

# LLM-ForcedAligner: A Non-Autoregressive and Accurate LLM-Based Forced Aligner for Multilingual and Long-Form Speech

Anonymous ACL submission

## Abstract

Forced alignment (FA) predicts start and end timestamps for words or characters in speech, but existing methods are language-specific and prone to cumulative temporal shifts. The multilingual speech understanding and long-sequence processing abilities of speech large language models (SLLMs) make them promising for FA in multilingual, crosslingual, and long-form speech settings. However, directly applying the next-token prediction paradigm of SLLMs to FA results in hallucinations and slow inference. To bridge the gap, we propose **LLM-ForcedAligner**, reformulating FA as a slot-filling paradigm: timestamps are treated as discrete indices, and special timestamp tokens are inserted as slots into the transcript. Conditioned on the speech embeddings and the transcript with slots, the SLLM directly predicts the time indices at slots. During training, causal attention masking with non-shifted input and label sequences allows each slot to predict its own timestamp index based on itself and preceding context, with loss computed only at slot positions. Dynamic slot insertion enables FA at arbitrary positions. Moreover, non-autoregressive inference is supported, avoiding hallucinations and improving speed. Experiments across multilingual, crosslingual, and long-form speech scenarios show that LLM-ForcedAligner achieves a 69%~78% relative reduction in accumulated averaging shift compared with prior methods. The checkpoint and inference code will be released later.

## 1 Introduction

In speech processing, the objective of forced alignment (FA) is to estimate the start and end timestamps of each word or character in a speech signal, given its corresponding transcript. FA is indispensable in numerous applications, including large-scale speech corpus construction and cleaning, automatic subtitling and word-level highlighting, as well as duration modeling and prosody analysis in

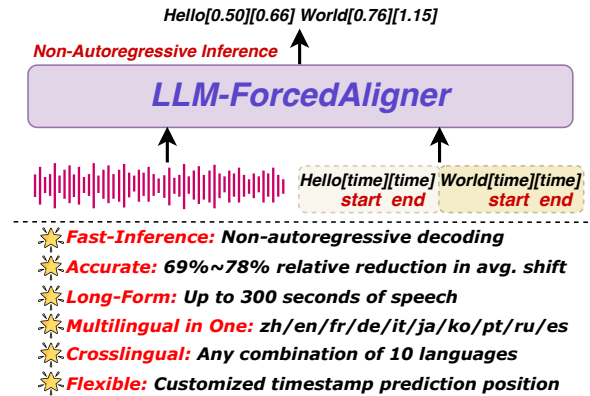


Figure 1: Highlights of LLM-ForceAligner.

speech synthesis. With the continuous advancement of multilingual and multimodal applications, efficient and accurate FA has become increasingly essential (Tseng et al., 2021; Liu et al., 2025).

Existing FA methods can be broadly categorized into two main groups: traditional hybrid systems (McAuliffe et al., 2017) and end-to-end models (Rastorgueva et al., 2023; Kürzinger et al., 2020; Shi et al., 2023; Bain et al., 2023). Montreal forced aligner (MFA) is typically a hybrid Gaussian Mixture Model–Hidden Markov Model (GMM–HMM) framework, computing the timestamps via Viterbi decoding for frame-level phoneme-to-text alignment paths (McAuliffe et al., 2017). Connectionist temporal classification (CTC) is a common end-to-end FA method that leverages frame-to-token alignment computed by CTC-based automatic speech recognition (ASR) models, employing dynamic programming to identify the optimal path that aligns with the text sequence within constrained search paths (Kürzinger et al., 2020; Rastorgueva et al., 2023). Continuous integrate-and-fire (CIF) predicts a weight for each encoder output frame and integrates these weights over time. When the accumulated weight exceeds a threshold, a fire event is triggered, at which point a weighted

070 sum of the accumulated frame-level acoustic vec- 122  
071 tors is computed to generate an acoustic embed- 123  
072 ding aligned with an output token, thereby en- 124  
073 abling the assignment of a corresponding times- 125  
074 tamp to each token (Shi et al., 2023). WhisperX 126  
075 employs a lightweight end-to-end phoneme recog- 127  
076 nition model to perform frame-level phoneme clas- 128  
077 sification on speech, and then aligns the result- 129  
078 ing phoneme sequence with the transcript using 130  
079 dynamic time warping (DTW), thereby obtaining 131  
080 word-level timestamps by aggregating phoneme- 132  
081 level timestamps (Bain et al., 2023). 133

082 However, the aforementioned FA methods are 134  
083 tied to language-specific phonemes, lexicons, or 135  
084 structural designs, which means that in multilin- 136  
085 gual scenarios, deployment typically involves a col- 137  
086 lection of independent systems with disparate struc- 138  
087 tures, leading to engineering costs and maintenance 139  
088 complexity that grow linearly with the number of 140  
089 languages. Furthermore, previous FA methods can 141  
090 be summarized as a process of calculating local 142  
091 acoustic similarities followed by a monotonic path 143  
092 search. While these methods can produce reason- 144  
093 ably accurate boundaries for short segments, they 145  
094 frequently accumulate significant systematic tem- 146  
095 poral shifts in long-form speech. 147

