

LFA: Long-speech Forced Alignment with Phoneme-based Reference Assignment

Anonymous submission to Interspeech 2026

Abstract

Forced alignment of short utterances is well-solved, but aligning long speech recordings against reference transcripts remains challenging in practical applications. We propose LFA, a modular divide-and-conquer approach that decomposes long recordings into short chunks via ASR, then assigns reference transcript portions to each chunk using word-level dynamic time warping with edit distance on phonemized word representations. The modular design allows us to use any alignment method on the resulting chunks. We evaluate LFA on noise-augmented TIMIT with file lengths from 5 min to 9.7 h. On these recordings, MFA fails or degrades severely, and NeMo baselines degrade with length. Depending on the per-chunk alignment backend, LFA achieves up to 52 ms median onset error or 68.4% onset accuracy@100 ms on a single 9.7-hour file, maintaining stable quality across all durations.

Index Terms: forced alignment, automatic speech recognition, dynamic time warping, CTC alignment, long speech, reference transcript assignment

1. Introduction

Forced alignment (FA) is the task of determining precise time boundaries for words or phonemes given a speech recording and a transcript. It is fundamental to speech science, corpus linguistics, and language documentation. While FA for short, clean utterances is well-solved, e.g., with the Montreal Forced Aligner (MFA) [1] achieving 97.4% word-level accuracy at 100 ms tolerance on TIMIT [2], alignment of long speech recordings (minutes to hours) remains challenging in practical applications. Early work addressed long recordings using anchor-point methods based on Hidden Markov Models (HMM) [3, 4], but to the best of our knowledge, no prior work has systematically varied recording length to measure how alignment accuracy scales.

Parliamentary corpora span hours per session [5, 6, 7], and oral history and language documentation archives regularly contain multi-hour recordings. These applications require word-level alignment of recorded speech against existing reference transcripts prepared by experts and following specific conventions, rather than using ASR output directly.

Beyond alignment itself, reliable long-speech FA enables large-scale data mining for ASR training: vast quantities of paired speech and text exist as audiobooks, lectures, and parliamentary proceedings, and adding word-level time annotations would turn these into fine-grained training data [8].

Existing FA systems face two scalability problems on long recordings. First, CTC-based exact Viterbi decoding requires dynamic programming over a trellis whose size scales with $O(T \times L)$ (acoustic frames \times label sequence length), which can

become impractical for hour-long recordings, motivating linear-memory CTC formulations [9], linear-memory edit distance on decoded phone sequences [10], and divide-and-conquer approaches that keep T small by chunking [3, 4, 11]. Second, HMM-based aligners such as MFA can fail entirely on noisy recordings. In our experiments, standalone MFA produces no useful output on noise-augmented speech.

We address this through a modular divide-and-conquer approach called LFA: we use ASR to decompose the long-speech alignment problem into short, independently solvable chunks, then align each chunk with any FA method of choice. Our contributions are:

1. *Divide-and-conquer reference assignment via phoneme DTW:* A method that converts both ASR hypotheses and reference words to phoneme sequences via grapheme-to-phoneme (G2P) conversion, then uses dynamic time warping (DTW) with the Damerau-Levenshtein distance on phonemized word pairs to partition the reference across audio chunks. This tolerates ASR errors because phonemically similar words yield low cost even when text matching fails; a cumulative assignment strategy further handles word splits and merges (Section 3.3, Figure 1).
2. *Controlled scaling evaluation on noisy speech:* We construct noise-augmented TIMIT length variants spanning 5 min to 9.7 h with identical speech content, providing the first controlled experiment systematically varying recording length on challenging audio.
3. *Flexible use of FA methods:* We show that our decomposition enables us to use any forced alignment method to process recordings of arbitrary length. We evaluate LFA with three different backends — HMM-based (MFA), wav2vec2 CTC, and NeMo Conformer. All variants maintain stable accuracy from 5 min to 9.7 h, with quality determined by the backend’s short-segment performance rather than recording duration.

2. Related Work

2.1. Forced Alignment Systems

The Montreal Forced Aligner (MFA) [1] uses Kaldi-based GMM-HMM acoustic models with pronunciation dictionaries and remains one of the most widely used FA tools in phonetics research [2]. However, MFA requires language-specific acoustic models and pronunciation dictionaries, and its behavior on long or noisy recordings has not been systematically studied.

The NeMo toolkit provides two alignment tools. NeMo Forced Aligner (NFA) [13] applies Viterbi decoding to CTC log-probabilities; for long recordings, the model forward pass may need to be chunked to fit memory. NeMo’s CTC segmentation [14] uses dynamic programming with optional zero-cost

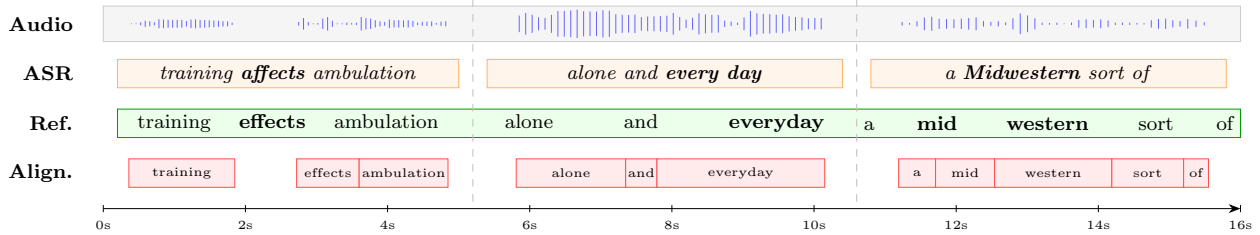


Figure 1: ASR error types from Buckeye [12]: substitution (“affects” for “effects”; $c = 0.17$), word split (“every day” for “everyday”), and word merge (“Midwestern” for “mid western”). Word-level DTW assigns reference words despite these errors via low phonemic cost (substitutions) and cumulative assignment (splits/merges). Bottom row: final aligned word boundaries.

95 blank transitions, scaling to longer recordings but with poten- 143
 96 tially reduced accuracy. 144

97 WhisperX [15] is not a forced alignment system in the strict 145
 98 sense but is closely related. It combines Whisper ASR with 146
 99 VAD-based segmentation and wav2vec2 alignment. However, 147
 100 WhisperX aligns its own ASR output rather than an external re- 148
 101 ference transcript, making it primarily an enhanced ASR system 149
 102 with word timestamps rather than a tool for aligning a given 150
 103 reference transcript to audio. Our system builds on the VAD 151
 104 chunking idea but addresses this latter problem.

