aligner: A tool for forced alignment of laboratory speech. Canadian Acoustics, 39(3), 192-193."
5 URL: http://prosodylab.org/tools/aligner/
6
7 # basic features
8 samplerate: 16000 # in Hz
9 phoneset: [AA0, AA1, AA2, AE0, AE1, AE2, AH0, AH1, AH2, AO0, AO1, AO2,
10           AW0, AW1, AW2, AY0, AY1, AY2, EH0, EH1, EH2, ER0, ER1, ER2,
11           EY0, EY1, EY2, IH0, IH1, IH2, IY0, IY1, IY2, OW0, OW1, OW2,
12           OY0, OY1, OY2, UH0, UH1, UH2, UW0, UW1, UW2,
13           B, CH, D, DH, F, G, HH, JH, K, L, M, N, NG, P, R,
14           S, SH, T, TH, V, W, Y, Z, ZH]
15
16 # specs for feature extractor; change at your own risk
17 HCopy:
18     SOURCEKIND: WAVEFORM
19     SOURCEFORMAT: WAVE
20     TARGETRATE: 100000.0
21     TARGETKIND: MFCC_D_A_0
22     WINDOWSIZE: 250000.0
23     PREEMCOEF: 0.97
24     USEHAMMING: T
25     ENORMALIZE: T
26     CEPLIFTER: 22
27     NUMCHANS: 20
28     NUMCEPS: 12
29
30 # pruning parameters, to use globally; change at your own risk
31 pruning: [250, 100, 5000]
32
33 # specs for flat start; change at your own risk
34 HCompV:
35     F: .01
36
37 # specs for estimation; change at your own risk
38 HERest:
39     TARGETRATE: 100000.0
40     TARGETKIND: MFCC_D_A_0
41     WINDOWSIZE: 250000.0
42     PREEMCOEF: 0.97
43     USEHAMMING: T
44     ENORMALIZE: T
45     CEPLIFTER: 22
46     NUMCHANS: 20
47     NUMCEPS: 12
48
49 # specs for the decoder; change at your own risk
50 HVite:
51     SFAC: 5
52
```

WHAT YOU NEED

Requirements

Audio files (.wav)

Transcription files (.lab)

Configuration file (.yaml)

Dictionary file

- Provides pronunciation
- Uses “phones” listed in .yaml file
- Follows prescribed format
- North American English example included in download (others available)

Example

```
115276 TRANSISTORS T R AE0 N Z IH1 S T ER0 Z
115277 TRANSIT T R AE1 N Z IH0 T
115278 TRANSITED T R AE1 N Z IH0 T IH0 D
115279 TRANSITING T R AE1 N Z IH0 T IH0 NG
115280 TRANSITION T R AE0 N Z IH1 SH AH0 N
115281 TRANSITIONAL T R AE0 N S IH1 SH AH0 N AH0 L
115282 TRANSITIONAL T R AE0 N Z IH1 SH AH0 N AH0 L
115283 TRANSITIONING T R AE0 N Z IH1 SH AH0 N IH0 NG
115284 TRANSITIONS T R AE0 N Z IH1 SH AH0 N Z
115285 TRANSITORY T R AE1 N Z AH0 T AO2 R IY0
115286 TRANSITS T R AE1 N Z IH0 T S
115287 TRANSKEI T R AE1 N Z K EY2
115288 TRANSLATE T R AE0 N S L EY1 T
115289 TRANSLATE T R AE0 N Z L EY1 T
115290 TRANSLATED T R AE0 N S L EY1 T IH0 D
115291 TRANSLATED T R AE0 N Z L EY1 T AH0 D
115292 TRANSLATES T R AE0 N Z L EY1 T S
115293 TRANSLATES T R AE1 N S L EY2 T S
115294 TRANSLATING T R AE0 N Z L EY1 T IH0 NG
115295 TRANSLATING T R AE1 N S L EY2 T IH0 NG
115296 TRANSLATION T R AE0 N S L EY1 SH AH0 N
115297 TRANSLATION T R AE0 N Z L EY1 SH AH0 N
115298 TRANSLATIONS T R AE0 N S L EY1 SH AH0 N Z
```


WHAT YOU NEED



What we used

Audio files (.wav)

- Word list readings
- Collected in the field
- Recorded with lavalier mics and Zoom H4n digital recorder
- 16 bit, 44.1 kHz (did not resample)
- Some files “cleaned” in Praat (22 files, 1:41:30), others left “dirty” with only extraneous speech removed in Audacity (16 files, 2:39:44 + 5 “dirty” versions of clean files 1:23:01)

WHAT YOU NEED

What we used

Audio files (.wav)

Transcription files (.lab)

- Originally created in Elan using controlled vocabulary
- Transcriptions of “clean” files: extracted non-empty intervals, and concatenated in Praat, then exported and formatted in Word and Notepad++
- Transcriptions of “dirty” files exported from Elan and prepared in Word and Notepad++.
- Used Tongan orthography (with ? instead of ‘)

1 KOTOA KOTOA PEA MANU EFUEFU TUJA KOWI KILIPI KOE?UHI KETE MANUPUNA UPU ?ULI?ULI TOTO ANGI HUI HUHU MANAVA TUTU TAMASIPI ?AO MOMOKO HAU TONU LAU TUFUSI MATE KELI ?ULI KULI INU MOMOA PEKU EFU TELINGA FONUA MAHA MATA TO MAMA?O TAMAI NGAKO MANAVARE SIPI NGE?ESI NIMA AFI IKA NIMA TETE LA?I ?AKAU TAFE PUNA KAKAPU VAZE VAO POTO FOTI ?AKAU FONU FOAKI LELEI MUSIE LANUMATA NGAKAU LOU ?ULU NIMA ?ULU FANONGO MAPU MAMAPA IA TA SEU FEFE TULI MANU HOA TANGATA AU ?AISI KAPAU LOTO TAMATE?I TUI ?ILO ANO LAHI KATA LAU HEMA VAZE TOKOTO M?UUI ?ATE L?LOA KUTU TANGATA LAHI KAKANO MAHINA FA?E MO?UNGA NGUTU HINGOA LAUSI?I OFI KIA FO?OU PO IHU MOTU?A TASA TOKO TASA VA?INGA FUSI TEKE ?UBA KULOKULA MATAPU VAI TAFE BALA AKA MASA PALA FUOPOTOPOTO OLO MASIMA ?ONE?ONE PEHE ?AHI TENGA TUITUI MASILA NOUNOU HIVA TANGUTU KILI LANGI MORE NAMU?I ?AHU MOLEMOLE NGATA SINGU HA ?APANU FAHI HOKA TUTU FETUTU TOKOTOKO TOTONU MISI LA?A ISU KINAUTOLU MATOLU FAKAKAUKAU MANIFI KIMCUTOLU KOE TOLU LI NONO?O ?ELELO NIFO FUFU ?AKAU FULIHI UA LUA VELA VIRU HA ?APE FE HINEHINA HAI HA LAULAHU HOA FEFINE MATANGI KAPAKAU KAPAKAU KELEMUTU HALA ?APU ENGEENGA KOTOA PEA MANU EFUEFU TUJA KOWI KILIPI KOE?UHI KETE MANUPUNA UPU ?ULI?ULI TOTO ANGI HUI HUHU MANAVA TUTU TAMASIPI ?AO MOMOKO HAU TONU LAU TUFUSI MATE KELI ?ULI KULI INU MOMOA PEKU EFU TELINGA FONUA MAHA MATA TO MAMA?O TAMAI NGAKO MANAVARE SIPI NGE?ESI NIMA AFI IKA NIMA TETE LA?I ?AKAU TAFE PUNA KAKAPU VAZE VAO POTO FOTI ?AKAU FONU FOAKI LELEI MUSIE LANUMATA NGAKAU LOU ?ULU NIMA ?ULU FANONGO MAPU MAMAPA IA TA SEU FEFE TULI MANU HOA TANGATA AU ?AISI KAPAU LOTO TAMATE?I TUI ?ILO ANO LAHI KATA LAU HEMA LAU HEMA VAZE TOKOTO M?UUI ?ATE L?LOA KUTU TANGATA LAHI KAKANO MAHINA FA?E MO?UNGA NGUTU HINGOA LAUSI?I OFI KIA FO?OU PO IHU MOTU?A TASA TOKO TASA VA?INGA FUSI TEKE ?UBA KULOKULA MATAPU VAI TAFE BALA AKA MASA PALA FUOPOTOPOTO OLO MASIMA ?ONE?ONE PEHE ?AHI TENGA TUITUI MASILA NOUNOU HIVA TANGUTU KILI LANGI MORE NAMU?I ?AHU MOLEMOLE NGATA SINGU HA ?APANU FAHI HOKA TUTU FETUTU TOKOTOKO TOTONU MISI LA?A ISU KINAUTOLU MATOLU FAKAKAUKAU MANIFI KIMCUTOLU KOE TOLU LI NONO?O ?ELELO NIFO FUFU ?AKAU FULIHI UA LUA VELA VIRU HA ?APE FE HINEHINA HAI HA LAULAHU HOA FEFINE MATANGI KAPAKAU KAPAKAU KELEMUTU HALA ?APU ENGEENGA

1 KOTOA KOTOA PEA MANU EFUEFU TU?A KOVI KILI?I KOE?UHI KETE MANUPUNA U?U ?ULI?ULI TOTO ANGI HUI HUHU I
INU MŌMOA PEKU EFU TELINGA FONUA MAHA MATA TŌ MAMA?O TAMAI NGAKO MANAVAHĒ SI?I NGE?ESI NIMA AFI IKA
FOAKI LELEI MUSIE LANUMATA NGĀKAU LOU ?ULU NIMA ?ULU FANONGO MAFU MAMAFI IA TĀ SEU FĒFĒ TULI MANU HO
VA?E TOKOTO MO?UI ?ATE LŌLOA KUTU TANGATA LAHI KAKANO MĀHINA FA?Ē MO?UNGA NGUTU HINGOA LAUSI?I OFI I
MATA?U VAI TAFE HALA AKA MAEA PALA FUOPOTOPOTO OLO MĀSIMA ?ONE?ONE PEHĒ TAHI TENGA TUITUI MĀSILA NOU
?A?ANU FAHI HOKA TU?U FETU?U TOKOTOKO TOTONU MISI LA?Ā IKU KINAUTOLU MATOLU FAKAKAUKAU MANIFI KIMOU
?AFĒ FĒ HINEHINA HAI HĀ LAULAHĒ HOA FEFINE MATANGI KAPAKAU FEFINE KELEMUTU HALA TA?U ENGEENGA KOTOA
ANGI HUI HUHU MĀNAVA TUTU TAMASI?I ?AO MOMOKO HA?U TONU LAU TU?USI MATE KELI ?ULI KULĪ INU MŌMOA PE
NGE?ESI NIMA AFI IKA NIMA TĒTĒ LA?I ?AKAU TAFE PUNA KAKAPU VA?E VAO POTO FO?I ?AKAU FONU FOAKI LELE
SEU FĒFĒ TULI MANU HOA TANGATA AU ?AISI KAPAU LOTO TĀMATE?I TUI ?ILO ANO LAHI KATA LAU HEMA LAU HEM
MO?UNGA NGUTU HINGOA LAUSI?I OFI KIA FO?OU PŌ IHU MOTU?A TAHA TOKO TAHA VA?INGA FUSI TEKE ?UHA KULO
PEHĒ TAHI TENGA TUITUI MĀSILA NOUNOU HIVA TANGUTU KILI LANGI MOHE NĀMU?I ?AHU MOLEMOLE NGATA SINOU I
FAKAKAUKAU MANIFI KIMOUTOLU KOE TOLU LĪ NONO?O ?ELELO NIFO FU?U ?AKAU FULIHI UA LUA VELA VIKU HĀ ?A
KELEMUTU HALA TA?U ENGEENGA KOTOA PEA MANU EFUEFU TU?A KOVI KILI?I KOE?UHI KETE MANUPUNA U?U ?ULI?U
TU?USI MATE KELI ?ULI KULĪ INU MŌMOA PEKU EFU TELINGA FONUA MAHA MATA TŌ MAMA?O TAMAI NGAKO MANAVAHĒ
VAO POTO FO?I ?AKAU FONU FOAKI LELEI MUSIE LANUMATA NGĀKAU LOU ?ULU NIMA ?ULU FANONGO MAFU MAMAFI IA
ANO LAHI KATA LAU HEMA VA?E TOKOTO MO?UI ?ATE LŌLOA KUTU TANGATA LAHI KAKANO MĀHINA FA?Ē MO?UNGA NGU
FUSI TEKE ?UHA KULOKULA MATA?U HALA AKA MAEA PALA FUOPOTOPOTO OLO MĀSIMA ?ONE?ONE PEHĒ TAHI TENGA T
NGATA SINOU HA ?A?ANU FAHI HOKA TU?U FETU?U TOKOTOKO TOKOTOKO TOTONU MISI LA?Ā IKU KINAUTOLU MATOLU
FULIHI UA LUA VELA VIKU HĀ ?AFĒ FĒ HINEHINA HAI HĀ LAULAHĒ HOA FEFINE MATANGI KAPAKAU KAPAKAU KELEM

WHAT YOU NEED

```
1 # for human reading only
2 authors: Kyle Gorman
3 language: Tongan
4 citation: "K. Gorman, J. Howell, and M. Wagner. 2011. Prosodylab-Aligner: A tool
5 URL: http://prosodylab.org/tools/aligner/
6
7 # basic features
8 # samplerate: 16000 # in Hz
9 # CJ - modified to 44100 to match the recordings and avoid a downsample
10 samplerate: 44100 # in Hz
11
12 phoneset: [a, ă, e, ê, i, î, o, ô, u, û, f, h, k, l, m, n, ng, p, s, t, v, ?]