096 Large language models (LLMs) have demon- 148  
097 strated powerful abilities in multilingual text un- 149  
098 derstanding and long-sequence processing tasks (Tou- 150  
099 vron et al., 2023a,b; Yang et al., 2024, 2025), offer- 151  
100 ing a new possibility for FA that supports multilin- 152  
101 gual, crosslingual, and long-form speech. Increas-  
102 ing studies have explored integrating speech en-  
103 coders with LLMs to build Speech LLMs (SLLMs)  
104 that process speech and text within a unified frame-  
105 work. Nevertheless, existing SLLMs have mainly  
106 achieved success in high-level semantic tasks such  
107 as ASR (Geng et al., 2024; Mu et al., 2026), speech  
108 understanding (Geng et al., 2025b; Chu et al., 2023,  
109 2024), speech synthesis (Wang et al., 2025a; Du  
110 et al., 2025), and spoken dialogue (Wang et al.,  
111 2025b; Geng et al., 2025a). For FA, which is more  
112 sensitive to acoustic characteristics, these SLLMs  
113 typically treat it as a by-product of ASR and gen-  
114 erate word-level or character-level timestamps via  
115 next-token prediction. Such a paradigm is suscepti-  
116 ble to temporal non-monotonic hallucinations and  
117 incurs substantial inference latency.

118 In this work, we propose a new FA framework  
119 named **LLM-ForcedAligner** that reformulates FA  
120 as a slot-filling paradigm: the start and end times-  
121 tamps of each word or character are treated as dis-

crete time indices, and dedicated special tokens  
are inserted into the transcript as slots so that, con-  
ditioned on the speech embeddings and the tran-  
script augmented with these slots, the SLLM can  
directly predict the corresponding time indices at  
the designated slots. This new FA paradigm ef-  
fectively leverages the LLM’s strengths in slot fill-  
ing and long-context processing, extends conven-  
tional purely acoustic, phoneme-level alignment to  
semantic-boundary-aware, character-level or word-  
level alignment. To perform slot filling, we apply  
causal attention masking during training without  
introducing any shift between the input and label se-  
quences, allowing each slot to predict its own time  
index based on itself and the preceding context, and  
we compute the loss function only at the slot posi-  
tions. In addition, we adopt a dynamic slot inser-  
tion strategy that randomly decides whether to in-  
sert special tokens for each word or character in the  
transcript, enabling LLM-ForcedAligner to predict  
timestamps for arbitrary words or characters. Dur-  
ing inference, LLM-ForcedAligner enables non-  
autoregressive decoding, completely avoiding hal-  
lucinations compared with autoregressive decod-  
ing and achieving faster speed. Experimental re-  
sults show that, in multilingual, cross-lingual, and  
long-form speech scenarios of up to 300 seconds,  
LLM-ForcedAligner achieves a 69%~78% relative  
reduction in accumulated averaging shift (AAS)  
compared with prior FA methods, while incurring  
only a slight increase in real-time factor (RTF).

## 2 Related Work 153

**Traditional Hybrid Systems for FA.** These meth- 154  
ods (McAuliffe et al., 2017) build GMM–HMM 155  
acoustic models based on Kaldi (Povey et al., 2011), 156  
map text to phoneme sequences using lexicons and 157  
pronunciation dictionaries, and then apply Viterbi 158  
decoding to search for the most probable frame- 159  
level alignment path under constrained transitions. 160  
Owing to the inherent monotonic temporal struc- 161  
ture of HMMs and their fine-grained phoneme mod- 162  
eling, such methods typically achieve high accu- 163  
racy at phoneme-level boundaries. However, they 164  
rely heavily on language-specific lexicons, pronun- 165  
ciation dictionaries, and phoneme sets, with each 166  
language typically requiring a separately trained or 167  
adapted GMM–HMM model. 168

**End-to-End Models for FA.** These methods no 169  
longer rely on explicit lexicon mapping; instead, 170  
they use the CTC or internal attention weights 171