105 Aeneas [16] aligns text fragments against audio using TTS 152
 106 synthesis and MFCC-level DTW, handling long recordings at 153
 107 the acoustic level without ASR. It is widely deployed but oper- 154
 108 ates on frame-level features rather than word-level alignment. 155

109 Moreno et al. [3] proposed a recursive algorithm that uses 156
 110 ASR to identify anchor points and recursively aligns segments 157
 111 between anchors. Hazen [4] adopted a similar anchor-based 158
 112 approach but replaced the recursion with a fixed three-stage 159
 113 pipeline, using an FST-based pseudo-forced alignment that tol- 160
 114 erates errors in approximate transcripts. SailAlign [11] com- 161
 115 bined iterative ASR with text matching to find reliable anchors, 162
 116 then performed per-chunk FA. DSAlign [17] used DeepSpeech 163
 117 ASR with character-level Smith-Waterman local alignment in a 164
 118 recursive divide-and-conquer strategy: longer ASR fragments 165
 119 are matched first to anchor the alignment, and shorter frag- 166
 120 ments are resolved recursively within the remaining intervals. 167
 121 Candidate reference windows are ranked by shared character 168
 122 3-gram count before full Smith-Waterman scoring. DSAlign 169
 123 is the closest predecessor to our work: it shares the same 170
 124 VAD-ASR-text-alignment pipeline structure, but uses ortho- 171
 125 graphic character matching and depends on the discontinued 172
 126 Mozilla DeepSpeech. Bordel et al. [10] aligned 3-hour bilingual 173
 127 parliamentary recordings by running a language-independent 174
 128 HMM phone decoder, then aligning the decoded phone se- 175
 129 quence against the reference using Hirschberg’s linear-memory 176
 130 edit distance at the individual phone level — explicitly avoiding 177
 131 language models or lexical resources to handle code-switching 178
 132 between Basque and Spanish. Álvarez et al. [18] extended this 179
 133 approach for automatic subtitling, showing that predefined non- 180
 134 binary phone-relatedness matrices outperform binary matching. 181
 135 Both operate at the individual phone level, producing align- 182
 136 ment matrices orders of magnitude larger than our word-level 183
 137 approach. In the parliamentary speech domain, Ljubešić et al. 184
 138 [7] aligned over 5,000 hours of Slavic parliamentary speech us- 185
 139 ing word-histogram sliding windows with Levenshtein distance 186
 140 for text-to-text matching, achieving a 74% word yield for Croa- 187
 141 tian, with the remaining 26% of words unmatched due to ASR 188
 142 errors. Unlike these phone-level predecessors, LFA operates

143 at the word level and separates reference assignment from per- 144
 145 chunk alignment, allowing any aligner on the resulting short 146
 147 chunks.

148 Rousso et al. [2] compared MFA, WhisperX, and MMS on 149
 150 TIMIT and Buckeye, finding that MFA is superior on short ut- 151
 152 terances. Deng et al. [19] provided a large-scale FA bench- 153
 154 mark including long conversational audio, confirming that long- 155
 156 speech FA remains challenging. 157

158 2.2. CTC-based Alignment 159

160 CTC [20] provides a natural framework for FA through its 161
 162 frame-level output probabilities. Doras et al. [9] devel- 163
 164 oped a linear-memory CTC alignment algorithm by adapting 164
 165 Hirschberg-style divide-and-conquer to the CTC Viterbi ma- 165
 166 trix, reducing memory from $O(T \times L)$ to $O(T + L)$ at the 166
 167 cost of additional computation. This solves the memory bot- 167
 168 tleneck for single-file CTC alignment, but requires the full text 168
 169 to be aligned in one pass. Our approach instead decomposes 169
 170 the problem via ASR-guided chunking, then uses CTC on each 170
 171 short chunk. 171

172 Recent work has improved CTC alignment accuracy. 172
 173 Huang et al. [21] showed that label priors reduce CTC’s peaky 173
 174 output distributions, improving boundary accuracy by 12–40%. 174
 175 Meta’s MMS [22] provides multilingual wav2vec2 models 175
 176 across 1,130 languages with a forced alignment tool. Charsiu 176
 177 [23] offers wav2vec2-based phone-level FA. Mu et al. [24] re- 177
 178 formulated FA as a slot-filling task using a Qwen3-based LLM 178
 179 with an audio encoder, predicting discrete timestamps non- 179
 180 autoregressively. Although titled “long-form,” their evaluation 180
 181 concatenates short utterances up to 5 min with 80 ms timestamp 181
 182 resolution, whereas our method at 20 ms CTC frame resolution 182
 183 (wav2vec2 FA) works on recordings up to 9.7 h. Rehman et 183
 184 al. [25] presented a CTC-based multilingual aligner focused on 184
 185 phoneme-level temporal boundaries. 185

186 2.3. Sequence Alignment with Dynamic Programming 187

188 Dynamic time warping (DTW) [26] is a classic algorithm 189
 189 for aligning temporal sequences, extensively used in mu- 190
 190 sic information retrieval [27] and speech processing. For 191
 191 text-to-audio alignment, González-Carrasco et al. [28] used 192
 192 word-level Needleman-Wunsch sequence alignment with ortho- 193
 193 graphic word comparison for live subtitle synchronization, later 194
 194 extended by Masiello-Ruiz et al. [29] with adaptive selection 195
 195 among NW, Levenshtein, and Smith-Waterman algorithms. Our 196
 196 approach differs by (1) using DTW rather than alignment based 197
 197 on the edit distance, so that every transition maps an ASR word 198
 198 to a reference word without free gap operations, and (2) operat-

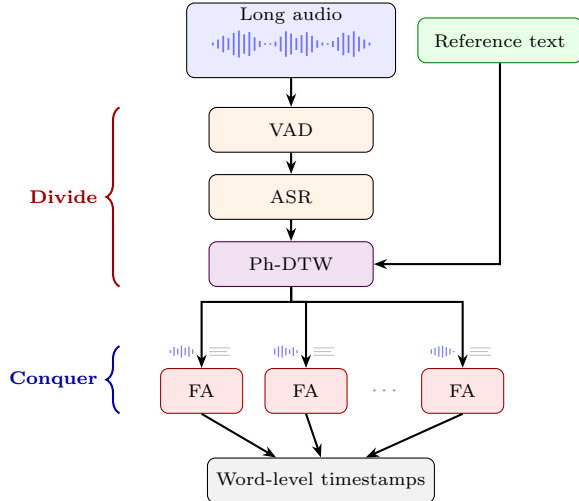


Figure 2: LFA divide-and-conquer approach. *Divide*: split audio into chunks via VAD, transcribe each chunk with ASR, assign reference transcript portions to chunks via phoneme DTW. *Conquer*: Align each chunk independently using forced alignment.

ing at the word level with a phoneme edit distance cost, rather than aligning individual phone sequences. Jin et al. [30] independently used DTW with phoneme-based costs for on-device forced alignment, but at the frame-to-phoneme level for short utterances rather than word-level reference assignment.

While phoneme similarity has been used at the phone level for long speech alignment [18, 10], our approach operates at the word level: we apply DTW with phoneme-based edit distance on ASR word hypotheses against reference words. This yields an $N_{\text{ASR}} \times K$ cost matrix (ASR words \times reference words), which is orders of magnitude smaller than the $T \times L$ CTC frame-label matrix or phone-level alignment matrices, making it tractable even for multi-hour recordings. The cumulative assignment strategy (Section 3.3.4) then maps phoneme DTW paths to chunk-level reference partitions.