13
14 # specs for feature extractor; change at your own risk
15 HCopy:
16   SOURCEKIND: WAVEFORM
17   SOURCEFORMAT: WAVE
18   TARGETRATE: 100000.0
19   TARGETKIND: MFCC_D_A_0
20   WINDOWSIZE: 250000.0
21   PREEMCOEF: 0.97
22   USEHAMMING: T
23   ENORMALIZE: T
24   CEPLIFTER: 22
25   NUMCHANS: 20
26   NUMCEPS: 12
27
28 # pruning parameters, to use globally; change at your own risk
29 pruning: [250, 100, 5000]
30
31 # specs for flat start; change at your own risk
32 HCompV:
33   F: .01
34
35 # specs for estimation; change at your own risk
36 HEst:
37   TARGETRATE: 100000.0
38   TARGETKIND: MFCC_D_A_0
39   WINDOWSIZE: 250000.0
40   PREEMCOEF: 0.97
41   USEHAMMING: T
42   ENORMALIZE: T
43   CEPLIFTER: 22
44   NUMCHANS: 20
45   NUMCEPS: 12
46
47 # specs for the decoder; change at your own risk
48 HVite:
49   SFAC: 5
50
51
```

What we used

Audio files (.wav)

Transcription files (.lab)

Configuration file (.yaml)

- “Phone list” uses 1 digraph (ng) and one Unicode IPA character (?)
- Changed 1 setting (targetrate) to prevent crash.

```
1 # for human reading only
2 authors: Kyle Gorman
3 language: Tongan
4 citation: "K. Gorman, J. Howell, and M. Wagner. 2011. Prosodylab-Aligner: A tool for
5 URL: http://prosodylab.org/tools/aligner/
6
7 # basic features
8 # samplerate: 16000 # in Hz
9 # CJ - modified to 44100 to match the recordings and avoid a downsample
10 samplerate: 44100 # in Hz
11
12 phoneset: [a, ā, e, ē, i, ī, o, ō, u, ū, f, h, k, l, m, n, ng, p, s, t, v, ?]
13
14
15 # specs for feature extractor; change at your own risk
16 HCopy:
17     SOURCEKIND: WAVEFORM
18     SOURCEFORMAT: WAVE
19     TARGETRATE: 100000.0
20     TARGETKIND: MFCC_D_A_0
21     WINDOWSIZE: 250000.0
22     PREEMCOEF: 0.97
23     USEHAMMING: T
24     ENORMALIZE: T
25     CEPLIFTER: 22
26     NUMCHANS: 20
27     NUMCEPS: 12
28
29 # pruning parameters, to use globally; change at your own risk
30 pruning: [250, 100, 5000]
```

WHAT YOU NEED

```
1 AFI a f i
2 AKA a k a
3 ANGI a n g i
4 ANO a n o
5 AU a u
6 EFU e f u
7 EFUEFU e f u e f u
8 ENGEENGA e n g e e n g a
9 FAHI f a h i
10 FAKAKAUKAU f a k a k a u k a u
11 FANONGO f a n o n g o
12 FAʔĒ f a ʔ ē
13 FEFINE f e f i n e
14 FETUʔU f e t u ʔ u
15 FĒ f ē
16 FĒFĒ f ē f ē
17 FOAKI f o a k i
18 FONU f o n u
19 FONUA f o n u a
20 FOʔI f o ʔ i
21 FOʔOU f o ʔ o u
22 FULIHI f u l i h i
23 FUOPOTOPOTO f u o p o t o p o t o
24 FUSI f u s i
25 FUʔU f u ʔ u
26 HA h a
27 HAI h a i
28 HALA h a l a
29 HAʔU h a ʔ u
30 HĀ h ā
31 HEMA h e m a
32 HINEHINA h i n e h i n a
33 HINGOA h i n g o a
34 HIVA h i v a
35 HOA h o a
36 HOKA h o k a
37 HUH U h u h u
38 HUI h u i
39 IA i a
40 IHU i h u
```

What we used

Audio files (.wav)

Transcription files (.lab)

Configuration file (.yaml)

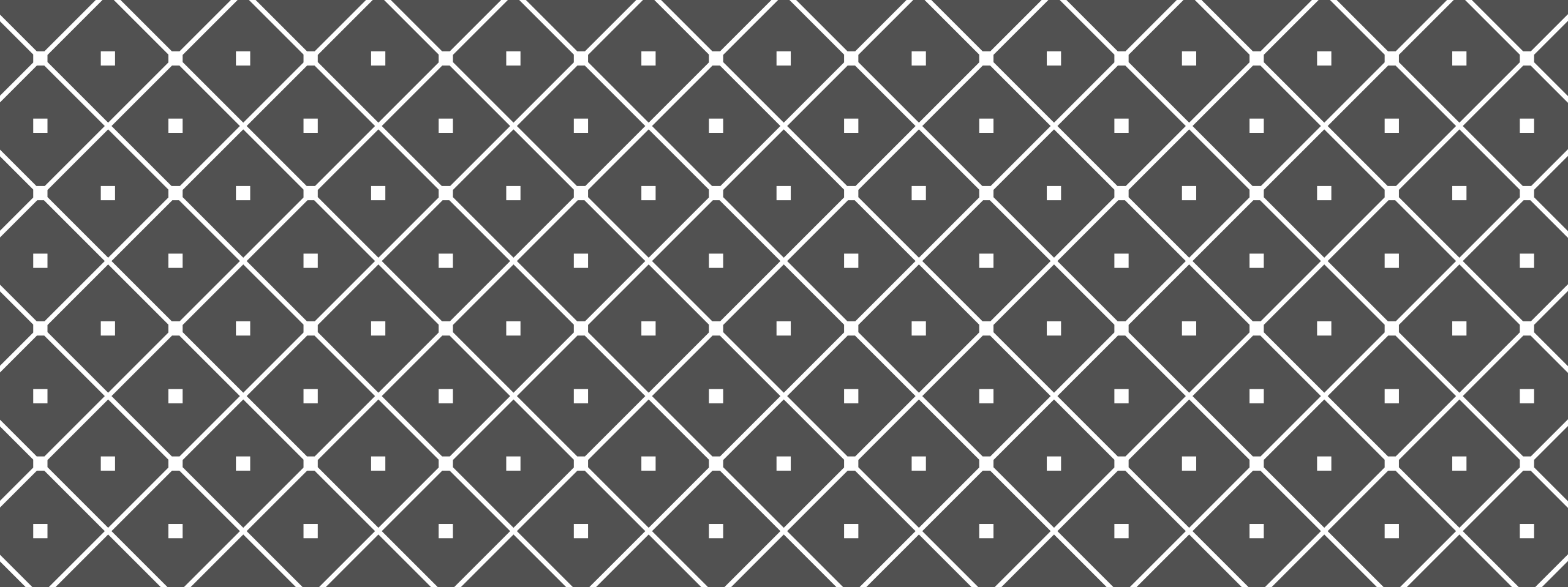
Dictionary file

- Created for this project from word list.
- Pronunciations based on orthography
- No alternate pronunciations included

1	AFI	a	f	i
2	AKA	a	k	a
3	ANGI	a	ng	i
4	ANO	a	n	o
5	AU	a	u	
6	EFU	e	f	u
7	EFUEFU	e	f	u e f u
8	ENGEENGA	e	ng	e e ng a
9	FAHI	f	a	h i
10	FAKAKAUKAU	f	a	k a k a u k a u
11	FANONGO	f	a	n o ng o
12	FAʔĒ	f	a	ʔ ē
13	FEFINE	f	e	f i n e
14	FETUʔU	f	e	t u ʔ u
15	FĒ	f	ē	
16	FĒFĒ	f	ē	f ē
17	FOAKI	f	o	a k i
18	FONU	f	o	n u
19	FONUA	f	o	n u a
20	FOʔI	f	o	ʔ i
21	FOʔOU	f	o	ʔ o u
22	FULIHI	f	u	l i h i
23	FUOPOTOPOTO	f	u	o p o t o p o t o

ISSUES AND SOLUTIONS

- Files must be saved as UTF-8 without “byte order mark” (BOM or “signature”)
- May need to check for extra spaces and carriage returns at the end of the text file
- Dictionary file must be sorted in Python’s sort order (script included)
- Apostrophes can be problematic
- May need to check for hidden .txt extensions



TONGAN TESTS

Training and Alignment

TRAINING

- Produces a acoustic model by which alignments can be created.
- Requires pairs of audio (.wav) files and transcription (.lab) files in the same folder.
- Is accomplished in three cycles, with a set number of iterations (“epochs”) in each cycle.
- Is executed by entering a Python script (command line) into Terminal (Mac) or Command Prompt (PC).

Key elements of command line:

- -c lang.yaml (configuration file path)
- -d lang.dict (dictionary file path)
- -e 5 (number of epochs)
- -t lang/ (path of folder containing training data)
- -w lang-mod.zip (zip file to which model will be written)

ISSUES AND SOLUTIONS

- Problems can be difficult to diagnose and resolve.
- Problems with script syntax or file prep/organization cause process to fail.
- It can be hard to determine which files contain out-of-dictionary words.
- Some HTK error codes are not included in the PL-A docs or HTK Book. (“ERROR [+7390] StepAlpha: Alpha prune failed”)
- Added a few lines of diagnostics to the code.
- Follow instructions carefully.
- Added code to include this information in the output. (We can make this available.)
- Had to Google the error to see how others had solved the problem. increased Targetrate setting in configuration file from 100000 (default) to 125000; feature measurements extracted every 12.5 ms rather than every 10 ms.

TRAINING TESTS

Test ID #	Type and Number of Audio Files	# of Epochs	Targetrate	Name of Acoustic Model Created	Runtime
TonT001	clean (22 files)	5	100000	ton-001-mod.zip	1:04:43
TonT002	clean (22 files)	10	100000	ton-002-mod.zip	0:28:45
TonT003	clean (22 files)	15	100000	ton-003-mod.zip	1:00:49
TonT004	dirty (16 files)	5	125000	ton-004-mod.zip	1:11:05
TonT005	clean & dirty (38 files)	5	125000	ton-005-mod.zip	1:44:00
TonT006	clean (22 files)	5	125000	ton-006-mod.zip	0:17:52
TonT010	clean (17 files)	5	100000	ton-010-mod.zip	0:16:00
TonT011	dirty (11 files)	5	125000	ton-011-mod.zip	0:18:00

ALIGNMENT

- Produces aligned TextGrids based on a previously created acoustic model.
- Requires pairs of audio (.wav) files and transcription (.lab) files in the same folder.
- Is executed by entering a Python script (command line) into Terminal (Mac) or Command Prompt (PC).

Key elements of command line:

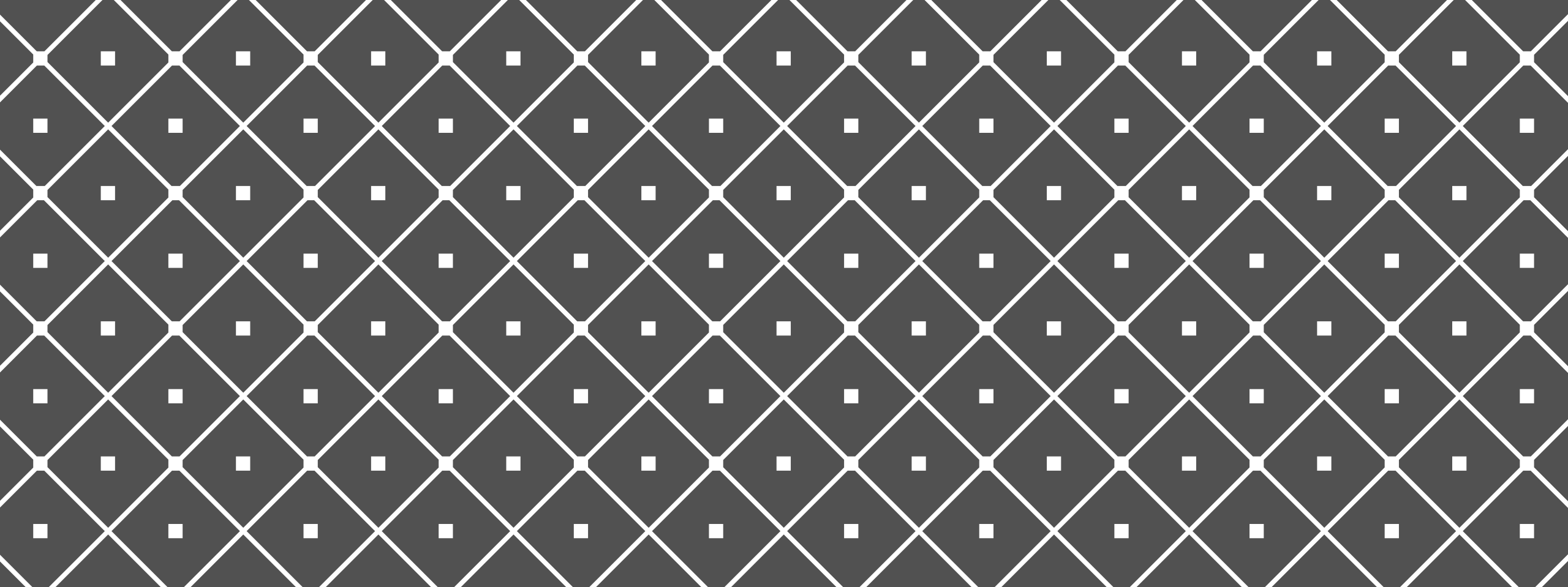
- -r lang-mod.zip ('read': path to language model)
- -a data/ ('align': directory containing files to be aligned)
- -d lang.dict (dictionary file path)

ISSUES AND SOLUTIONS

- Program produces no output to show progress through the process.