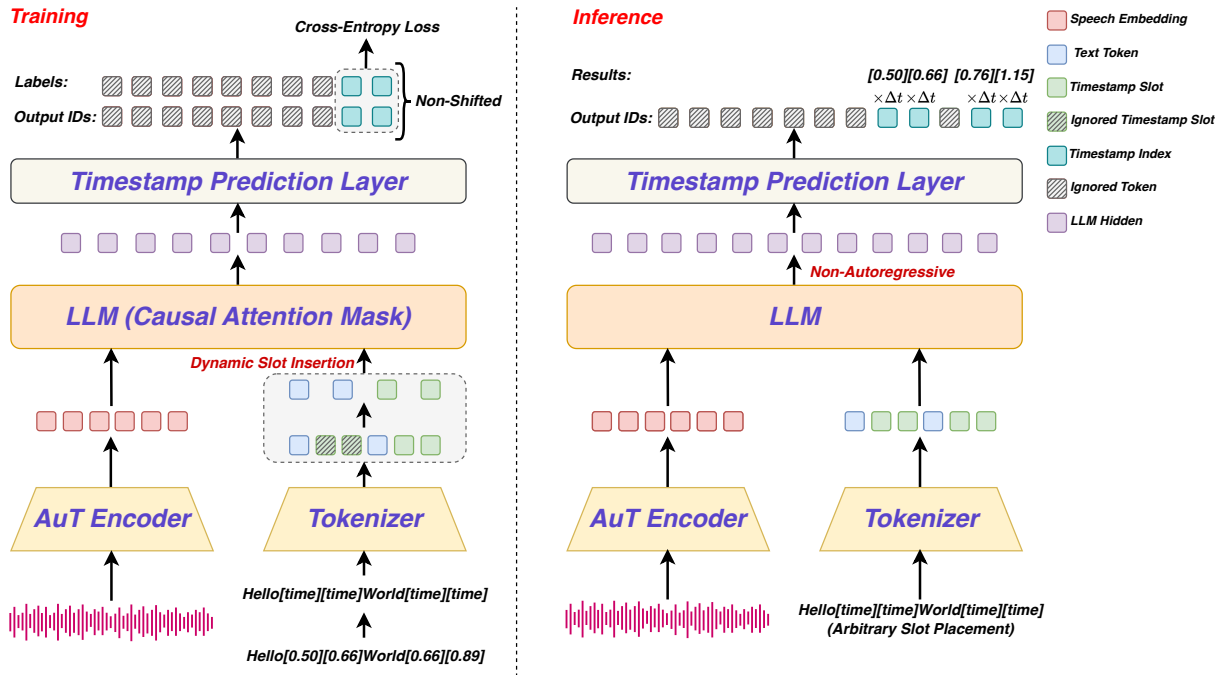


Figure 2: Overview of LLM-ForcedAligner. **Left:** During training, we replace the word-level or character-level start and end timestamps in the transcript with a special token `[time]` to serve as slots, and use dynamic slot insertion to randomly determine which word or character slots to ignore. The speech embeddings produced by the AuT encoder are then concatenated and fed into the LLM for training with causal attention masking. When computing the cross-entropy loss, the output IDs and the labels are non-shifted, and the loss is computed only at the positions of the retained slots. **Right:** During inference, users can insert the special token `[time]` at arbitrary positions in the transcript and rapidly obtain the corresponding start and end timestamps through non-autoregressive decoding.

172 to learn a direct mapping between acoustic fea- 172  
 173 tures and phoneme sequences, achieving align- 173  
 174 ment by computing emission probabilities. How- 174  
 175 ever, they require adaptation for each language 175  
 176 and tend to produce systematic temporal shifts. 176  
 177 CTC-based NeMo forced aligner (NFA) (Rastorgueva 177  
 178 et al., 2023) requires language-specific lexicons, 178  
 179 and multilingual CTC models suffer from trade- 179  
 180 offs between different languages. NFA relies on 180  
 181 sparse peak frames, which often lag the actual on- 181  
 182 set of speech, leading to overall backward shifts 182  
 183 in token boundaries that accumulate into signifi- 183  
 184 cant systemic temporal shifts in long-form speech. 184  
 185 CIF (Shi et al., 2023) couples with language charac- 185  
 186 teristics, resulting in much lower FA performance 186  
 187 in word-boundary languages like English than in 187  
 188 character-boundary languages like Chinese. In CIF, 188  
 189 each token is typically accumulated over approxi- 189  
 190 mately a fixed number of encoder frames rather 190  
 191 than adapting to its actual duration, leading to sys- 191  
 192 tematically delayed end times for long vowels, long 192  
 193 syllables, or long words. Although post-processing 193  
 194 such as scaled-CIF and fire-delay can partially 194  
 195 alleviate this issue, such structural shifts remain 195

196 difficult to eliminate in long-form speech. Whis- 196  
 197 perX (Bain et al., 2023) requires switching between 197  
 198 phoneme recognition models to accommodate dif- 198  
 199 ferent languages, and remains essentially based on 199  
 200 DTW with local similarity. In long-form speech, lo- 200  
 201 cal alignment errors can similarly accumulate into 201  
 202 global temporal shifts. 202  
 203 **SLLMs for FA.** Existing SLLMs (Chu et al., 2023; 203  
 204 Geng et al., 2025b) treat FA as a by-product of 204  
 205 ASR, processing timestamps as natural language, 205  
 206 which prevents guaranteed temporal monotonicity 206  
 207 during autoregressive inference and can result in 207  
 208 excessive timestamps due to LLM hallucinations; 208  
 209 in long-form speech scenarios, this significantly 209  
 210 increases inference time. Therefore, the next-token 210  
 211 prediction paradigm of conventional SLLM train- 211  
 212 ing and inference is not well-suited for FA, whereas 212  
 213 LLM-ForcedAligner restructures the FA, enabling 213  
 214 accurate and fast timestamp prediction. 214

### 3 LLM-ForcedAligner 215

#### 3.1 Overall Architecture 216

217 LLM-ForcedAligner formulates FA as a slot-filling 217  
 218 paradigm: given the speech signal and a transcript 218

augmented with special tokens  $[time]$  that denote word-level or character-level start and end time slots, the SLLM directly predicts the corresponding discrete timestamp indices for each slot. Unlike previous FA methods, which first perform frame-level or phoneme-level alignment and then aggregate the results into word-level or character-level timestamps, LLM-ForcedAligner directly predicts word-level or character-level timestamp indices.

Training LLM-ForcedAligner requires word-level or character-level timestamp labels for a large number of speech–transcript pairs; however, because manual annotation is prohibitively expensive, we adopt pseudo-timestamp labels generated by MFA, which is the most accurate among existing FA methods. It is important to emphasize that MFA pseudo-labels inherently contain noise and systematic shifts. LLM-ForcedAligner does not merely replicate the MFA outputs; instead, it distills and smooths these pseudo-labels using the SLLM, achieving more stable, lower-shift timestamp predictions. Given a speech signal  $x$  and corresponding transcript sequence  $w = (w_1, \dots, w_N)$ , MFA is used to obtain the start and end times of each word or character and incorporate them into the transcript, resulting in a sequence  $wt = (w_1, t_{1,start}^{MFA}, t_{1,end}^{MFA}, \dots, w_N, t_{N,start}^{MFA}, t_{N,end}^{MFA})$ .