3. Method

Our approach, called LFA, follows a divide-and-conquer strategy (Figure 2): divide the long recording into short chunks using VAD, transcribe each chunk with ASR, assign reference transcript portions to each chunk via phoneme DTW, and solve each chunk independently with forced alignment. The four stages are: (1) VAD, (2) ASR, (3) Phoneme DTW Reference Assignment, and (4) Per-Chunk Forced Alignment.

3.1. Voice Activity Detection

The input audio is segmented into speech regions using Pyannote VAD [31] with onset threshold 0.2 and offset threshold 0.05 (low thresholds to favor recall over precision, since missed speech is more harmful than false alarms for alignment). Long speech regions are further split into chunks of at most C seconds (default $C = 10$) at detected pause boundaries. This produces a sequence of audio chunks $\mathbf{S} = (s_1, \dots, s_N)$ with corresponding time spans.

3.2. Automatic Speech Recognition

Each chunk s_i is transcribed using Faster Whisper [32, 33] (large-v3). The ASR produces word-level hypotheses $\hat{\mathbf{w}}_i = (\hat{w}_{i,1}, \dots, \hat{w}_{i,M_i})$ for each chunk with approximate time-stamps.

3.3. Phoneme DTW Reference Assignment

The main problem is assigning the correct portion of the reference transcript $\mathbf{R} = (r_1, \dots, r_K)$ to each audio chunk. The simplest approach would use orthographic matching, but this may fail when ASR errors cause divergence from the reference. We exploit that phonemic forms often differ less than orthographic forms under noisy ASR conditions by comparing words through their phonemic representations (e.g., homophones such as “through”/“threw”, both $/\theta ru:/$).

3.3.1. Grapheme-to-Phoneme Conversion

Both ASR and reference words are converted to IPA phoneme strings using Epitran [34], a rule-based G2P system supporting 61 languages. Rule-based G2P can produce errors on irregular words, but the edit-distance cost function (Eq. 1) tolerates such errors: even imperfect phonemizations produce lower costs for correct word pairs than for unrelated words. For each word w , we obtain its phoneme string $\phi(w)$ — for example, $\phi(\text{“effects”}) = \text{“ifekts”}$.

3.3.2. Word-Level Cost Matrix

We construct an $N_{\text{ASR}} \times K$ cost matrix where N_{ASR} is the total number of ASR words across all chunks and K is the number of reference words. The cost for word pair (i, j) uses the Damerau-Levenshtein (DL) similarity [35] on their phoneme strings:

$$c(i, j) = 1 - \text{DL}_{\text{sim}}(\phi(\hat{w}_i), \phi(r_j)) \quad (1)$$

where DL_{sim} returns normalized similarity (0 = completely different, 1 = identical), so $c = 0$ for identical words and $c = 1$ for completely different ones. DL extends Levenshtein distance with transpositions, additionally handling phoneme metathesis. By comparing phonemic forms rather than orthography, homophones and near-homophones receive low cost despite orthographic differences. Figure 3 shows a word-level cost matrix example.

The matrix dimensions are $N_{\text{ASR}} \times K$ words, not $T \times L$ frames as in CTC alignment. For a 2-hour recording with $\sim 11,000$ words per side, this is an $11\text{K} \times 11\text{K}$ word matrix, i.e., orders of magnitude smaller than the millions of frames that make full-file CTC alignment infeasible.

3.3.3. DTW Alignment

We use DTW [26] to find the optimal alignment path through the cost matrix. The DTW recurrence is:

$$D(i, j) = c(i, j) + \min(D(i-1, j-1), D(i-1, j), D(i, j-1)) \quad (2)$$

We use DTW rather than edit-distance-based alignment (Needleman-Wunsch). In DTW, all three transitions (diagonal, vertical, horizontal) add the local cost $c(i, j)$: there are no free gap operations. This means every step on the path maps an ASR word to a reference word. Diagonal moves give 1:1 alignment; vertical moves handle many-to-one cases where ASR splits a word (“every day” \rightarrow “everyday”); horizontal moves handle

Algorithm 1 Phoneme DTW Reference Assignment

Require: ASR words \hat{w} with chunk IDs, reference words \mathbf{R} , phonemizer ϕ

Ensure: Reference word assignment per chunk

- 1: Compute cost $c(i, j) = 1 - \text{DL}_{\text{sim}}(\phi(\hat{w}_i), \phi(r_j))$ for all i, j
 - 2: Initialize $D(0, 0) = 0$; all other $D = \infty$
 - 3: **for** $i = 1$ to N_{ASR} **do**
 - 4: **for** $j = 1$ to K **do**
 - 5: $D(i, j) \leftarrow \min(D(i-1, j-1), D(i-1, j), D(i, j-1)) + c(i, j)$
 - 6: **end for**
 - 7: **end for**
 - 8: $i^* \leftarrow \arg \min_{1 \leq i \leq N_{\text{ASR}}} D(i, K)$
 - 9: Backtrack from (i^*, K) to $(1, 1)$ to obtain path π
 - 10: **for** each reference word r_j **do**
 - 11: Collect aligned ASR positions $\mathcal{A}(j) = \{(i, 1 - c(i, j)) : (i, j) \in \pi\}$
 - 12: Assign r_j to $\arg \max_k \sum_{(i,s) \in \mathcal{A}(j), \text{chunk}(i)=k} S$
 - 13: **end for**
-

272 one-to-many cases where ASR merges words (“Midwestern”
273 \rightarrow “mid western”).

274 We use phoneme DTW with open-end termination: the path
275 must start at $(1, 1)$ but may terminate at any row in the last
276 column, i.e., at (i^*, K) for any $i^* \leq N_{\text{ASR}}$. This ensures
277 that all K reference words participate in the alignment while
278 allowing surplus ASR words at the end (e.g., from hallucina-
279 tions or noise-induced insertions) to remain unmatched. For
280 scalability, we apply a Sakoe-Chiba band constraint [26] with
281 width $B = 0.05 \cdot \max(N_{\text{ASR}}, K)$, restricting computation to
282 a diagonal band and reducing memory from $O(N_{\text{ASR}} \times K)$ to
283 $O(N_{\text{ASR}} \times B)$. Algorithm 1 shows the procedure (banding omit-
284 ted for clarity).

285 3.3.4. Cumulative Assignment

286 The DTW path aligns each reference word to one or more ASR
287 words. Since each ASR word belongs to a known chunk, we
288 assign each reference word to the chunk whose aligned ASR
289 words contribute the highest cumulative similarity. Formally,
290 for reference word r_j , let $\mathcal{A}(j) = \{(i, s_i) : (i, j) \in \pi\}$ be its
291 aligned ASR positions with similarities $s_i = 1 - c(i, j)$. We
292 assign r_j to chunk $k^* = \arg \max_k \sum_{(i,s_i) \in \mathcal{A}(j), \text{chunk}(i)=k} S_i$.
293 This handles word splits and merges: when ASR splits a refer-
294 ence word across chunks, cumulative scoring assigns it to the
295 chunk with the stronger match. When ASR is highly accurate,
296 most reference words align unambiguously to a single chunk
297 but under noisy conditions, the cumulative scoring can resolve
298 ambiguous cross-boundary alignments.