- Problems with script syntax or file prep/organization cause process to fail.
- Unicode characters display properly in word tier of output TextGrid but as number codes in phone tier.
- Used Task Manager (processes tab) to monitor process.
- Follow instructions carefully
- Can search and replace in TextGrid, but the characters are unique consistent so the intended Unicode character is clear.

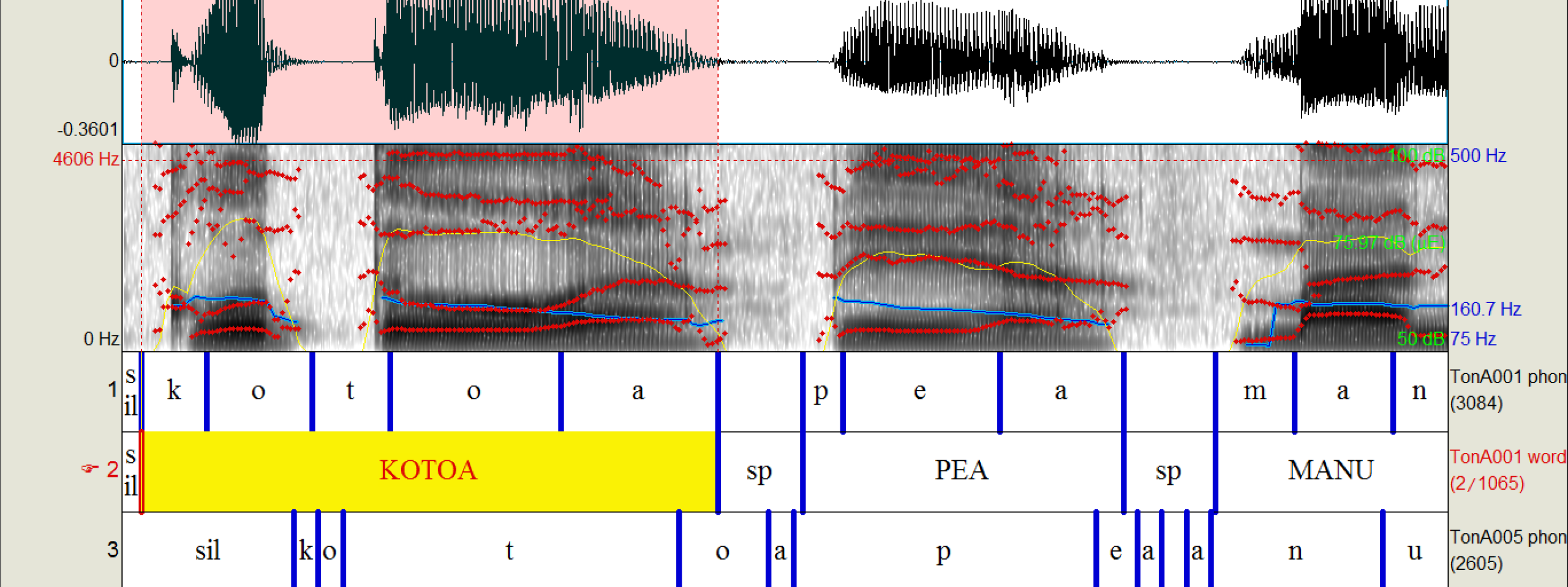
ALIGNMENT TESTS

Test ID #	Type and Number of Aligned Files	Acoustic Model Used in Alignment (and Type of Training Files)	# of Epochs	Targetrate	Runtime
TonA001	clean (22)	ton-001-mod.zip (trained on clean)	5	100000	0:13:19
TonA002	clean (22)	ton-002-mod.zip (trained on clean)	10	100000	0:12:45
TonA003	clean (22)	ton-003-mod.zip (trained on clean)	15	100000	0:20:20
TonA004	dirty (16)	ton-001-mod.zip (trained on clean)	5	100000	0:36:58
TonA005	clean & dirty (38)	ton-004-mod.zip (trained on dirty)	5	125000	0:30:45
TonA006	clean & dirty (38)	ton-005-mod.zip (trained on clean & dirty)	5	125000	0:51:50
TonA007	dirty (16)	ton-002-mod.zip (trained on clean)	10	100000	0:25:50
TonA008	dirty (16)	ton-003-mod.zip (trained on clean)	15	100000	0:26:20
TonA009	dirty (5)	ton-001-mod.zip (trained on clean)	5	100000	0:10:57
TonA010	dirty (5)	ton-002-mod.zip (trained on clean)	10	100000	0:13:06
TonA011	dirty (5)	ton-003-mod.zip (trained on clean)	15	100000	0:13:38
TonA012	dirty (5)	ton-004-mod.zip (trained on dirty)	5	125000	0:14:25
TonA013	dirty (5)	ton-005-mod.zip (trained on clean & dirty)	5	125000	0:16:46
TonA014	clean & dirty (43)	ton-006-mod.zip (trained on clean)	5	125000	0:12:02