The speech encoder in LLM-ForcedAligner is from the Audio Transformer (AuT), and its output speech embeddings based on speech signal  $x$  are sampled at 12.5Hz, meaning that each speech embedding frame corresponds to 80ms of the speech signal (Xu et al., 2025). Before feeding  $wt$  into the tokenizer, we replace all timestamps with the special token  $[time]$  to represent slots, resulting in the input transcript sequence  $ws = (w_1, [time], [time], \dots, w_N, [time], [time])$ . Moreover, we discretize the timestamps in the sequence  $wt$  into timestamp indices:

$$\tau_{i,start} = \left\lfloor \frac{t_{i,start}^{MFA}}{\Delta t} \right\rfloor, \tau_{i,end} = \left\lfloor \frac{t_{i,end}^{MFA}}{\Delta t} \right\rfloor, \quad (1)$$

where  $\Delta t$  is 80ms chosen to align with the frame rate of the AuT encoder, resulting in a transcript sequence with discrete timestamp indices  $wl = (w_1, \tau_{1,start}, \tau_{1,end}, \dots, w_N, \tau_{N,start}, \tau_{N,end})$ . The speech embeddings from the AuT encoder and the text embeddings from  $ws$  are fed into the LLM, and a linear layer with 3,750 (i.e., 500s/80ms) output classes is used to predict timestamp indices for the entire input sequence.

The multilingual and crosslingual capabilities of LLM-ForcedAligner are jointly provided by the speech encoder and the multilingual LLM. Specifically, the AuT encoder, pre-trained on a large-scale multilingual corpus, generates effective frame-level speech embeddings for multiple languages, while the multilingual LLM handles semantic information across different languages. Moreover, the special token  $[time]$  and timestamp prediction layer do not rely on language-specific phoneme sets or lexicons. Therefore, LLM-ForcedAligner can process multilingual and crosslingual speech–transcript pairs, overcoming the language-specific limitations of previous FA methods.

### 3.2 Training Strategy

SLLMs commonly adopt a training scheme in which the last token of the output ID sequence and the first token of the label sequence are removed, creating a one-position shift between the two sequences; the cross-entropy loss is then computed, achieving the standard next-token prediction paradigm. However, this paradigm is not suitable for filling timestamp slots. Instead, we adopt causal training with the output ID and label sequences non-shifted, enabling the LLM-ForcedAligner to explicitly perceive timestamp slots during training and predict the timestamps to be filled into those slots. Moreover, causal training enables the LLM-ForcedAligner to incorporate prior context when predicting the timestamp for the current slot, ensuring global consistency in timestamp prediction.

During training, we compute the cross-entropy loss only at timestamp slot positions, thereby focusing the training objective of LLM-ForcedAligner on timestamp slot filling. Given the joint input sequence  $x$  and  $ws$ , the LLM produces a hidden-state sequence  $h$  under causal attention masking. For each timestamp slot position  $j \in \mathcal{I}_{ts}$ , a discrete timestamp index distribution is obtained from the timestamp prediction layer:

$$p(\hat{\tau}_j | x, ws) = \text{softmax}(\text{TPL}(h_j)), \quad (2)$$

where TPL denotes the timestamp prediction layer, and position  $j$  corresponds to a start or end timestamp slot for an arbitrary word or character. The loss function of LLM-ForcedAligner is defined as:

$$\mathcal{L} = -\frac{1}{|\mathcal{I}_{ts}|} \sum_{j \in \mathcal{I}_{ts}} \log p(\hat{\tau}_j = \tau_j | x, ws), \quad (3)$$

where  $\tau_j$  is the discrete timestamp index from  $wl$ .

Furthermore, consistently inserting start and end timestamp slots for every word or character during training would cause LLM-ForcedAligner to rely excessively on previously predicted timestamps. We propose a dynamic slot insertion strategy during training to enhance the generalization capability of LLM-ForcedAligner. Specifically, for each sample, we determine with a 50% probability whether to apply dynamic slot insertion. When dynamic slot insertion is applied, each word or character in the sample has a 50% chance of having start and end timestamp slots inserted after it. This strategy continuously varies the set of timestamp slot positions  $\mathcal{I}_{ts}$ , enabling LLM-ForcedAligner to predict start and end timestamps for words or characters at arbitrary positions.

### 3.3 Non-Autoregressive Inference

As the output ID and label sequences are kept non-shifted during training, LLM-ForcedAligner can predict the timestamp indices for all slots in a transcript simultaneously using a non-autoregressive decoding. Specifically, for a speech–transcript pair, users can customize start and end timestamp slots after any word or character. Given a user-defined timestamp slot position  $j \in \mathcal{I}_{ts}$ , LLM-ForcedAligner predicts its timestamp index via non-autoregressive decoding, which is then converted to a millisecond-level timestamp  $\hat{t}_j$ :

$$\hat{t}_j = \hat{\tau}_j \cdot \Delta t, \quad (4)$$

where  $\Delta t$  is 80ms.