299 3.4. Per-Chunk Forced Alignment

300 Each chunk, with its DTW-assigned reference words, is aligned
301 independently using an FA method. Any aligner capable of pro-
302 cessing short chunks (typically ≤ 20 s) can be used. We evaluate
303 three options:

304 *CTC alignment with wav2vec2* [36, 37]: The au-
305 dio is passed through the encoder to obtain frame-level
306 CTC log-probabilities (~ 20 ms per frame). Reference
307 words are decomposed into characters, and Viterbi decod-
308 ing finds the optimal frame-to-label assignment. Charac-
309 ter timestamps are merged into word boundaries. We use

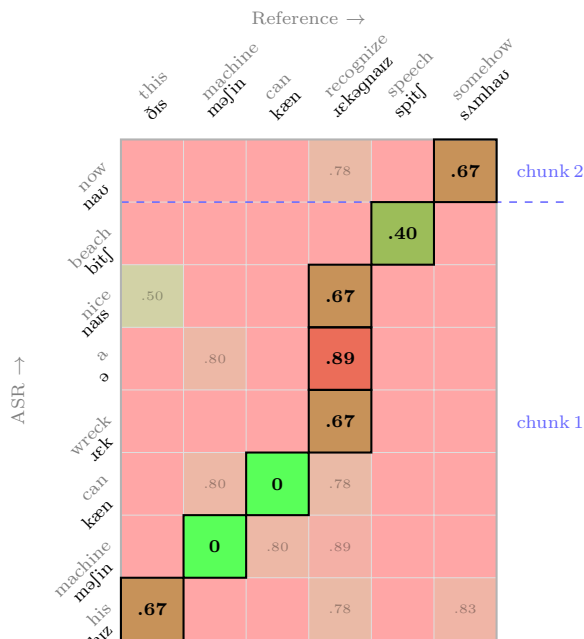


Figure 3: DTW cost matrix. Cell color encodes distance on IPA bold values mark the optimal path. Dashed lines: Chunk boundaries after ASR step.

wav2vec2-xls-r-1b-english¹ (1B parameters) as de-
fault; the CTC model choice is the single largest factor affecting
alignment quality (Section 4.5).

MFA alignment [1]: Each chunk is aligned using MFA’s
GMM-HMM acoustic model with a pronunciation dictionary.
MFA provides phoneme-level HMM duration modeling, which
often yields finer boundaries than CTC’s frame-level resolution.

NeMo Conformer alignment: Each chunk is
aligned using NeMo’s Conformer CTC model
(stt_en_conformer_ctc_large) with Viterbi decod-
ing. This provides an additional CTC-based baseline that
differs from wav2vec2 in model architecture and training data.

For CTC alignment, several post-processing steps improve
quality. The wav2vec2 feature extractor expects zero-mean,
unit-variance input. We apply this normalization per-chunk
rather than globally, yielding more consistent CTC probabil-
ities across chunks with varying volume levels. CTC align-
ment tends to extend word boundaries into silence regions at
chunk edges. We detect silence by thresholding the RMS en-
ergy at chunk boundaries and adjust the first/last word bound-
aries inward to the speech onset/offset. We prepend a blank
column to the CTC trellis before Viterbi decoding, allowing the
alignment path to remain in silence before the first token rather
than absorbing preceding frames into the first word, and extend
word onset boundaries backward through preceding low-blank-
probability frames. Each segment is also extended by 0.5 s of
acoustic context at its left boundary before alignment. For MFA
alignment, we merge tokens that MFA splits at apostrophes and
hyphens back to match the reference tokenization.

¹<https://huggingface.co/jonatasgrosman/wav2vec2-xls-r-1b-english>

Table 1: *TIMIT-LN evaluation datasets derived from TIMIT. All variants contain identical audio content; file length differs.*

Dataset	Files	Words	Mean (min)	Range (min)
5 min	119	457	4.9	1.1–5.2
10 min	59	921	9.9	8.3–10.4
20 min	30	1,812	19.4	6.8–20.8
35 min	17	3,197	34.2	26.2–35.1
1 h	10	5,436	58.3	45.8–60.4
2 h	5	10,871	116.5	104.7–120.2
full	1	54,357	582.0	582.0

4. Experiments

4.1. Datasets

We derive a long and noisy benchmark dataset from the TIMIT corpus [38] (read American English, with manually verified time annotations) which we will refer to as TIMIT-LN. This dataset simulates real-world scenarios such as long parliamentary sessions, lecture recordings, oral history interviews or speeches, where word-precise alignment is necessary. We reorganize TIMIT corpus into acoustically challenging length variants. All 6,300 TIMIT utterances (train + test, 630 speakers) are shuffled so that each file contains a mix of speakers, then concatenated with three per-utterance degradations: room reverberation (RT60 0.3–0.6 s) using pyroomacoustics [39], random gain modulation (± 4 –6 dB), and random inter-utterance pauses (0–5 s). Mixed additive noise from the MUSAN corpus [40] is then applied at a global 15 dB SNR (computed over the entire file; representing moderate noise typical of meeting or field recording environments) with three components: babble speech ($\sim 60\%$), environmental noise ($\sim 25\%$), and white Gaussian noise ($\sim 15\%$), with per-file energy proportions drawn from a Dirichlet distribution. Ground truth annotations are shifted by cumulative duration offsets. This procedure yields identical speech content and word sequences across all variants; differences in accuracy are attributable solely to recording length.

Table 1 summarizes the seven evaluation datasets in TIMIT-LN, spanning 5 min to 9.7 h. The longest variant concatenates all 6,300 utterances into a single file.

TIMIT-LN is a controlled benchmark dataset under realistic acoustic conditions: at short durations (5 min), most FA systems still perform adequately despite noise; as file length increases through 1 h and 2 h, systems are increasingly stressed by memory growth, error propagation, and the challenge of aligning tens of thousands of words in a single pass.

Since all variants contain the same 54,357 words, differences in accuracy directly measure each system’s sensitivity to recording length.

We build on TIMIT rather than using naturally noisy corpora for two reasons. First, TIMIT provides accurate manually verified word-level time annotations, whereas alternatives have significant annotation limitations: Buckeye [12] leaves interviewer speech unannotated, and AMI [41] annotations were themselves generated by forced alignment and contain frequent disfluencies (filled pauses, hesitations). Second, constructing noise-augmented TIMIT variants avoids the cost of manually annotating 10 h of naturally noisy speech while providing a controlled variable (audio length) that cannot be easily isolated in existing corpora.