TonA017	clean (5)	ton-010-mod.zip (trained on clean)	5	100000	0:04:00
TonA018	dirty (5)	ton-011-mod.zip (trained on dirty)	5	125000	0:05:00



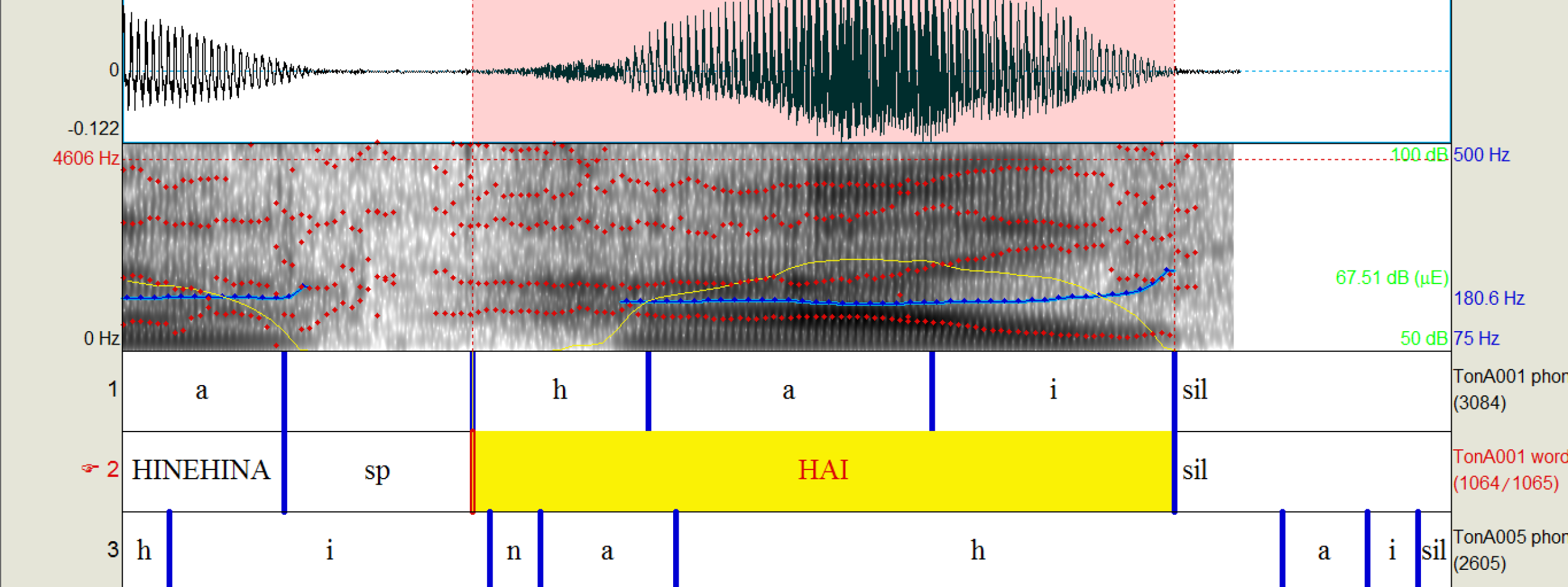
ALIGNMENT COMPARISONS

Reliability and Validity



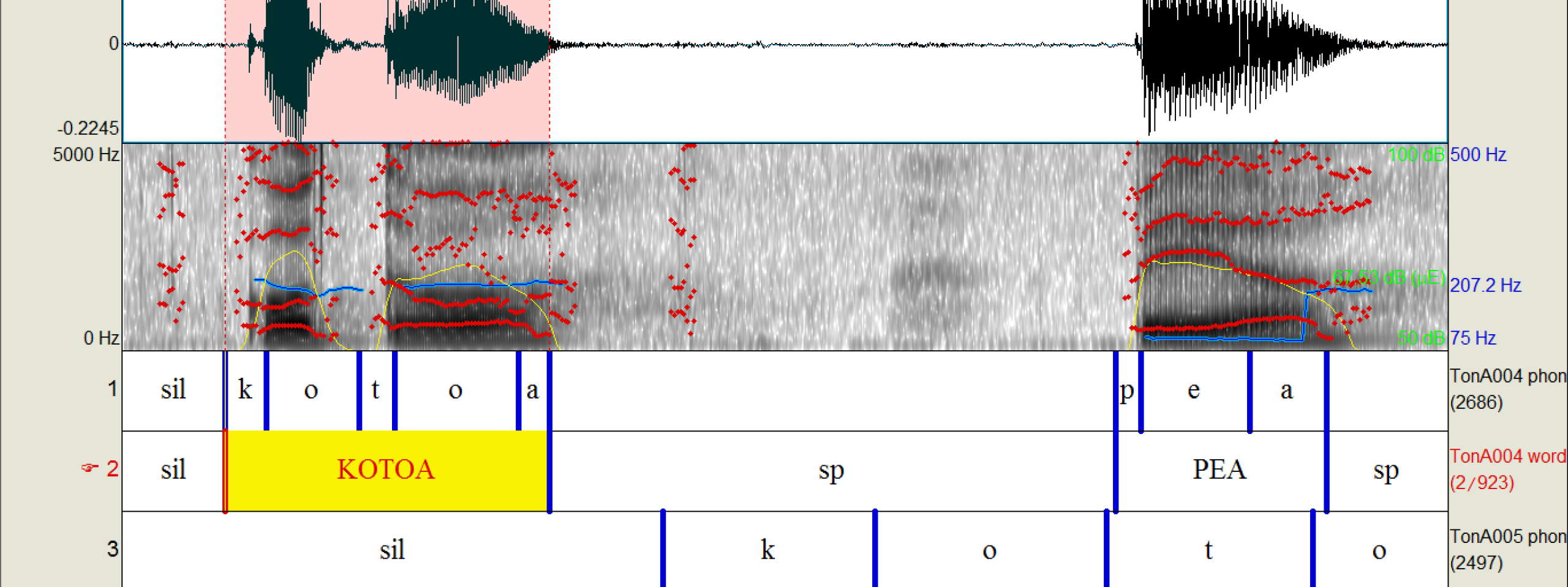
TRAINED ON CLEAN VS. ON DIRTY

TonA001 vs. TonA005
Clean File. Beg. of Recording



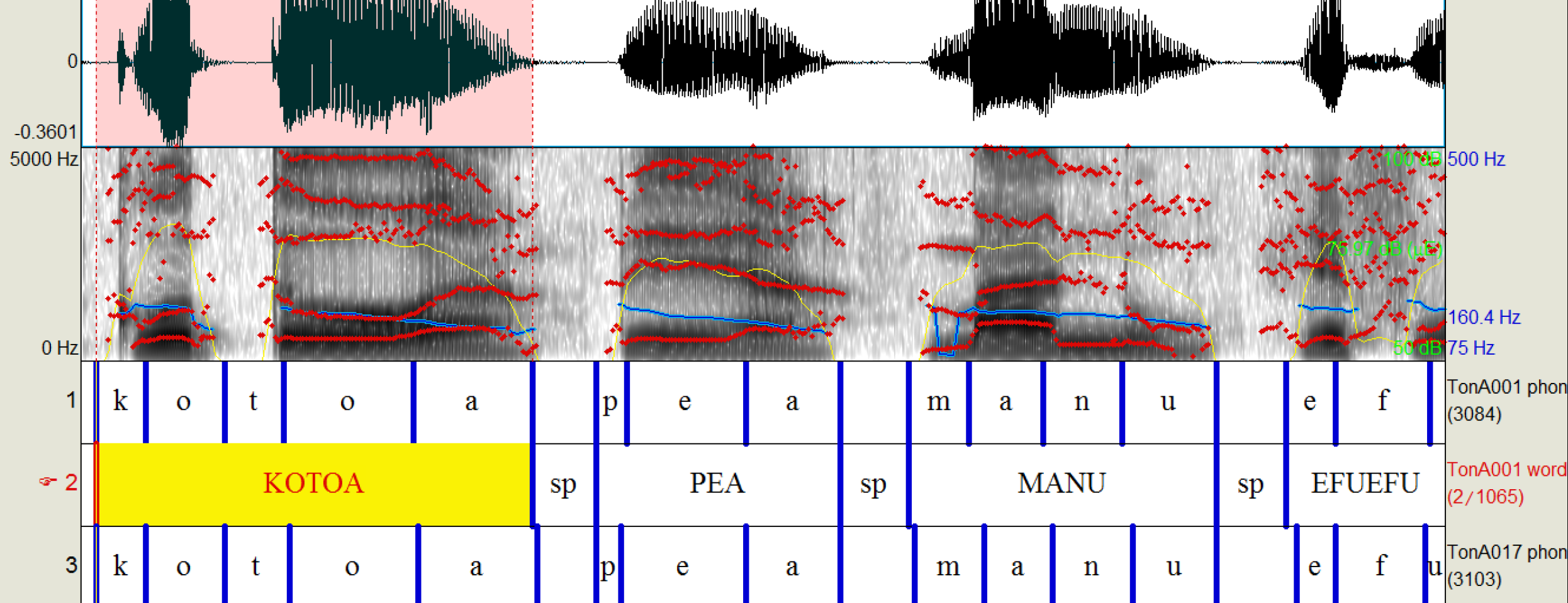
TRAINED ON CLEAN VS. ON DIRTY

TonA001 vs. TonA005
Clean File. End of Recording



TRAINED ON CLEAN VS. ON DIRTY

TonA001 vs. TonA005
Dirty File. Beg. of Recording



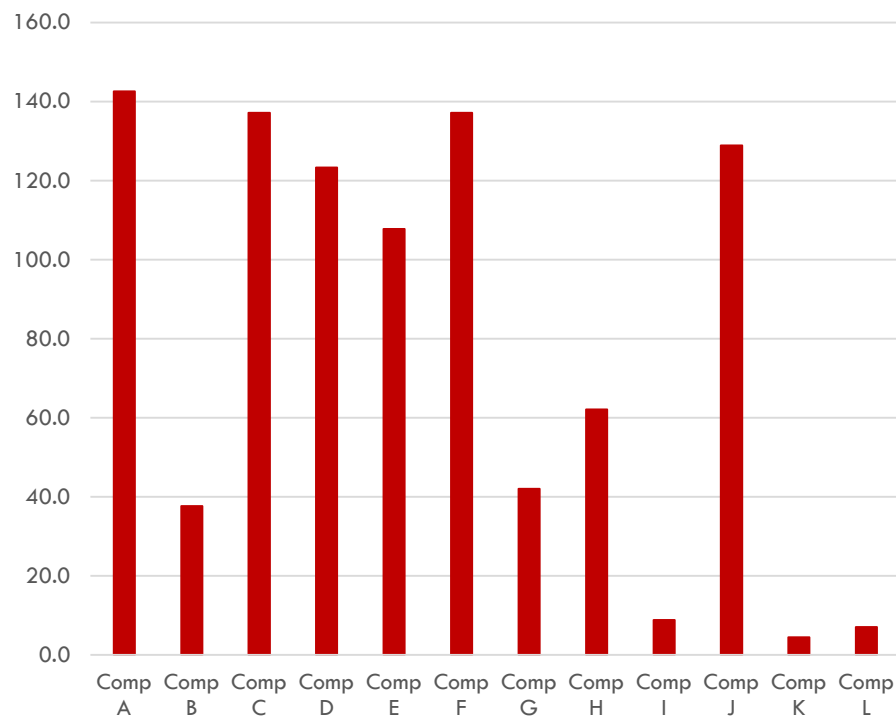
TRAINED ON FILES TO BE ALIGNED?

TonA001: Yes

TonA017: No

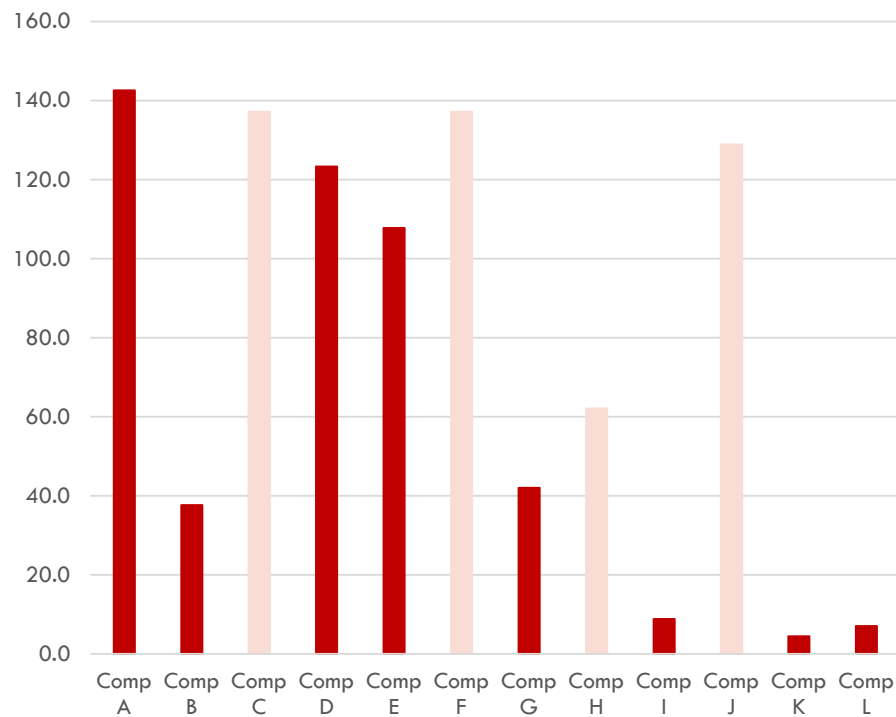
QUANTITATIVE MODEL COMPARISON

Mean Euclidean Distance for Model Comparisons
(F1 & F2, Averaged across All Speakers and Vowels)



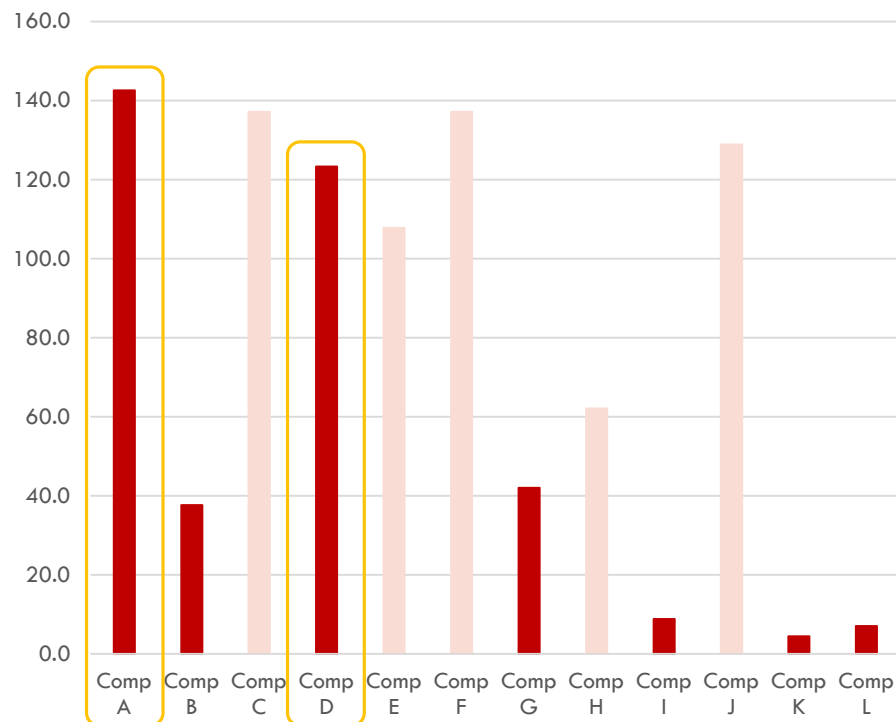
QUANTITATIVE MODEL COMPARISON

Mean Euclidean Distance for Model Comparisons
(F1 & F2, Averaged across All Speakers and Vowels)



QUANTITATIVE MODEL COMPARISON

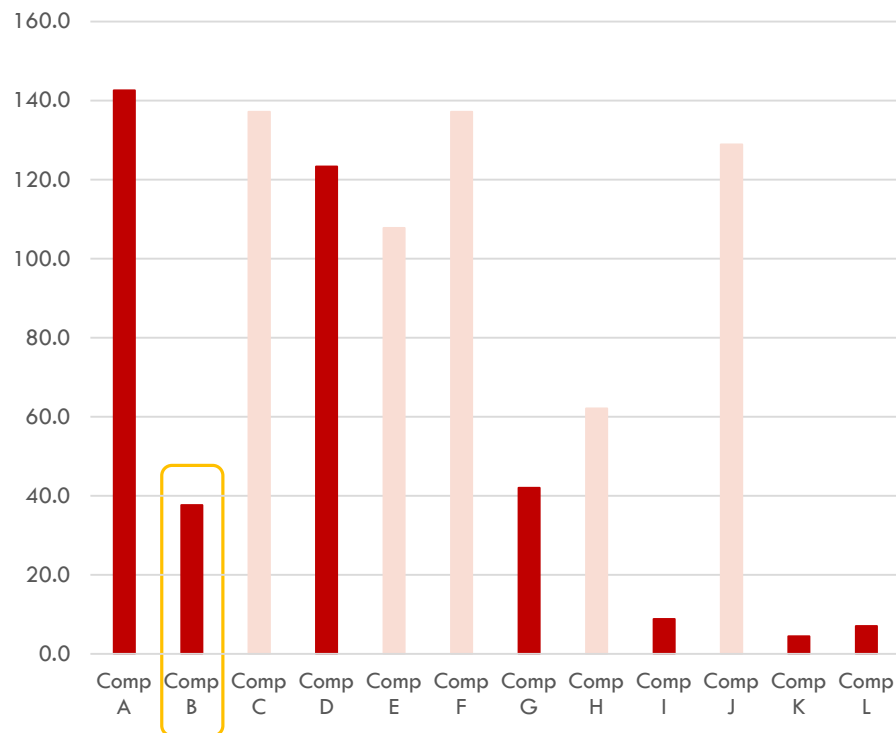
Mean Euclidean Distance for Model Comparisons
(F1 & F2, Averaged across All Speakers and Vowels)



Cleaning up the files used for training the acoustic models had a large effect on the alignments.
(Comparisons A and D)

QUANTITATIVE MODEL COMPARISON

Mean Euclidean Distance for Model Comparisons
(F1 & F2, Averaged across All Speakers and Vowels)

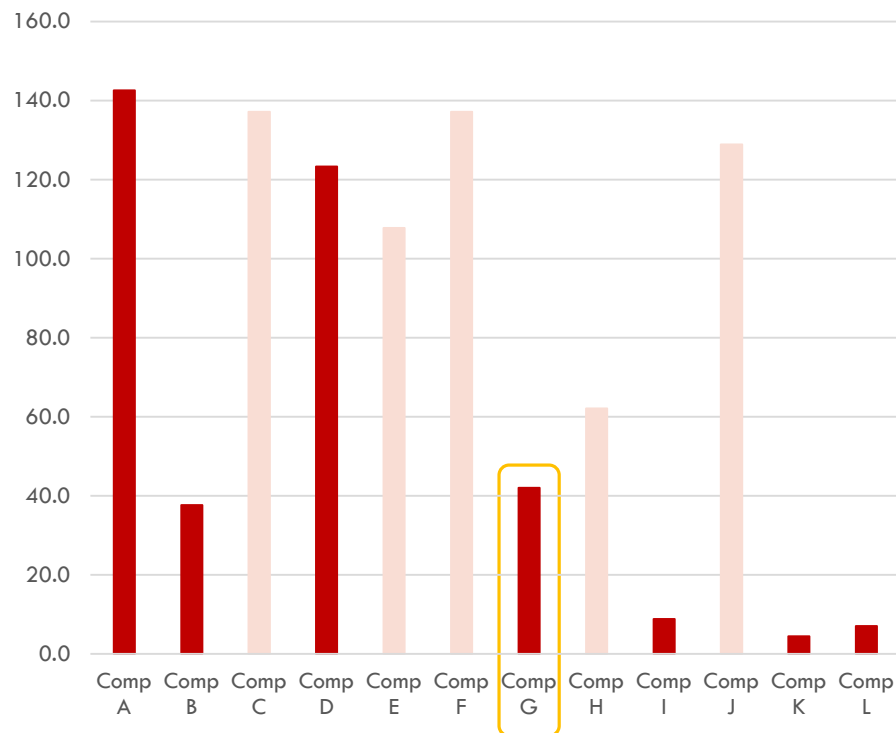


Cleaning up the files used for training the acoustic models had a large effect on the alignments.