## 4 Experimental Setup

### 4.1 Dataset

We conduct our experiments on 56,000 hours of data containing 10 languages: Chinese, English, French, German, Italian, Japanese, Korean, Portuguese, Russian, and Spanish. These data are drawn from a combination of internal resources and open-source datasets, covering a wide range of scenarios such as read speech, conversational speech, podcasts, and meetings, including Aishell (Bu et al., 2017; Du et al., 2018), WenetSpeech (Zhang et al., 2022), Aidatatang<sup>1</sup>, Magicdata<sup>2</sup>, KeSpeech (Tang et al., 2021), LibriSpeech (Panayotov et al., 2015), GigaSpeech (Chen et al., 2021), LibriTTS (Zen et al., 2019), Emilia (He et al., 2024), and MLS (Pratap et al., 2020). All training and test

<sup>1</sup><http://openslr.magicdatatech.com/62/>

<sup>2</sup><https://www.openslr.org/68/>

datasets are annotated with pseudo-timestamps generated by MFA, and we additionally assess the performance of LLM-ForcedAligner on an internal Chinese test dataset with manually annotated timestamps. The transcripts in the training dataset are obtained from either manual annotations or ASR model predictions, which enhances the generalization ability of LLM-ForcedAligner to varying transcript quality. Details of the training and test datasets are provided in the A.1.

### 4.2 Implementation Details

The AuT encoder in LLM-ForcedAligner contains 316.42M parameters and is initialized from the AuT encoder of Qwen3-Omni<sup>3</sup>. The LLM uses Qwen3-0.6B<sup>4</sup>, and the timestamp prediction layer is a single linear layer with 3,750 output timestamp classes and contains 3.84M parameters. During training, the AuT encoder, the LLM, and the timestamp prediction layer are jointly optimized using the Adam optimizer with a warm-up scheduler, which increases the learning rate to a peak of 0.0003 after 1,000 steps.

### 4.3 Evaluation Metrics

We use accumulated averaging shift (AAS) to measure the performance of timestamp prediction (Shi et al., 2023). Lower AAS indicates better timestamp prediction performance. Specifically, AAS computes the average shift for each timestamp slot, defined as the mean absolute difference between the predicted timestamps and the ground-truth timestamps across all slots in the test dataset:

$$AAS = \frac{\sum_{k=1}^K |\hat{k}_i - k_i|}{K}, \quad (5)$$

where  $K$  is the total number of timestamp slots in the test dataset,  $\hat{k}_i$  is computed according to Eq. 4, and  $k_i$  can be timestamps obtained from MFA or manual annotations.

## 5 Experimental Results

### 5.1 Main Results

Table 1 and Table 2 compares LLM-ForcedAligner with other FA methods on MFA-labeled test datasets across multilingual, crosslingual, and long-form speech scenarios. Monotonic-Aligner<sup>5</sup> sup-

<sup>3</sup><https://huggingface.co/Qwen/Qwen3-Omni-30B-A3B-Thinking>

<sup>4</sup><https://huggingface.co/Qwen/Qwen3-0.6B>

<sup>5</sup>[https://modelscope.cn/models/iic/speech\\_timestamp\\_prediction-v1-16k-offline](https://modelscope.cn/models/iic/speech_timestamp_prediction-v1-16k-offline)

Table 1: AAS (ms) ↓ of LLM-ForcedAligner and other FA methods on **MFA-labeled** test datasets. The test dataset for each language consists of the raw speech from both the open-source and internal test datasets of that language.

| Language   | Monotonic-Aligner | NFA   | WhisperX | LLM-ForcedAligner |
|------------|-------------------|-------|----------|-------------------|
| Chinese    | 161.1             | 109.8 | -        | <b>33.1</b>       |
| English    | -                 | 107.5 | 92.1     | <b>37.5</b>       |
| French     | -                 | 100.7 | 145.3    | <b>41.7</b>       |
| German     | -                 | 122.7 | 165.1    | <b>46.5</b>       |
| Italian    | -                 | 142.7 | 155.5    | <b>75.5</b>       |
| Japanese   | -                 | -     | -        | <b>42.4</b>       |
| Korean     | -                 | -     | -        | <b>37.2</b>       |
| Portuguese | -                 | -     | -        | <b>38.4</b>       |
| Russian    | -                 | 200.7 | -        | <b>40.2</b>       |
| Spanish    | -                 | 124.7 | 108.0    | <b>36.8</b>       |
| Avg.       | 161.1             | 129.8 | 133.2    | <b>42.9</b>       |

Table 2: AAS (ms) ↓ of LLM-ForcedAligner and other FA methods on **MFA-labeled** test datasets. The test dataset for each language consists of concatenated speech up to 300 seconds in duration from the raw speech in each language’s open-source and internal test datasets. “Mixed-Crosslingual” is the concatenated speech with a duration of up to 300 seconds from arbitrary languages’ open-source and internal test datasets.