4.2. Forced Alignment Systems

We compare our approach against three widely-used baselines:

Montreal Forced Aligner (MFA) [1]: MFA v2.2 with pretrained `english_mfa` acoustic model, pronunciation dictionary, G2P for out-of-vocabulary (OOV) handling, and beam width 1000 (increased to 4000 on failure). MFA processes each file as a single unit without internal chunking. On our noisy multi-speaker datasets, MFA in its default mode fails to produce any alignments across all beam widths tested. In single-speaker mode, MFA produces output but with rapidly degrading accuracy (on@100 drops from 51% at 5 min to 32% at 10 min, to effectively unusable 5.5% at 20 min).

NeMo CTC Segmentation [14]: CTC-based forced alignment with QuartzNet (`stt_en_quartznet15x5`) and Conformer (`stt_en_conformer_ctc_large`). Uses dynamic programming over CTC log-probabilities where blank transitions between words have zero cost, finding the most likely segmentation boundaries.

NeMo Forced Aligner (NFA) [13]: CTC-based forced alignment with Citrinet (`stt_en_citrinet_1024`) and Conformer models. Uses Viterbi decoding (also dynamic programming, but over the full CTC state space including blank probabilities). For long recordings we chunk audio into 10 min segments for the model forward pass.

LFA (our approach): Pyannote VAD² (onset=0.2, offset=0.05), Faster Whisper large-v3³, open-end DTW with cumulative assignment (Epitran G2P, Damerau-Levenshtein phoneme distance), and per-chunk forced alignment with chunk size $C=10$ s. We evaluate with three per-chunk aligners: MFA, wav2vec2 CTC (`xlsr-1b-en`, 1B parameters), and NeMo Conformer (Conformer).

4.3. Evaluation Metrics

Following standard practice [13, 2], we report mean and median absolute onset and offset error in milliseconds, and accuracy at thresholds $\tau \in \{25, 50, 100, 200\}$ ms (on@ τ for onset, off@ τ for offset), defined as the fraction of reference words with error $\leq \tau$.

4.4. Results

Figure 4 and Table 2 show how alignment accuracy scales with recording length on TIMIT-LN variants.

Standalone MFA fails on noisy audio: in default mode, it produces no usable output at any beam width. In single-speaker mode, accuracy degrades from 51% on@100 at 5 min to 32% at 10 min and 5.5% at 20 min. Among CTC-based baselines, NFA Conformer and NFA Citrinet range from 38–50% on@100, with NFA Citrinet maintaining 50.3% at 5 min but dropping to 42.9% at full; NFA Conformer shows a non-monotonic dip at 1 h (38.2%) driven by catastrophic drift in a subset of files (3 of 10 files exhibit onset means above 10 s), partially recovering at 2 h (44.5%), though still below its 35 min level (47.4%). CTC-Seg Conformer performs similarly (48% at 5 min, 41% at full). CTC-Seg QuartzNet degrades catastrophically from 16.4% at 5 min to 1.1% at full, as its weaker acoustic model accumulates alignment drift over long recordings.

In contrast, all three LFA variants maintain stable accuracy across the full range from 5 min to 9.7 h. LFA with wav2vec2

²<https://huggingface.co/pyannote/voice-activity-detection>

³<https://huggingface.co/openai/whisper-large-v3>

Table 2: Onset accuracy@100 ms (%) across noisy TIMIT-LN length variants. All datasets contain the same 54,357 words.

System	On@100 ms (%) by dataset						
	5 min	10 min	20 min	35 min	1 h	2 h	full
<i>Baselines</i>							
CTC-Seg (QuartzNet)	16.4	12.1	10.4	4.4	7.2	0.6	1.1
CTC-Seg (Conformer)	48.0	47.1	44.6	43.6	43.3	40.7	40.7
NFA (Citrinet)	50.3	50.4	50.0	49.4	47.9	46.7	42.9
NFA (Conformer)	49.4	49.0	46.9	47.4	38.2	44.5	45.3
MFA [†]	51.1	31.5	5.5	1.8	—	—	—
<i>LFA (ours)</i>							
LFA (NeMo Conf., C=10)	45.4	45.6	45.5	44.7	45.1	44.4	46.5
LFA (MFA, C=10)	63.1	64.0	64.1	62.9	63.2	62.7	63.5
LFA (wav2vec2, C=10)	67.7	68.9	68.9	68.1	66.8	65.6	68.4

[†]Single-speaker mode; default mode produces no output. At 5 min, MFA produces matching word counts on only 9 of 119 files (3,801 words); accuracy is computed over these files only, not the full 54,357 words. No results beyond 35 min.

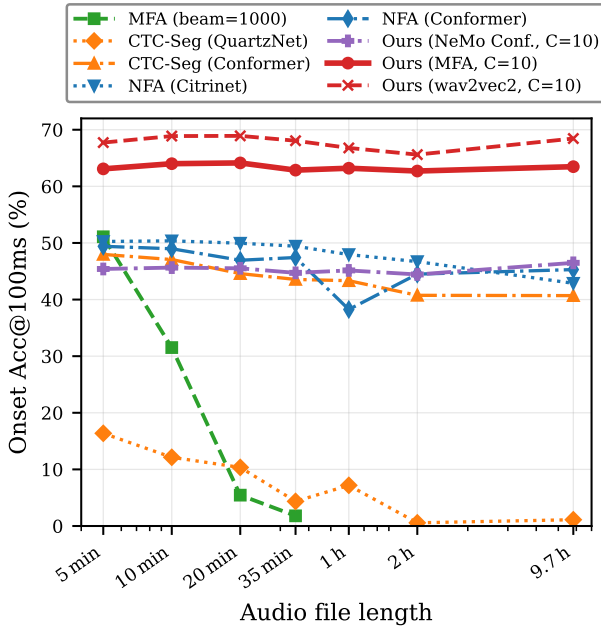


Figure 4: Onset accuracy@100 ms vs. file length (log scale) on TIMIT-LN. LFA configurations maintain stable accuracy from 5 min to 9.7 h, while baselines generally decline at longer durations.

442 achieves 65.6–68.9% on@100 (median onset 62–66 ms), LFA
 443 with MFA achieves 62.7–64.1% (median onset 52–54 ms), and
 444 LFA with NeMo achieves 44.4–46.5%. There is no scale degradation: LFA wav2vec2 scores 68.4% on the 9.7-hour file, within
 445 0.7 pp of its 5-minute score (67.7%).
 446

447 Table 3 provides a detailed breakdown on TIMIT-LN-full
 448 (1 file, 9.7 h). LFA with MFA per-chunk alignment achieves
 449 the best median onset (52 ms) with 63.5% on@100. LFA
 450 with wav2vec2 achieves the best on@100 (68.4%) with slightly
 451 higher median onset (62 ms). LFA NeMo exceeds standalone
 452 NFA Conformer on on@100 (46.5% vs. 45.3%) and achieves
 453 lower mean onset error (504 vs. 5,797 ms), confirming that
 454 phoneme DTW assignment improves alignment stability re-
 455 gardless of the per-chunk aligner.

4.5. Ablation Studies

We ablate two design choices: maximum chunk size (Table 4) and CTC alignment model (Table 5).