(Comparisons A and D)

Including dirty files along with clean files in the training data had a moderate effect on the alignments. (Comparison B)

QUANTITATIVE MODEL COMPARISON

Mean Euclidean Distance for Model Comparisons
(F1 & F2, Averaged across All Speakers and Vowels)



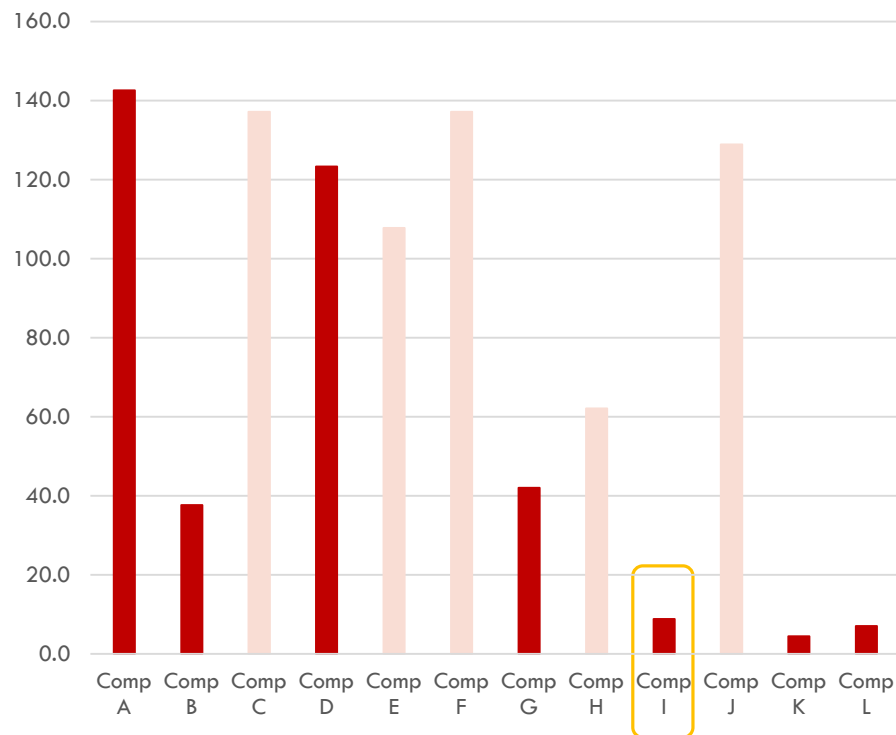
Cleaning up the files used for training the acoustic models had a large effect on the alignments.
(Comparisons A and D)

Including dirty files along with clean files in the training data had a moderate effect on the alignments. (Comparison B)

Changing the Targetrate setting from 100000 to 125000 had some effect on the alignments.
(Comparison G)

QUANTITATIVE MODEL COMPARISON

Mean Euclidean Distance for Model Comparisons
(F1 & F2, Averaged across All Speakers and Vowels)



Cleaning up the files used for training the acoustic models had a large effect on the alignments. (Comparisons A and D)

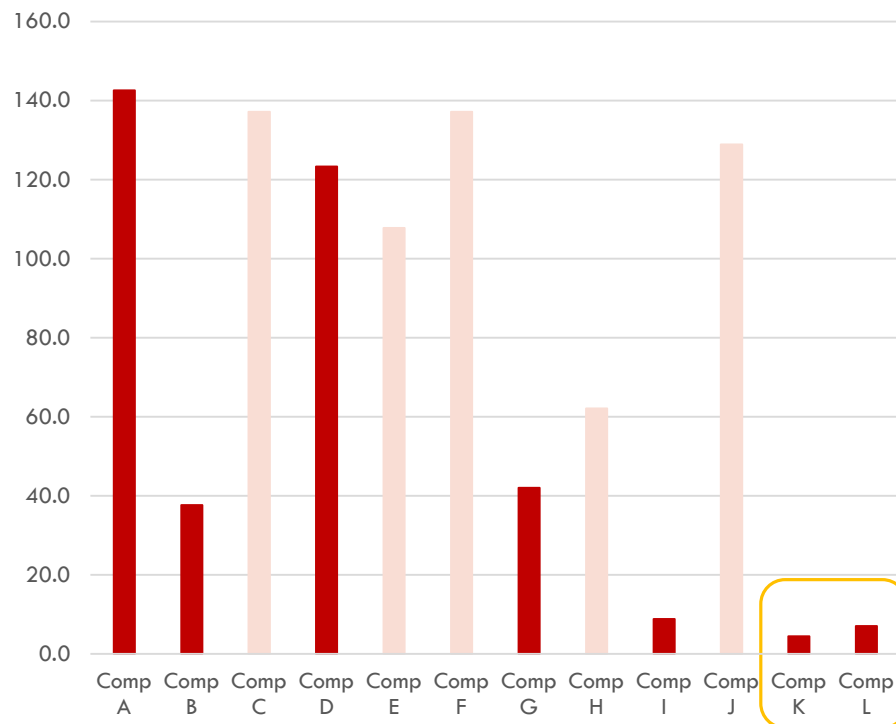
Including dirty files along with clean files in the training data had a moderate effect on the alignments. (Comparison B)

Changing the Targetrate setting from 100000 to 125000 had some effect on the alignments. (Comparison G)

Whether the data used to train the acoustic model included the exact files to be aligned had little effect on the alignments. (Comparison I)

QUANTITATIVE MODEL COMPARISON

Mean Euclidean Distance for Model Comparisons
(F1 & F2, Averaged across All Speakers and Vowels)



Cleaning up the files used for training the acoustic models had a large effect on the alignments. (Comparisons A and D)

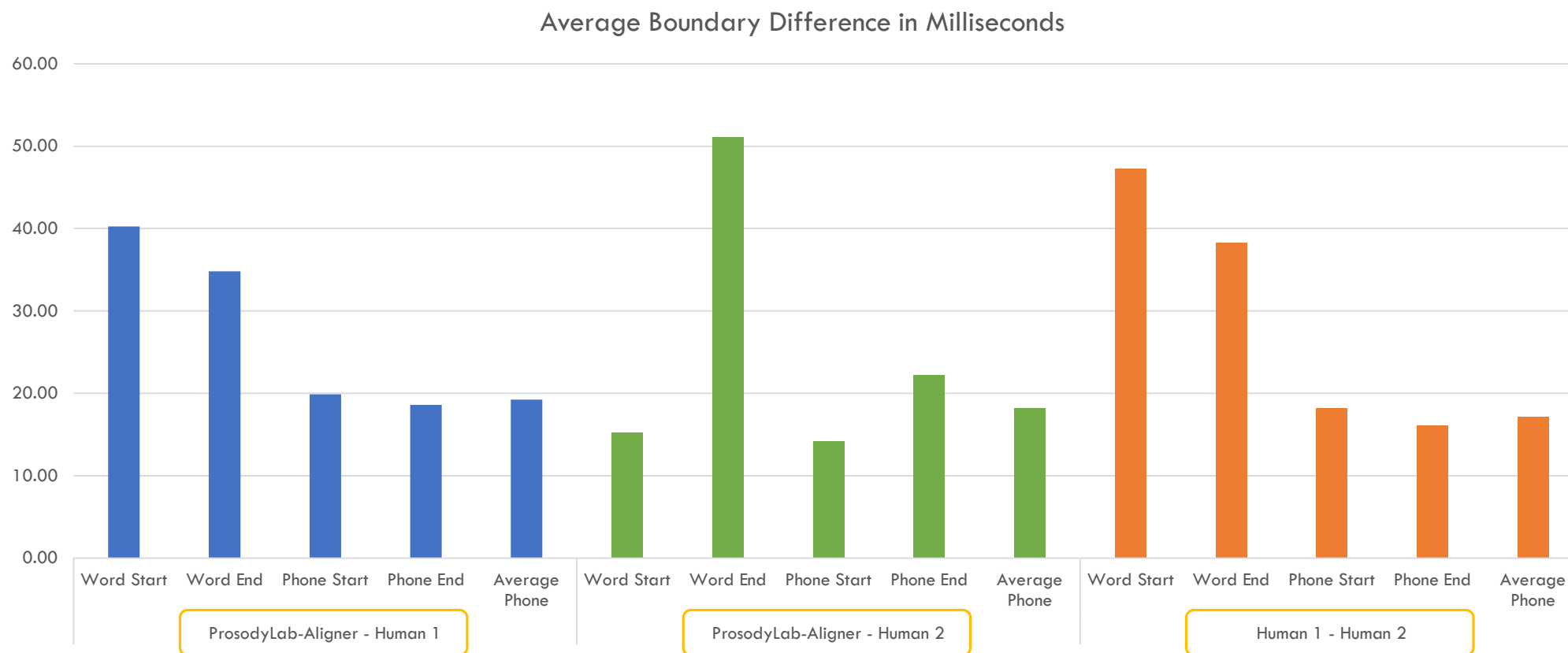
Including dirty files along with clean files in the training data had a moderate effect on the alignments. (Comparison B)

Changing the Targetrate setting from 100000 to 125000 had some effect on the alignments. (Comparison G)

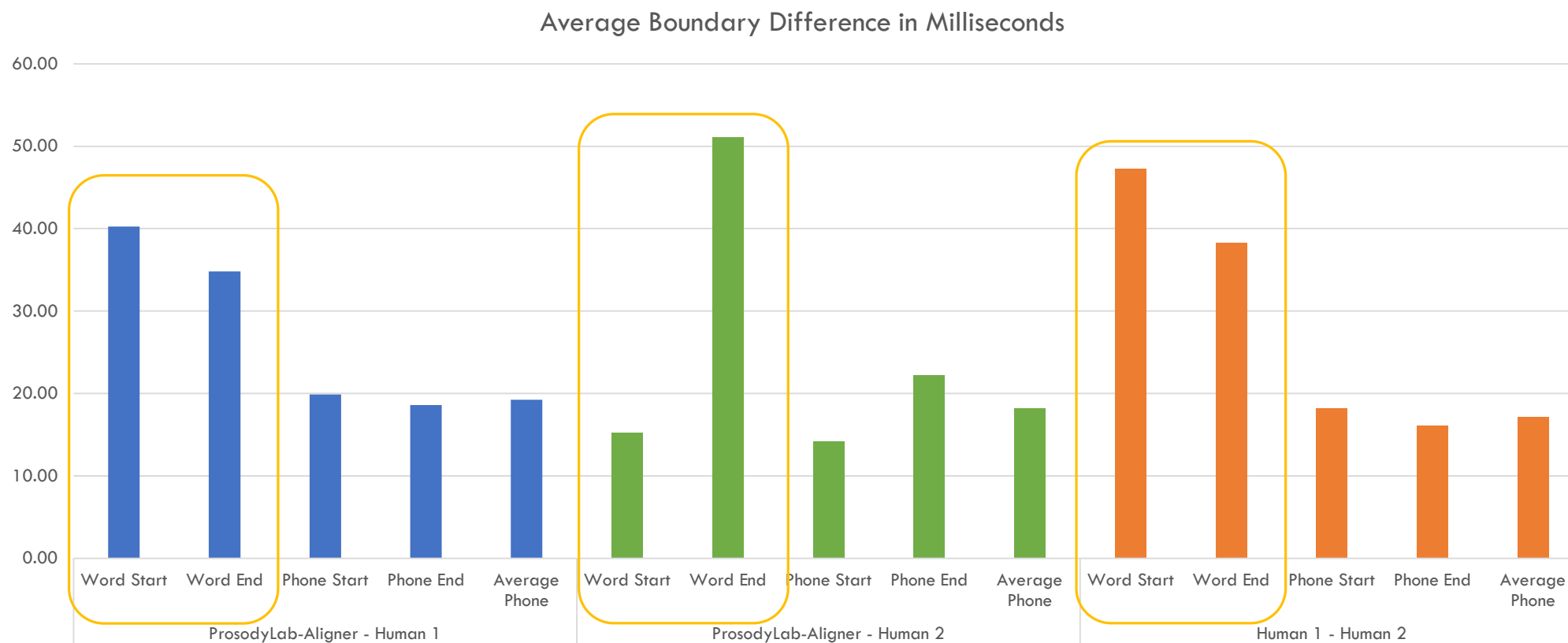
Whether the data used to train the acoustic model included the exact files to be aligned had little effect on the alignments. (Comparison I)

The number of epochs in each cycle of the acoustic model training process had little effect on the final alignments. (Comparisons K and L)

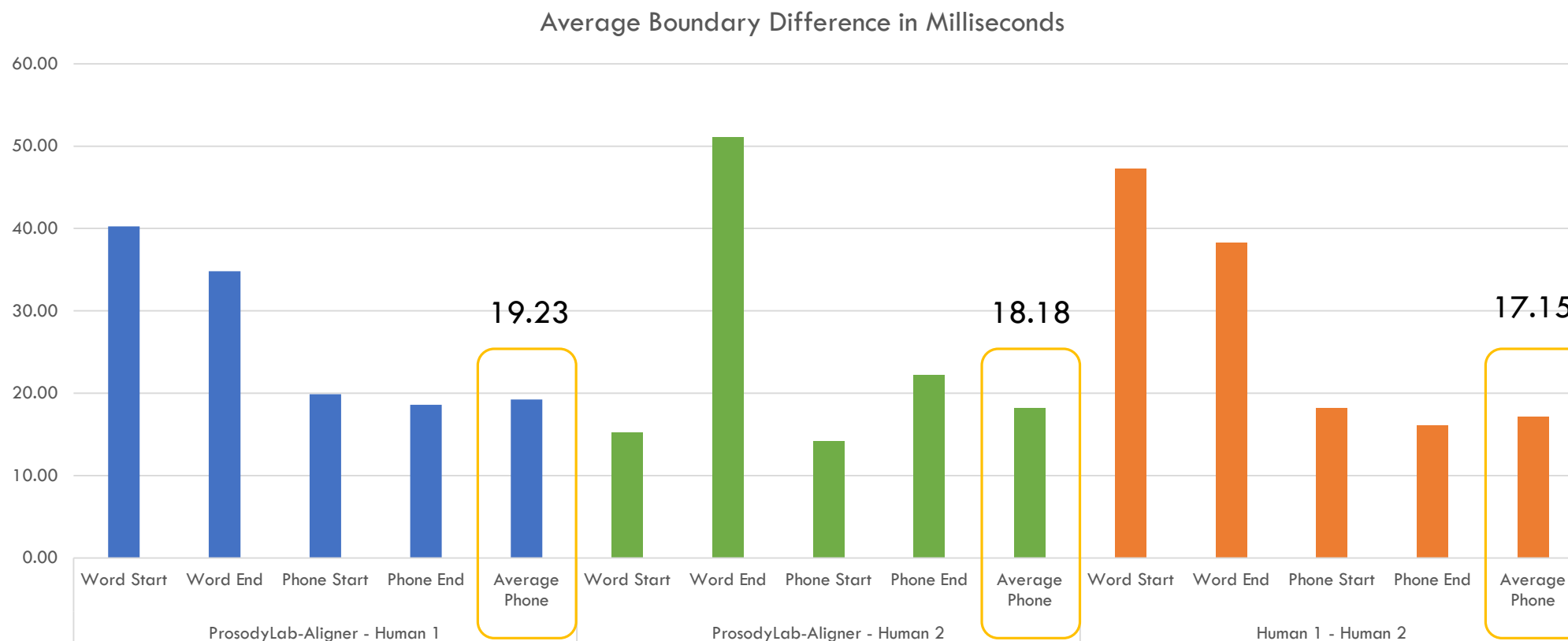
DIFFERENCE BETWEEN PL-A AND HUMAN ALIGNERS



DIFFERENCE BETWEEN PL-A AND HUMAN ALIGNERS



DIFFERENCE BETWEEN PL-A AND HUMAN ALIGNERS



PL-A SUMMARY AND RECOMMENDATIONS

Removing background noise from files used to train acoustic models seems to improve alignments, whether the files to be aligned contain background noise or not.