| Language           | Monotonic-Aligner | NFA   | WhisperX | LLM-ForcedAligner |
|--------------------|-------------------|-------|----------|-------------------|
| Chinese            | 1742.4            | 235.0 | -        | <b>36.5</b>       |
| English            | -                 | 226.7 | 227.2    | <b>58.6</b>       |
| French             | -                 | 230.6 | 2052.2   | <b>53.4</b>       |
| German             | -                 | 220.3 | 993.4    | <b>62.4</b>       |
| Italian            | -                 | 290.5 | 5719.4   | <b>81.6</b>       |
| Japanese           | -                 | -     | -        | <b>81.3</b>       |
| Korean             | -                 | -     | -        | <b>42.2</b>       |
| Portuguese         | -                 | -     | -        | <b>50.0</b>       |
| Russian            | -                 | 283.3 | -        | <b>43.0</b>       |
| Spanish            | -                 | 240.2 | 4549.9   | <b>39.6</b>       |
| Mixed-Crosslingual | -                 | -     | -        | <b>34.2</b>       |
| Avg.               | 1742.4            | 246.7 | 2708.4   | <b>52.9</b>       |

ports only Chinese, while NFA<sup>6</sup> and WhisperX<sup>7</sup> require switching models for different languages. Details of compared FA methods are provided in the A.2. We observe that LLM-ForcedAligner not only supports multiple languages without requiring model switching, but also achieves a 66%~73% relative reduction in average AAS on multilingual raw speech compared with other FA methods. Moreover, LLM-ForcedAligner achieves extremely low average AAS on multilingual and crosslingual long-form speech, a setting where other FA methods fail to perform effectively.

<sup>6</sup>[https://github.com/NVIDIA-NeMo/NeMo/tree/main/tools/nemo\\_forced\\_aligner](https://github.com/NVIDIA-NeMo/NeMo/tree/main/tools/nemo_forced_aligner)

<sup>7</sup><https://github.com/m-bain/whisperX>

Table 3 compares LLM-ForcedAligner with other FA methods on human-labeled Chinese test datasets, covering noisy, crosslingual, and long-form speech scenarios. We find that the average AAS of LLM-ForcedAligner on the human-labeled test datasets achieves a 68%~78% relative reduction compared with other FA methods, showing that training LLM-ForcedAligner with MFA-labeled data generalizes well to real-world conditions.

In addition, in Table 2, the average AAS for long-form speech on the MFA-labeled test datasets is slightly higher than that for raw speech, reflecting the systematic shifts of MFA on long-form speech. In contrast, in Table 3, on the human-labeled test datasets, the average AAS for long-

403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414

415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429

Table 3: AAS (ms) ↓ of LLM-ForcedAligner and other FA methods on **human-labeled** test datasets. “Raw” is the raw speech in the dataset; “Raw-Noisy” is the raw speech with added background noise; “Mixed-60s” and “Mixed-300s” are concatenations of raw speech into durations with maximum lengths of 60 and 300 seconds; “Mixed-Crosslingual” is a concatenation of human-labeled raw speech and MFA-labeled multilingual speech.

| Type               | Monotonic-Aligner | NFA   | LLM-ForcedAligner |
|--------------------|-------------------|-------|-------------------|
| Raw                | 49.9              | 88.6  | <b>27.8</b>       |
| Raw-Noisy          | 53.3              | 89.5  | <b>41.8</b>       |
| Mixed-60s          | 51.1              | 86.7  | <b>25.3</b>       |
| Mixed-300s         | 410.8             | 140.0 | <b>24.8</b>       |
| Mixed-Crosslingual | -                 | -     | <b>42.5</b>       |
| Avg.               | 141.3             | 101.2 | <b>32.4</b>       |

Table 4: The average RTF ↓ of LLM-ForcedAligner and other compared FA methods during inference.

| Methods           | Avg. RTF |
|-------------------|----------|
| Monotonic-Aligner | 0.0079   |
| NFA               | 0.0067   |
| WhisperX          | 0.0113   |
| LLM-ForcedAligner | 0.0159   |

Table 5: AAS (ms) ↓ on **MFA-labeled** and **human-labeled** test datasets for the ablation study of different timestamp token durations.. “Raw” is raw speech, while “Mixed” is monolingual and crosslingual concatenated speech up to 300 seconds.

| Timestamp Dur. | MFA-Labeled |             | Human-Labeled |             |
|----------------|-------------|-------------|---------------|-------------|
|                | Raw         | Mixed       | Raw           | Mixed       |
| 120ms          | 51.1        | 62.9        | 35.5          | 34.5        |
| 80ms           | 41.7        | 52.9        | <b>27.8</b>   | <b>25.1</b> |
| 40ms           | <b>34.0</b> | <b>50.9</b> | 32.5          | 27.9        |

form speech is lower than that for raw speech, indicating that LLM-ForcedAligner does not simply replicate MFA’s timestamp predictions, but instead learns a more robust and reliable timestamp prediction that can effectively correct MFA labels in long-form scenarios. Furthermore, during long-form inference, LLM-ForcedAligner can leverage a longer historical context to predict timestamps for the current slots, resulting in superior performance on human-labeled long-form speech test datasets.

Table 4 reports the average RTF of LLM-ForcedAligner and other compared FA methods under identical inference conditions. As the number of model parameters increases, the RTF shows a slight increase. Due to the benefit of the non-autoregressive inference of LLM-ForcedAligner, it achieves a substantial reduction in average AAS with only a minimal increase in RTF. Users can select the most suitable FA method based on the average AAS-RTF trade-off.