Table 4 shows the effect of maximum chunk size on TIMIT-LN-1h and TIMIT-LN-2h. All three aligners perform best with C=10, with the largest drop for wav2vec2 (−2.6 pp on@100 from C=10 to C=20 at 1 h) and smallest for NeMo (−0.7 pp at 1 h). The effect is moderate across all aligners (≤2.6 pp), indicating low sensitivity to chunk size within the 10–20 s range. We use C=10 as the default for all experiments.

Table 5 compares four wav2vec2 CTC models at C=10 on TIMIT-LN-2h. The 1B-parameter xlsr-1b-en achieves 65.6% on@100 with 66 ms median onset, outperforming xlsr-53-en (300M) by 5.5 pp. Among 300M-parameter models, cross-lingually pretrained xlsr-53 outperforms English-only wav2vec2-960h (60.1% vs. 55.9% on@100). We use xlsr-1b-en as the default for all other experiments.

Table 6 shows per-stage timing for LFA with wav2vec2 alignment. The real-time factor (RTF) is approximately 0.054 on the full 9.7-hour file. On the full file, ASR dominates the runtime (48%), followed by alignment (38%) and DTW including G2P conversion (11%). The DTW reference assignment step adds moderate overhead, with G2P conversion accounting for roughly two-thirds of that time.

5. Discussion

The core idea behind LFA is that ASR, despite being imperfect, provides enough information to decompose the long-speech alignment problem into locally solvable subproblems. Each short chunk can then be aligned with standard CTC-based or HMM-based methods, which are typically more stable on short segments than on very long recordings. Per-chunk alignment quality is primarily determined by the aligner’s short-segment performance, while any degradation on long files tends to arise from boundary effects and transcript-to-chunk assignment errors.

We chose DTW over Needleman-Wunsch or Smith-Waterman because all three DTW transitions (diagonal, horizontal, vertical) consume the local cell cost — there are no free gap operations. Every path position pays for mismatches, and ASR word splits/merges naturally correspond to many-to-one and one-to-many moves (cf. Figure 1). In Needleman-Wunsch,

Table 3: Detailed results on TIMIT-LN-full (1 file, 9 h 42 min, 54,357 words). CTC-Seg QuartzNet and MFA omitted.

System	Onset (ms)		Offset (ms)		On@ τ				Off@ τ		
	Mean	Med	Mean	Med	@25	@50	@100	@200	@50	@100	@200
<i>Baselines (standalone)</i>											
CTC-Seg (Conformer)	546.9	119	512.5	111	7.1	16.7	40.7	71.7	17.6	44.2	78.3
NFA (Citrinet)	6,500	117	6,490	87	10.5	21.4	42.9	74.0	31.1	55.0	75.9
NFA (Conformer)	5,797	108	5,789	68	7.9	18.1	45.3	80.5	40.1	61.5	78.3
<i>LFA (ours)</i>											
LFA (MFA, $C=10$)	738.9	52	746.8	58	32.9	49.3	63.5	70.9	46.2	62.0	69.7
LFA (wav2vec2, $C=10$)	673.5	62	693.2	61	17.3	39.5	68.4	74.7	40.0	69.8	74.2
LFA (NeMo Conf., $C=10$)	504.4	106	491.9	64	8.3	18.7	46.5	82.7	42.3	63.4	80.0

Table 4: Effect of maximum chunk size C on onset accuracy (%) on TIMIT-LN-1h and TIMIT-LN-2h.

Aligner	C	1 h		2 h	
		@50	@100	@50	@100
MFA	10	48.5	63.2	48.2	62.7
	15	47.2	61.3	47.0	61.0
	20	47.5	61.8	47.0	61.1
wav2vec2	10	37.7	66.8	36.8	65.6
	15	36.5	64.7	36.1	64.4
	20	36.0	64.2	35.4	63.6
NeMo Conf.	10	17.6	45.1	16.8	44.4
	15	16.9	44.3	16.4	43.3
	20	17.0	44.4	16.2	43.0

Table 5: CTC alignment model comparison on TIMIT-LN-2h ($C=10$).

CTC Model	Params	Onset (ms)		On@ τ		
		Med	Mean	@25	@50	@100
xlsr-1b-en	1B	66	690	16.1	36.8	65.6
xlsr-53-en	300M	79	723	12.9	29.0	60.1
fb-960h (base)	95M	83	735	14.0	29.8	57.2
fb-960h (large)	300M	87	745	10.4	25.7	55.9

497 gap penalties allow skipping words entirely — appropriate for
 498 sequence alignment but problematic when every reference word
 499 must be assigned to some chunk. True ASR insertions and dele-
 500 tions are handled by the cumulative assignment strategy, which
 501 aggregates reference words between chunk boundaries regard-
 502 less of path structure.

503 On noisy audio, MFA produces no output in default mode
 504 and degrades from 51% to near-random accuracy (5.5%) within
 505 20 min in single-speaker mode. The best CTC-based baselines
 506 (NFA Conformer, NFA Citrinet) range from 38–50% on@100
 507 and degrade with file length. LFA avoids both problems by
 508 operating on short, well-segmented chunks where any standard
 509 aligner performs well.

510 Since the decomposition is independent of the per-chunk
 511 aligner, any FA backend can be used on long recordings with-
 512 out modification. All three backends we evaluate — HMM-
 513 based MFA, wav2vec2 CTC, and NeMo Conformer Viterbi —
 514 maintain stable accuracy from 5 min to 9.7 h, with quality dif-

Table 6: Per-stage timing (LFA wav2vec2, $C=10$, A100 GPU). Per-file average in seconds. RTF = total / audio duration. Total includes I/O and initialization overhead not shown in individual stages.

Dataset	VAD	ASR	DTW	Align	Total	RTF
5 min	0.3	8.0	1.9	6.2	19.1	0.066
10 min	0.6	15.9	4.9	12.6	36.8	0.062
20 min	1.3	30.8	13.1	24.4	72.4	0.062
35 min	2.2	53.7	9.3	42.7	111.1	0.054
1 h	3.7	92.5	16.4	72.5	189.0	0.054
2 h	9.3	322.6	34.1	148.8	519.6	0.074
full	37.6	910.2	209.6	727.4	1,895	0.054

515 ferences reflecting each backend’s short-segment performance
 516 rather than any interaction with recording duration. MFA gives
 517 the best median onset (52 ms) and wav2vec2 gives better accu-
 518 racy at coarser thresholds (68.4% on@100).

519 The LFA NeMo vs. standalone NFA Conformer compar-
 520 ison is the cleanest ablation of chunking, since both use
 521 the same CTC model (`stt_en_conformer_ctc_large`) —
 522 LFA NeMo processes 10s chunks while NFA runs on the full
 523 recording. On *TIMIT-LN-full*, LFA NeMo achieves 11 \times lower
 524 mean onset error (504 vs. 5,797 ms) with consistent gains at
 525 all thresholds: on@100 (+1.2 pp), on@200 (+2.2 pp), off@50
 526 (+2.2 pp), and off@100 (+1.9 pp). Chunking eliminates large-
 527 scale drift while also improving per-word alignment at every
 528 threshold.