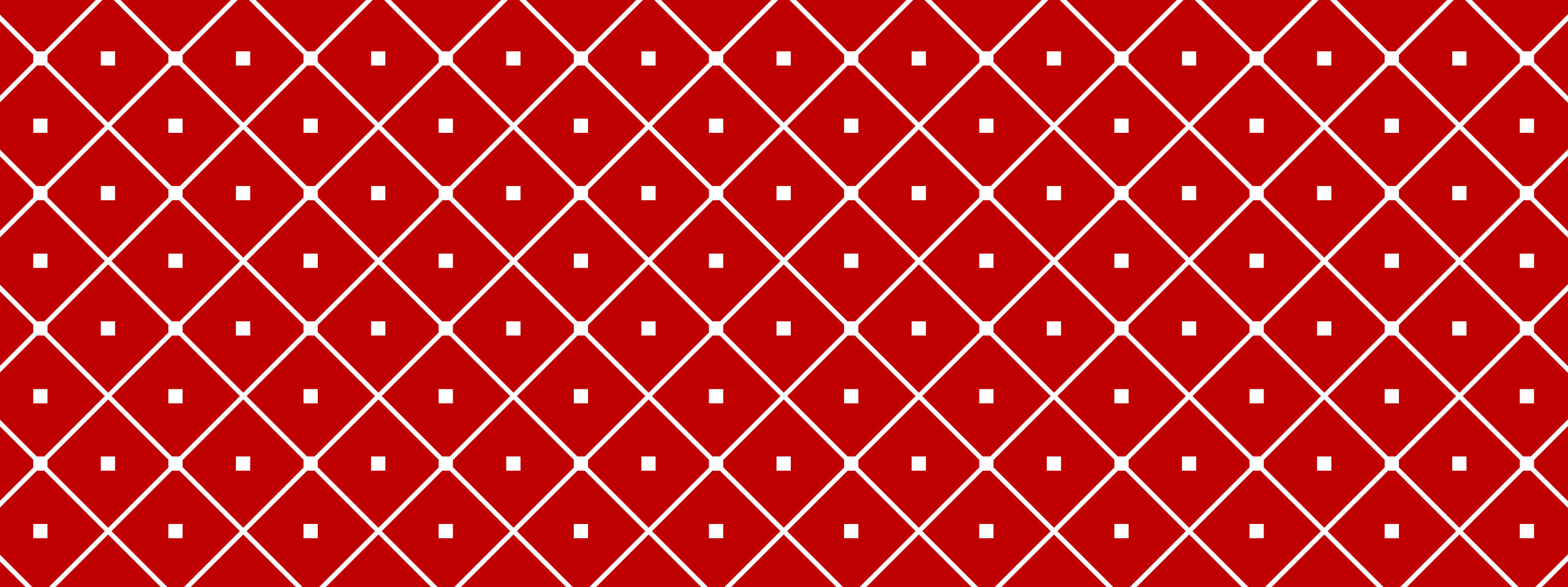
Cleaning files to be aligned also seems to improve performance, though not as much as cleaning the training files does.

It is better to use a smaller number of clean files than a larger number of mixed clean and dirty files when training acoustic models, even if the files to be aligned are dirty.

It is acceptable to use the same files in both the training and the alignment processes.

The default Targetrate setting of 100000 seems to produce better alignments than the adjusted 125000 setting.

Increasing the number of epochs used in the training process did not produce better alignments, though it did increase the time required to train the acoustic models.



MONTREAL FORCED ALIGNER

A New Alternative

MONTREAL FORCED ALIGNER

Created at the same lab as Prosodylab-Aligner

Like PL-A, can train and align same data or use pretrained acoustic model

Uses Python scripts like PL-A

Uses a different underlying technology:

- [Kaldi ASR](#) toolkit instead of HTK

Goes through three stages of training:

- First pass with monophone models
- Second pass using triphone models, which take into account the sound on both sides of the target phone
- Final pass that enhances triphone models by taking into account speaker differences

Has been used on:

- Bulgarian, Mandarin, Croatian, Czech, French, German, Hausa, Korean, Polish, Portuguese, Russian, Swahili, Spanish, Swedish, Thai, Turkish, Ukrainian, Vietnamese, English, Afrikaans, English, Ndebele, Xhosa, Zulu, Setswana, Sesotho sa Leboa, Sesotho, siSwati, Tshivenda, Xitsonga (working on Japanese)

ADVANTAGES OF MONTREAL FORCED ALIGNER

1. Accounts for interspeaker differences by considering speaker ID during acoustic model training.
2. Can align for multiple speakers in the same file
3. Can align without a dictionary if working from a fairly transparent and consistent orthography.
4. Does not crash when encountering out-of-dictionary words
unknown word marked as <unk> in the output and list of unknown words generated
5. Automatically strips punctuation from ends of words in transcripts and converts capital letters to lowercase.
6. Accepts two kinds of transcription inputs: PL-A format or Praat TextGrid format

MFA INPUT

Audio Files

Must be in .wav format

Any sampling rate above 16kHz* accepted—consistent sampling rate for each speaker

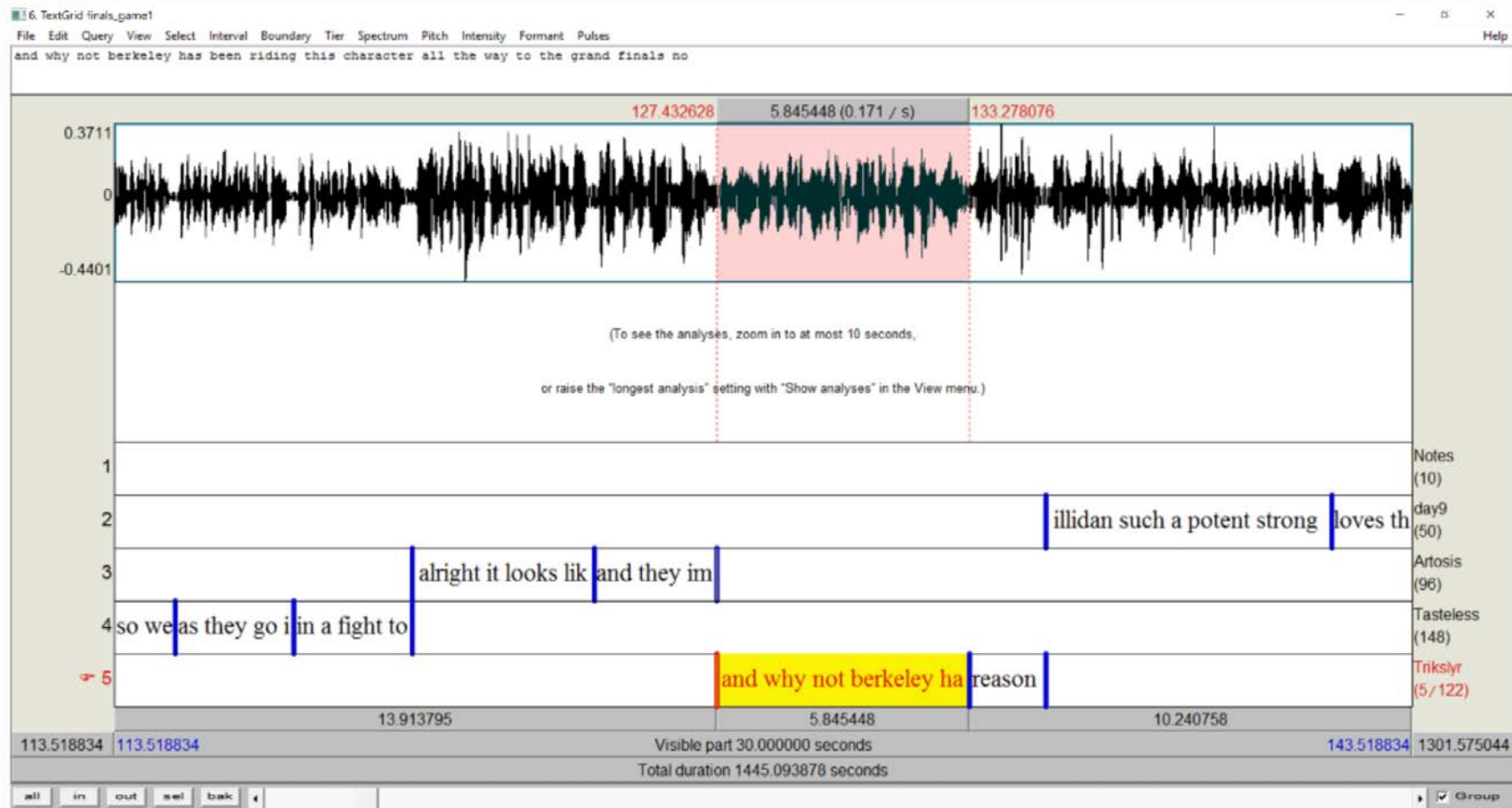
Audio “chunks” should be less than 30 seconds (sound files for PL-A format and intervals for Textgrid format)

Transcription Files

Two allowable formats:

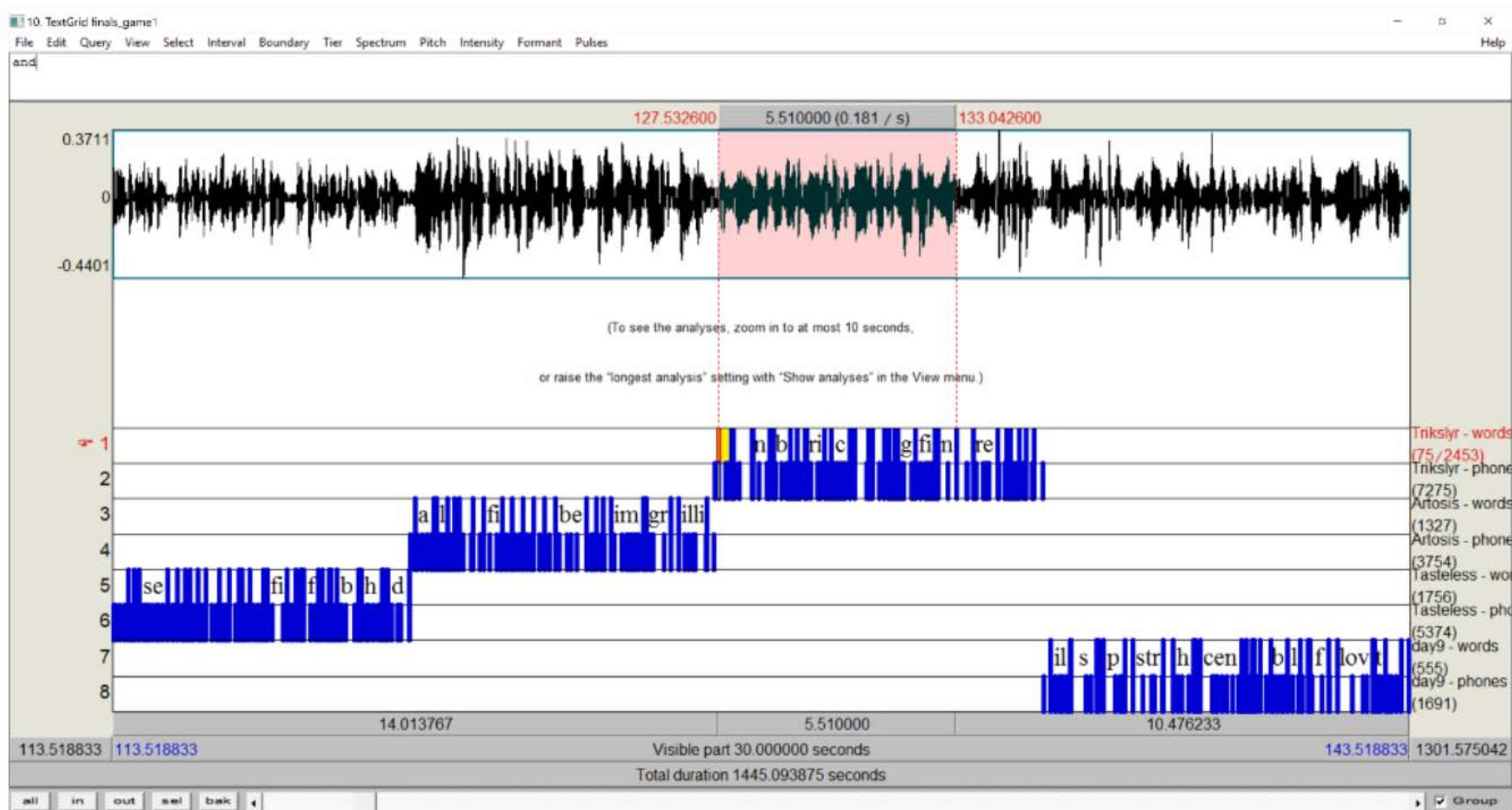
1. PL-A format (plain text, as described in previous slides)
2. TextGrid format (with transcribed “chunks” > 100ms and < 30 seconds)

INPUT EXAMPLE (TIER NAME = SPEAKER ID)



http://montreal-forced-aligner.readthedocs.io/en/stable/data_format.html#prosodylab-format

OUTPUT EXAMPLE (WORD AND PHONE TIERS)



TRAINING AND ALIGNMENT