## 5.2 Ablation Study

Table 5 shows the average AAS results of LLM-ForcedAligner trained with different timestamp token durations on the MFA-labeled and human-labeled test datasets. When the timestamp token duration is 120ms, the timestamp prediction layer has 2,500 classes (i.e., 300s/120ms); when the dura-

tion is 80ms, it has 3,750 classes (i.e., 300s/80ms); and when the duration is 40ms, it has 7,500 classes (i.e., 300s/40ms). As the timestamp token duration decreases, the AAS on the MFA-labeled test datasets steadily declines, indicating that finer-grained timestamp prediction better fits the MFA labels. However, on human-labeled test datasets, finer-grained timestamp prediction does not yield lower AAS because it better fits the MFA timestamp distribution, leading to reduced generalization. The timestamp token duration of 80ms is the optimal choice, as each frame of the AuT encoder’s output also represents 80ms of speech, which helps LLM-ForcedAligner better determine the start and end timestamps of words or characters based on speech boundaries.

Table 6 shows the average AAS results of LLM-ForcedAligner on the MFA-labeled and human-labeled test datasets, comparing results with and without dynamic slot insertion during training. Dynamic slot insertion randomly determines whether to insert timestamp slots after each word or character, enabling LLM-ForcedAligner to predict start and end timestamps at arbitrary positions and preventing it from relying excessively on previously

Table 6: AAS (ms) ↓ of the ablation study on dynamic slot insertion across **MFA-labeled** and **human-labeled** test datasets. “Raw” is raw speech, while “Mixed” is monolingual and crosslingual concatenated speech up to 300 seconds.

| Dynamic Slot Insertion | MFA-Labeled |             | Human-Labeled |             |
|------------------------|-------------|-------------|---------------|-------------|
|                        | Raw         | Mixed       | Raw           | Mixed       |
| w/o                    | 51.2        | 66.1        | 31.4          | 30.4        |
| w/                     | <b>41.7</b> | <b>52.9</b> | <b>27.8</b>   | <b>25.1</b> |

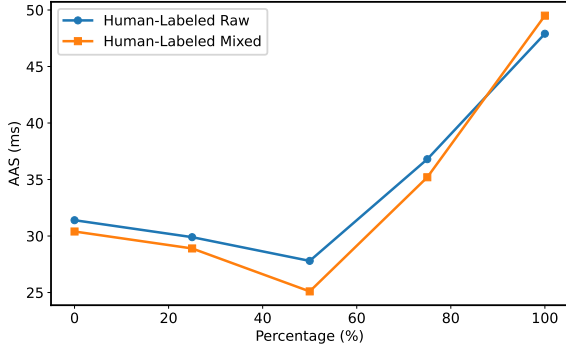


Figure 3: AAS (ms) of LLM-ForcedAligner on the **human-labeled** test datasets for different dynamic slot insertion percentages.

482 predicted timestamps. We find that dynamic slot  
483 insertion reduces AAS on both test datasets, with  
484 the improvement more pronounced for long-form  
485 speech. This phenomenon is because dynamic slot  
486 insertion, by randomly deciding whether to insert  
487 timestamp slots after each word or character, pre-  
488 vents LLM-ForcedAligner from excessively rely-  
489 ing on historically predicted timestamps, which  
490 can otherwise lead to systematic temporal shifts.  
491 Furthermore, dynamic slot insertion enables LLM-  
492 ForcedAligner to predict start and end timestamps  
493 for words or characters at arbitrary positions, sup-  
494 porting user-customizable timestamp prediction.

### 495 5.3 Visualization

496 Figure 3 shows the AAS results on the human-  
497 labeled test datasets when LLM-ForcedAligner is  
498 trained with different percentages of dynamic slot  
499 insertion. When the percentage of dynamic slot  
500 insertion is below 50% of the training samples,  
501 LLM-ForcedAligner achieves lower AAS, and the  
502 AAS continues to decrease as the percentage in-  
503 creases. However, when the percentage exceeds  
504 50% of the training samples, the AAS begins to  
505 increase, reaching its highest value at 100%. There-  
506 fore, selecting an appropriate dynamic slot inser-  
507 tion percentage of 50% is crucial for enhancing the

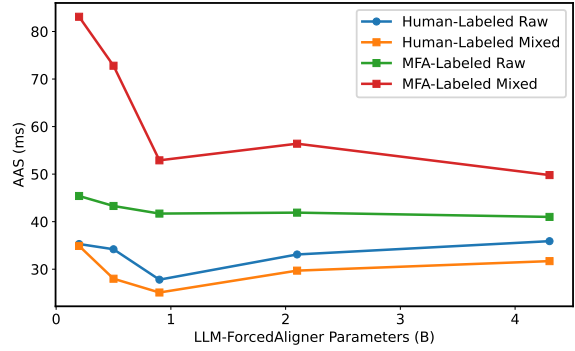


Figure 4: AAS (ms) on the **MFA-labeled** and **human-labeled** test datasets for LLM-ForcedAligner with different parameter settings.

508 generalization of LLM-ForcedAligner.

509 Figure 4 shows the AAS results on the MFA-  
510 labeled and human-labeled test datasets for LLM-  
511 ForcedAligner with different parameter settings.  
512 When the LLM-ForcedAligner parameter size is  
513 below 0.9B, timestamp prediction performance is  
514 limited by insufficient model capacity. When the  
515 parameter size exceeds 0.9B, the AAS on the MFA-  
516 labeled test datasets shows no significant change,  
517 while the AAS on the human-labeled test datasets  
518 increases, indicating that LLM-ForcedAligner over-  
519 fits the MFA timestamp distribution. A param-  
520 eter size of 0.9B is therefore optimal. At this  
521 scale, LLM-ForcedAligner does not strictly fit the  
522 MFA timestamp distribution, but instead learns a  
523 smoother and more robust timestamp prediction  
524 behavior with better generalization performance.