529 Our *TIMIT-LN-full* result (68.4% on@100 on a single 9.7 h
 530 file with 54,357 words) demonstrates that LFA scales to very
 531 long recordings. We use banded DTW with a band width of 5%
 532 of the matrix dimension, reducing memory from $O(N_{ASR} \times K)$
 533 to $O(N_{ASR} \times B)$ where $B \ll K$. The DTW step (including
 534 G2P conversion) takes 1.9s on 5-minute files, 16s on 1-hour
 535 files, and 210s on the 9.7-hour file — typically under 15% of
 536 total runtime. The system’s overall runtime is dominated by
 537 ASR and per-chunk alignment, both scaling linearly with audio
 538 duration.

539 LFA requires three language-dependent components: ASR,
 540 CTC alignment, and G2P. When MFA is used as the per-chunk
 541 aligner, a language-specific acoustic model and dictionary are
 542 additionally required. The CTC-only configuration avoids this
 543 dependency. LFA supports any language for which all three
 544 components are available; the overlap of Whisper (99 lan-
 545 guages), XLS-R wav2vec2 (over 50 fine-tuned), and Epitran
 546 (61 languages) covers several dozen languages, though we have

547 only evaluated on English.

548 LFA's reliance on ASR as an intermediate step introduces a
549 structural limitation: words present in the reference transcript
550 but absent from ASR output cannot be reliably assigned or
551 aligned. Whisper (our ASR backend) systematically omits filled
552 pauses and hesitations both by training convention (subtitle-
553 derived data) and explicit text normalization [33]. On conversa-
554 tional speech such as Buckeye [12], where fillers constitute
555 a substantial portion of annotated words, these tokens have no
556 corresponding ASR hypothesis; DTW assigns them based on
557 positional context alone, and CTC alignment is forced to place
558 them without acoustic evidence. The same mechanism affects
559 partial words, disfluencies, and interviewer speech that appears
560 in the audio but not in the ASR output. Whole-file acoustic
561 aligners such as NFA and MFA do not share this limitation, as
562 they align the reference directly against the audio signal using
563 models trained on such phenomena. Our TIMIT-LN benchmark
564 avoids this issue because TIMIT consists of read speech without
565 fillers, but applying LFA to conversational corpora would re-
566 quire a verbatim ASR system (e.g., CrisperWhisper [42]). For
567 aligning parliamentary and broadcast speech, where reference
568 transcripts are typically edited to remove disfluencies, this lim-
569 itation has minimal practical impact.

570 6. Conclusion

571 We presented LFA, a modular divide-and-conquer approach to
572 long-speech forced alignment that uses ASR to decompose the
573 problem into independently solvable chunks, with phoneme
574 DTW for reference transcript assignment. On noisy TIMIT
575 recordings spanning 5 min to 9.7 h, standalone MFA fails or de-
576 grades severely while the best CTC-based baselines range from
577 38–50% on@100. LFA maintains stable accuracy across all du-
578 rations: 68.4% on@100 (62 ms median onset) on a single 9.7-
579 hour file with wav2vec2, and 63.5% (52 ms median onset) with
580 MFA as per-chunk aligner. Any aligner capable of aligning 10 s
581 chunks can be used to process very long recordings by operating
582 on short, well-segmented chunks.

583 We release the LFA code as open-source and the benchmark
584 dataset.⁴

585 Future work includes evaluating LFA on conversational cor-
586 pora with overlapping speech and disfluencies and integrating
587 speech enhancement as a preprocessing stage to improve VAD
588 on severely degraded audio.

589 7. References

590 [1] M. McAuliffe, M. Socolof, S. Mihuc, M. Wagner, and M. Son-
591 deregger, "Montreal forced aligner: Trainable text-speech align-
592 ment using Kaldi," in *Proc. Interspeech*, 2017, pp. 498–502.

593 [2] R. Rousso, E. Cohen, J. Keshet, and E. Chodroff, "Tradition or
594 innovation: A comparison of modern ASR methods for forced
595 alignment," in *Proc. Interspeech*, 2024, pp. 1525–1529.

596 [3] P. J. Moreno, C. Joerg, J.-M. Van Thong, and O. Glickman, "A
597 recursive algorithm for the forced alignment of very long audio
598 segments," in *Proc. ICSLP*, 1998.

599 [4] T. J. Hazen, "Automatic alignment and error correction of human
600 generated transcripts for long speech recordings," in *Proc. Inter-
601 speech*, 2006.

602 [5] O. Aubert and J. Jäger, "Making parliamentary debates more
603 accessible: Aligning video recordings with text proceedings in
604 OpenParliamentTV," in *Proc. ParlaCLARIN IV Workshop*, 2024,
605 pp. 77–83.

⁴Code and data will be made available upon publication.

[6] J. Wirth and R. Peinl, "ASR Bundestag: A large-scale political
606 debate dataset in German," in *Proc. IntelliSys*. Springer, 2023,
607 pp. 190–202. 608

[7] N. Ljubešić, P. Rupnik, and D. Koržinek, "The ParlaSpeech col-
609 lection of automatically generated speech and text datasets from
610 parliamentary proceedings," in *Proc. SPECOM*. Springer, 2024,
611 pp. 137–150. 612

[8] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Lib-
613 riSpeech: An ASR corpus based on public domain audio books,"
614 in *Proc. ICASSP*, 2015, pp. 5206–5210. 615

[9] G. Doras, Y. Teytaut, and A. Roebel, "A linear memory CTC-
616 based algorithm for text-to-voice alignment of very long audio
617 recordings," *Applied Sciences*, vol. 13, no. 3, p. 1854, 2023. 618

[10] G. Bordel, S. Nieto, M. Penagarikano, L. J. Rodriguez-Fuentes,
619 and A. Varona, "A simple and efficient method to align very long
620 speech signals to acoustically imperfect transcriptions," in *Proc.
621 Interspeech*, 2012, pp. 1840–1843. 622

[11] A. Katsamanis, M. P. Black, P. G. Georgiou, L. Goldstein, and
623 S. Narayanan, "SailAlign: Robust long speech-text alignment," in
624 *Proc. Workshop on New Tools and Methods for Very-Large Scale
625 Phonetics Research*, 2011. 626

[12] M. A. Pitt, K. Johnson, E. Hume, S. Kiesling, and W. Raymond,
627 "The Buckeye corpus of conversational speech: Labeling conven-
628 tions and a test of transcriber reliability," *Speech Communication*,
629 vol. 45, no. 1, pp. 89–95, 2005. 630

[13] E. Rastorgueva, V. Lavrukhin, and B. Ginsburg, "NeMo forced
631 aligner and its application to word alignment for subtitle genera-
632 tion," in *Proc. Interspeech*, 2023, pp. 5257–5258. 633

[14] L. Kürzinger, D. Winkelbauer, L. Li, T. Watzel, and G. Rigoll,
634 "CTC-segmentation of large corpora for German end-to-end
635 speech recognition," in *Proc. Speech and Computer (SPECOM)*.
636 Springer, 2020, pp. 267–278. 637