Accomplished in one step using a Python script

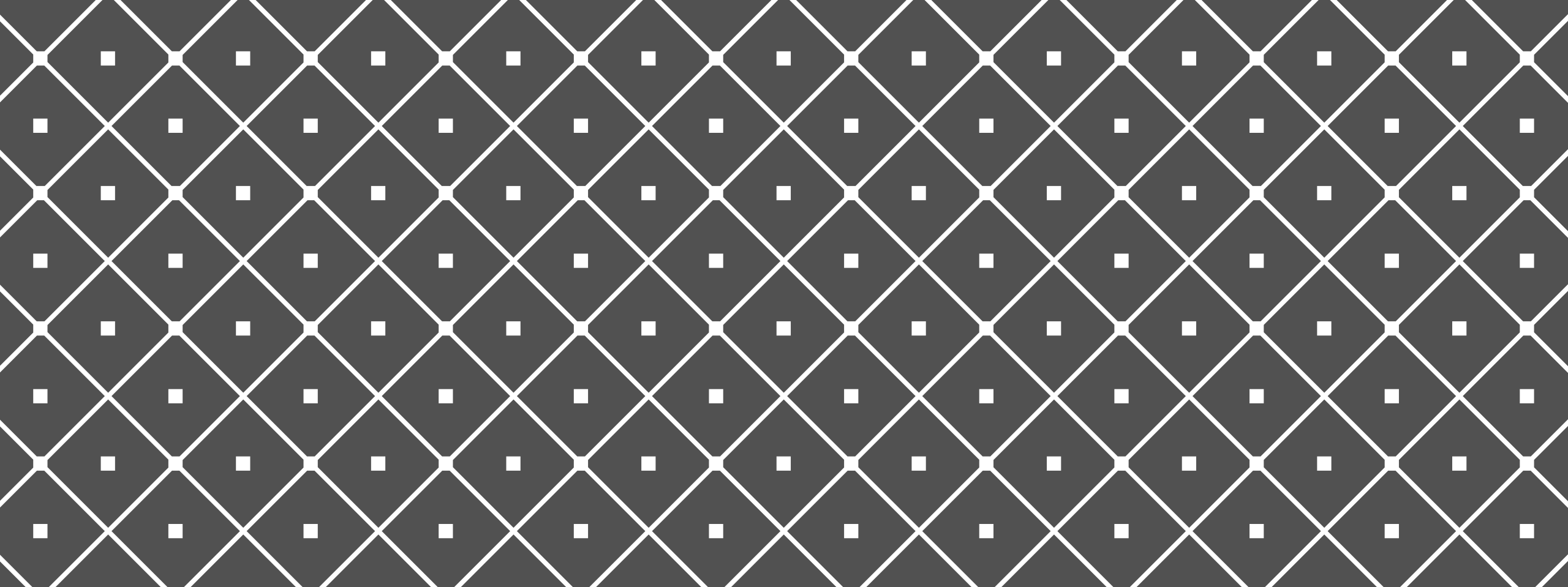
```
bin/mfa_train_and_align corpus_directory [dictionary_path] output_directory
```

Training can be skipped if aligning with a pretrained model

```
bin/mfa_align [model_path] corpus_directory output_directory
```

List of available options for both processes:

<http://montreal-forced-aligner.readthedocs.io/en/stable/aligning.html>

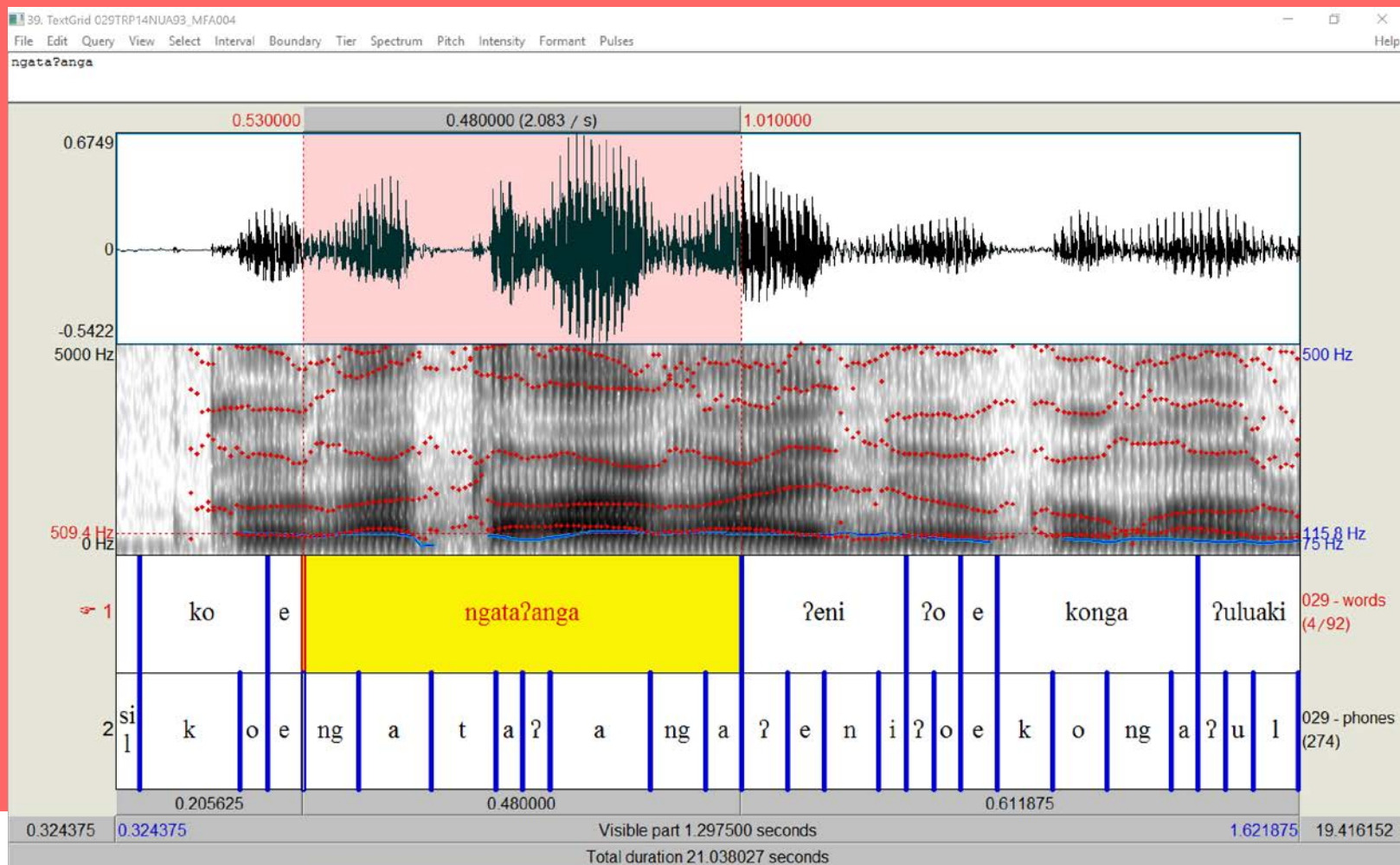


MFA TESTS

Model Training and Alignment

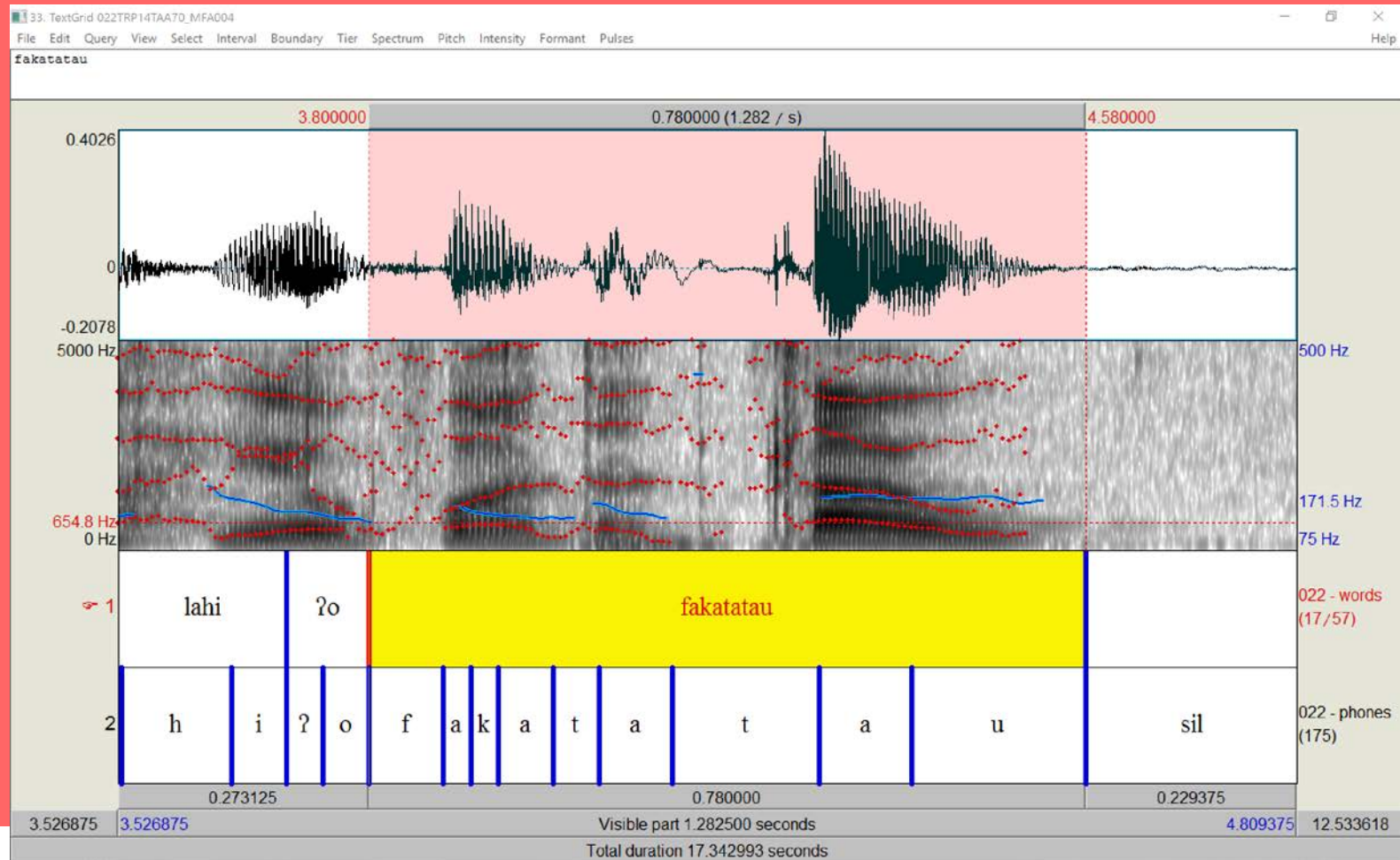
MFA TRAINING AND ALIGNMENT TESTS

Test ID #	Type and Number of Audio Files	Type of Transcription	Name of Acoustic Model Created	Runtime
MFA001	clean WL (same as PL-A) (22 files)	TextGrid	ton-001MFA.zip	:39
MFA002	clean WL (same as PL-A); dirty WL (20-second chunks) (38 files)	TextGrid	ton-002MFA.zip	1:14
MFA003	clean WL (same as PL-A); dirty WL (1-word chunks) (38 files)	TextGrid	ton-003MFA.zip	1:30
MFA004	clean (same as PL-A); dirty (1-word chunks); reading passage (59 files)	TextGrid	ton-004MFA.zip	1:30
MFA005	clean and dirty WL excerpts (10 files)	Text (PL-A)	(aligned only)	:02



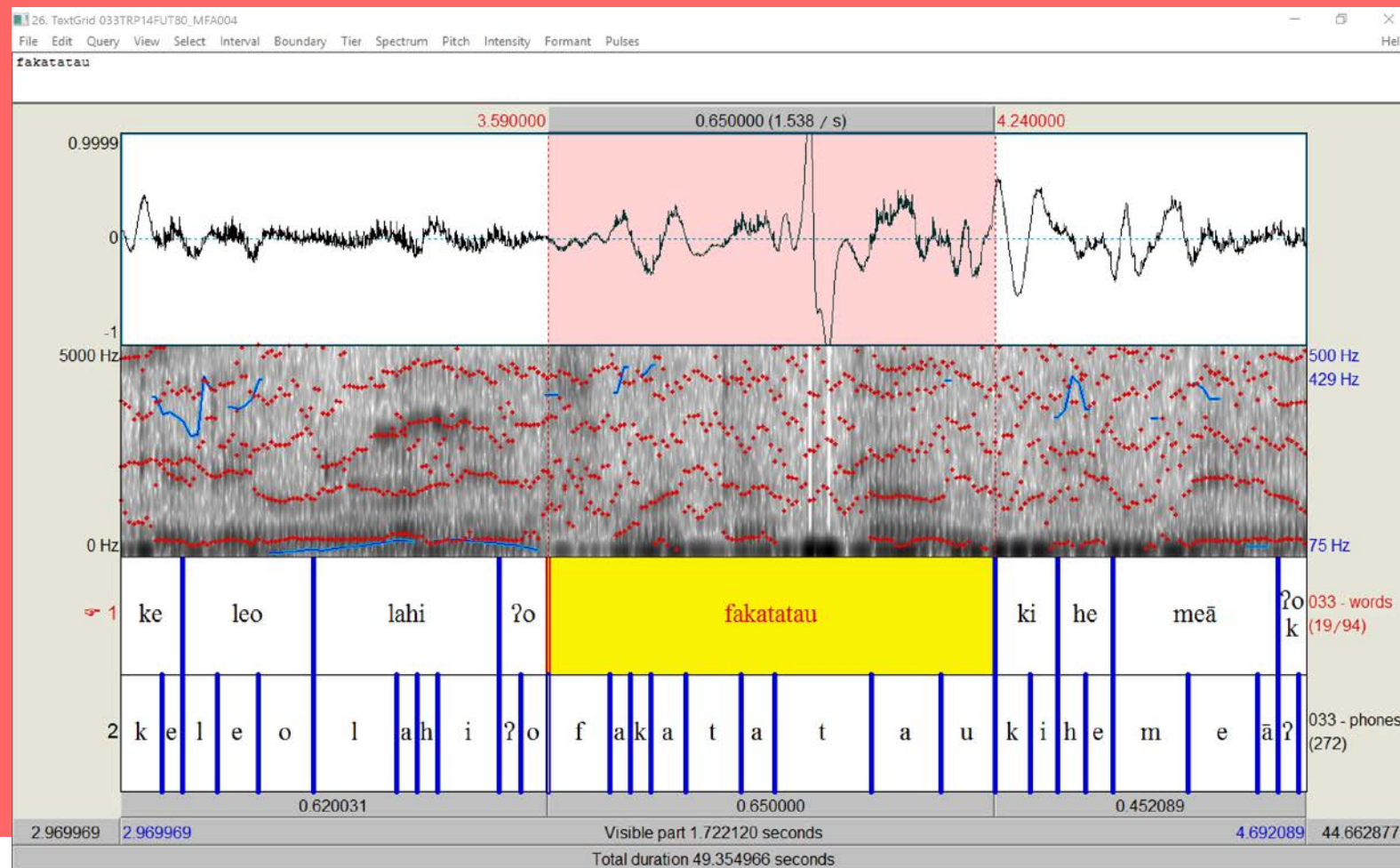
MFA ALIGNMENT

Connected Speech
(Speaker 029)



MFA ALIGNMENT

Connected Speech
(Speaker 22)



MFA ALIGNMENT

Connected Speech
(Speaker 33)

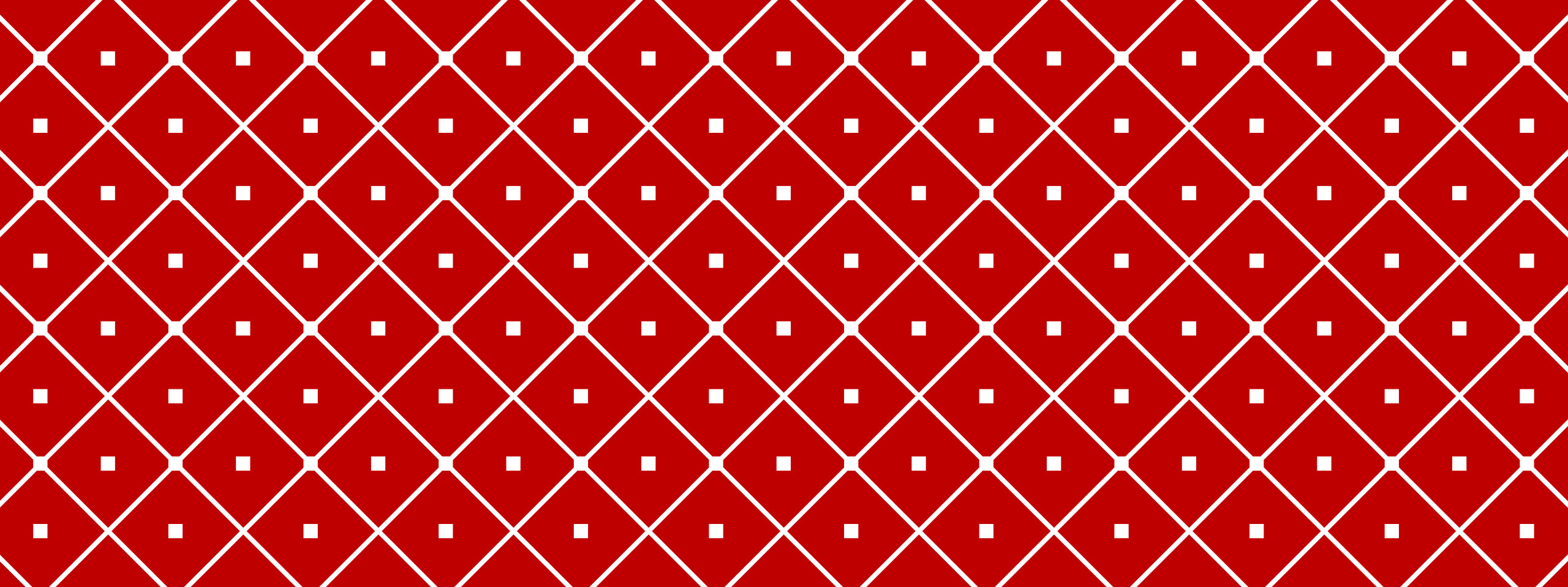
MFA SUMMARY AND CONCLUSIONS

Quality

- Using MFA TextGrid input seems to eliminate the dirty file effects we saw with PL-A.
- MFA produced good alignments with long recordings, allowing us to preserve token context for analysis.

Efficiency

- In our experience, MFA file preparation was much more efficient than PL-A file prep.
- MFA's "no dictionary" option will save considerable time when we begin to analyze free conversation and interview speech. (Note, this may not be as effective for languages with less transparent orthography.)
- MFA's ability to process speech from multiple speakers in the same file will save prep time and preserve discourse context.



IMPLICATIONS AND APPLICATIONS

Feasibility and Efficiency

SUMMARY AND RECOMMENDATIONS

Efficiency

- Forced alignment can greatly reduce the time required to prepare files for acoustic analysis.
- It is **possible** and **efficient** to **force align field recordings**, even with background noise.
 - TextGrid input using MFA produces good alignments with less clean-up time
- The amount of time saved will vary by language and the type of analysis planned

Reliability and Validity

- Forced alignment **may improve general consistency and replicability**
- It's necessary to make **manual boundary adjustments of TextGrids** output from forced alignment
 - Positive: It allows you to dig deeply into an understudied language early on



ACKNOWLEDGEMENTS

Bell, Adrian V. and Marianna Di Paolo. 2014. University Research Council Grant, University of Utah. “Language, Ethnic Markers, and the Adaptation of Tongan Immigrants to Utah.”

Holt, Carter. 2015. “Boosting Phonetics Research through Technology.” Undergraduate Research Opportunity Program. Summer 2015 Research Assistantship. (Marianna Di Paolo, mentor)

Special thanks to Kyle Gorman, Michael McAuliffe, and Craig Johnson

REFERENCES

- Albin, Aaron L. 2014. "PraatR: An architecture for controlling the phonetics software "Praat" with the R programming language." *The Journal of the Acoustical Society of America* 135 (4):2198-2199.
- Boersma, Paul and David Weenink. 2015. Praat: doing phonetics by computer [Computer program] 5.4.22.
- Cambridge University. 1989-2015. HTK Hidden Markov Model Toolkit.
- Carnegie Mellon University. 1993-2016. "CMU Pronouncing Dictionary." <http://www.speech.cs.cmu.edu/cgi-bin/cmudict>.
- Evanini, Keelan, Stephen Isard, and Mark Liberman. 2009. "Automatic formant extraction for sociolinguistic analysis of large corpora." *INTERSPEECH*.
- Goldman, Jean-Philippe. 2011. "Esyalalign: an automatic phonetic alignment tool under Praat." *Interspeech-2011*:3233-3236.
- Gorman, Kyle, Jonathan Howell and Michael Wagner. 2011. Prosodylab-Aligner: A Tool for Forced Alignment of Laboratory Speech. *Canadian Acoustics*. 39.3. 192–193.
- Kisler, Thomas, Florian Schiel, and Han Sloetjes. 2012. "Signal processing via web services: the use case WebMAUS." *Digital Humanities Conference 2012*. McAuliffe, Michael, Michaela Socolof, Sarah Mihuc, and Michael Wagner (2016). Montreal Forced Aligner [Computer program]. Version 0.5, retrieved 13 July 2016 from <http://montrealcorpus-tools.github.io/Montreal-Forced-Aligner/>.
- Povey, Daniel, Arnab Ghoshal, Giles Boulianne, Lukas Burget, Ondrej Glembek, Nagendra Goel, Mirko Hannemann, Petr Motlicek, Yanmin Qian, Petr Schwarz, Jan Silovsky, Georg Stemmer, and Karel Vesely. 2011. "The Kaldi Speech Recognition Toolkit." *IEEE 2011 Workshop on Automatic Speech Recognition and Understanding*.
- Reddy, Sravana, and James Stanford. 2015. "Toward completely automated vowel extraction: Introducing DARLA." *Linguistics Vanguard*.
- Rosenfelder, Ingrid. 2013. "Forced Alignment & Vowel Extraction (FAVE): An online suite for automatic vowel analysis." *University of Pennsylvania Linguistics Lab*, Last Modified December 8, 2013, accessed November 26, 2015. <http://fave.ling.upenn.edu/index.html>.
- Rosenfelder, Ingrid, Joe Fruehwald, Keelan Evanini, and Jiahong Yuan. 2011. FAVE (Forced Alignment and Vowel Extraction) Program Suite.
- Schiel, Florian, Christoph Draxler, Angela Baumann, Tania Ellbogen, and Alexander Steffen. 2012. "The Production of Speech Corpora."
- Young, Steve, Gunnar Evermann, Mark Gales, Thomas Hain, Dan Kershaw, Xunying (Andrew) Liu, Garth Moore, Julian Odell, Dave Ollason, Dan Povey, Valtcho Valtchev, and Phil Woodland. 1995-2006. *The HTK Book (for HTK Version 3.4)*. edited by Microsoft Corporation (1995-1999) Cambridge University Engineering Department (2001-2006).