[15] M. Bain, J. Huh, T. Han, and A. Zisserman, "WhisperX: Time-
638 accurate speech transcription of long-form audio," in *Proc. Inter-
639 speech*, 2023, pp. 4489–4493. 640

[16] A. Pettarin, "aeneas: automagically synchronize audio and text,"
641 <https://www.readbeyond.it/aeneas/>, 2017, version 1.7.3, open-
642 source tool. 643

[17] T. Kamp, "DSAlign: DeepSpeech based forced alignment tool,"
644 <https://github.com/mozilla/DSAlign>, 2020, mozilla, open-source
645 tool (unmaintained). 646

[18] A. Álvarez, H. Arzelus, and P. Ruiz, "Long audio alignment for
647 automatic subtitling using different phone-relatedness measures,"
648 in *Proc. ICASSP*, 2014, pp. 6280–6284. 649

[19] Y. Deng, F. Richardson, J. Steinberg, and P. Torres-Carrasquillo,
650 "Speech-to-text forced alignment benchmark," in *Proc. ISPA*,
651 2025. 652

[20] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, "Con-
653 nectionist temporal classification: Labelling unsegmented se-
654 quence data with recurrent neural networks," in *Proc. ICML*,
655 2006, pp. 369–376. 656

[21] R. Huang, X. Zhang, Z. Ni, L. Sun, M. Hira, J. Hwang,
657 V. Manohar, V. Pratap, M. Wiesner, S. Watanabe, D. Povey, and
658 S. Khudanpur, "Less peaky and more accurate CTC forced align-
659 ment by label priors," in *Proc. ICASSP*, 2024, pp. 11 831–11 835. 660

[22] V. Pratap, A. Tjandra, B. Shi, P. Tomasello, A. Babu, S. Kundu,
661 A. Elkahky, Z. Ni, A. Vyas, M. Fazel-Zarandi, A. Baevski, Y. Adi,
662 X. Zhang, W.-N. Hsu, A. Conneau, and M. Auli, "Scaling speech
663 technology to 1,000+ languages," *Journal of Machine Learning
664 Research*, vol. 25, no. 97, pp. 1–52, 2024. 665

[23] J. Zhu, C. Zhang, and D. Jurgens, "Phone-to-audio alignment
666 without text: A semi-supervised approach," in *Proc. ICASSP*,
667 2022, pp. 8167–8171. 668

[24] B. Mu, X. Shi, X. Wang, H. Liu, J. Xu, and L. Xie, "LLM-
669 ForcedAligner: A non-autoregressive and accurate LLM-based
670 forced aligner for multilingual and long-form speech," 2026. 671

- 672 [25] A. Rehman, J. Cai, J.-J. Zhang, and X. Yang, "BFA: Real-time
673 multilingual text-to-speech forced alignment," 2025.
- 674 [26] H. Sakoe and S. Chiba, "Dynamic programming algorithm op-
675 timization for spoken word recognition," *IEEE Transactions on*
676 *Acoustics, Speech, and Signal Processing*, vol. 26, no. 1, pp. 43–
677 49, 1978.
- 678 [27] M. Müller, *Fundamentals of music processing: Audio, analysis,*
679 *algorithms, applications.* Springer, 2015, vol. 5.
- 680 [28] I. González-Carrasco, L. Puente, B. Ruiz-Mezcua, and J. L.
681 López-Cuadrado, "Sub-Sync: Automatic synchronization of sub-
682 titles in the broadcasting of true live programs in Spanish," *IEEE*
683 *Access*, vol. 7, pp. 60968–60983, 2019.
- 684 [29] J. M. Masiello-Ruiz, B. Ruiz-Mezcua, P. Martinez, and
685 I. González-Carrasco, "Synchro-sub, an adaptive multi-algorithm
686 framework for real-time subtitling synchronisation of multi-type
687 TV programmes," *Computing*, vol. 105, pp. 1467–1495, 2023.
- 688 [30] K. Jin, S. Penke, and S. Alzubelli, "VoiceNet: Multilingual on-
689 device phoneme-to-audio alignment," in *Proc. Interspeech*, 2025,
690 pp. 2260–2264.
- 691 [31] H. Bredin, "pyannote.audio 2.1 speaker diarization pipeline: prin-
692 ciple, benchmark, and recipe," in *Proc. Interspeech*, 2023, pp.
693 1983–1987.
- 694 [32] G. Guillou, "Faster whisper transcription with CTranslate2,"
695 <https://github.com/SYSTRAN/faster-whisper>, 2023, software.
- 696 [33] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and
697 I. Sutskever, "Robust speech recognition via large-scale weak su-
698 pervision," in *Proc. ICML*, 2023, pp. 28492–28518.
- 699 [34] D. R. Mortensen, S. Dalmia, and P. Littell, "Epitrans: Precision
700 G2P for many languages," in *Proc. LREC*, 2018, pp. 2710–2714.
- 701 [35] F. J. Damerau, "A technique for computer detection and correction
702 of spelling errors," *Communications of the ACM*, vol. 7, no. 3, pp.
703 171–176, 1964.
- 704 [36] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0:
705 A framework for self-supervised learning of speech representa-
706 tions," in *Advances in Neural Information Processing Systems*,
707 vol. 33, 2020, pp. 12449–12460.
- 708 [37] A. Conneau, A. Baevski, R. Collobert, A. Mohamed, and M. Auli,
709 "Unsupervised cross-lingual representation learning for speech
710 recognition," in *Proc. Interspeech*, 2021, pp. 2426–2430.
- 711 [38] J. S. Garofolo, L. F. Lamel, W. M. Fisher, D. S. Pallett,
712 N. L. Dahlgren, V. Zue, and J. G. Fiscus, "TIMIT acoustic-
713 phonetic continuous speech corpus," Linguistic Data Consortium,
714 LDC93S1, 1993.
- 715 [39] R. Scheibler, E. Bezzam, and I. Dokmanić, "Pyroomacoustics: A
716 Python package for audio room simulation and array processing
717 algorithms," in *Proc. ICASSP*, 2018, pp. 351–355.
- 718 [40] D. Snyder, G. Chen, and D. Povey, "MUSAN: A music, speech,
719 and noise corpus," 2015.
- 720 [41] J. Carletta, S. Ashby, S. Bourban, M. Flynn, M. Guillemot,
721 T. Hain, J. Kadlec, V. Karaiskos, W. Kraaij, M. Kronenthal,
722 G. Lathoud, M. Lincoln, A. Lisowska, I. McCowan, W. Post,
723 D. Reidsma, and P. Wellner, "The AMI meeting corpus: A pre-
724 announcement," in *Proc. Machine Learning for Multimodal In-*
725 *teraction (MLMI)*. Springer, 2005, pp. 28–39.
- 726 [42] L. Wagner, B. Thallinger, and M. Zusag, "Crisperwhisper: Accu-
727 rate timestamps on verbatim speech transcriptions," in *Proc. In-*
728 *terspeech 2024*, 2024, pp. 1265–1269.