## 525 6 Conclusion

526 We propose **LLM-ForcedAligner**, a reformulation  
527 of FA as a slot-filling paradigm: timestamps are  
528 treated as discrete indices, and special timestamp  
529 tokens are inserted as slots into the transcript. Con-  
530 ditioned on the speech and the transcript with slots,  
531 LLM-ForcedAligner directly predicts the time in-  
532 dices at slots. During training, causal attention  
533 masking with non-shifted input and label sequences  
534 allows each slot to predict its own timestamp index  
535 based on itself and prior context, with loss com-  
536 puted only at slot positions. Dynamic slot insertion  
537 enables timestamp prediction at arbitrary positions.  
538 Non-autoregressive inference is supported, avoid-  
539 ing hallucinations and improving speed. Experi-  
540 ments show that LLM-ForcedAligner is an accu-  
541 rate, fast and customized FA method for multilin-  
542 gual, crosslingual, and long-form speech scenarios.

## 543 Limitations

544 Manually annotated millisecond-level timestamps  
545 also suffer from ambiguous temporal boundaries,  
546 making it difficult for any method to achieve highly  
547 precise alignment. Therefore, we train LLM-  
548 ForcedAligner on MFA-labeled pseudo-timestamp  
549 datasets and evaluate it on manually annotated  
550 timestamp datasets, thereby assessing whether the  
551 predicted timestamps exhibit perceptible devia-  
552 tions from human judgments. We generate pseudo-  
553 timestamp labels using the MFA method to con-  
554 struct most of the training and test datasets, which  
555 causes LLM-ForcedAligner to fit the relatively ac-  
556 curate MFA timestamp-prediction distribution. Al-  
557 though the manually annotated test dataset allows  
558 us to assess LLM-ForcedAligner’s performance  
559 in practical scenarios, it covers only Chinese and  
560 cannot verify the actual performance in other lan-  
561 guages. In addition, the training dataset exhibits an  
562 uneven language distribution, which may explain  
563 why the AAS for other languages is higher than that  
564 for Chinese and English. In future work, we plan to  
565 explore methods to improve LLM-ForcedAligner’s  
566 performance across other languages efficiently and  
567 to extend its application scenarios to more chal-  
568 lenging conditions, such as meetings, music, and  
569 film or television content.

## 570 References

571 Max Bain, Jaesung Huh, Tengda Han, and Andrew Zis-  
572 serman. 2023. WhisperX: Time-Accurate Speech  
573 Transcription of Long-Form Audio. In *Proc. Inter-  
574 speech*, pages 4489–4493.

575 Hui Bu, Jiayu Du, Xingyu Na, Bengu Wu, and Hao  
576 Zheng. 2017. AISHELL-1: An open-source Man-  
577 darin speech corpus and a speech recognition base-  
578 line. In *Proc. O-COCOSDA*, pages 1–5.

579 Guoguo Chen, Shuzhou Chai, Guan-Bo Wang, Jiayu  
580 Du, Wei-Qiang Zhang, Chao Weng, Dan Su, Daniel  
581 Povey, Jan Trmal, Junbo Zhang, Mingjie Jin, Sanjeev  
582 Khudanpur, Shinji Watanabe, Shuaijiang Zhao, Wei  
583 Zou, Xiangang Li, Xuchen Yao, Yongqing Wang,  
584 Zhao You, and Zhiyong Yan. 2021. GigaSpeech: An  
585 Evolving, Multi-Domain ASR Corpus with 10, 000  
586 Hours of Transcribed Audio. In *Proc. Interspeech*,  
587 pages 3670–3674.

588 Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei,  
589 Zhifang Guo, Yichong Leng, Yuanjun Lv, Jinzheng  
590 He, Junyang Lin, Chang Zhou, and Jingren Zhou.  
591 2024. Qwen2-Audio Technical Report. *CoRR*,  
592 abs/2407.10759.

Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shil-  
iang Zhang, Zhijie Yan, Chang Zhou, and Jingren  
Zhou. 2023. Qwen-Audio: Advancing Universal  
Audio Understanding via Unified Large-Scale Audio-  
Language Models. *CoRR*, abs/2311.07919. 593  
594  
595  
596  
597

Jiayu Du, Xingyu Na, Xuechen Liu, and Hui Bu. 2018.  
AISHELL-2: Transforming Mandarin ASR Research  
Into Industrial Scale. *CoRR*, abs/1808.10583. 598  
599  
600

Zhihao Du, Changfeng Gao, Yuxuan Wang, Fan Yu,  
Tianyu Zhao, Hao Wang, Xiang Lv, Hui Wang,  
Chongjia Ni, Xian Shi, Keyu An, Guanrou Yang,  
Yabin Li, Yanni Chen, Zhifu Gao, Qian Chen, Yue  
Gu, Mengzhe Chen, Yafeng Chen, and 3 others.  
2025. CosyVoice 3: Towards In-the-wild Speech  
Generation via Scaling-up and Post-training. *CoRR*,  
abs/2505.17589. 601  
602  
603  
604  
605  
606  
607  
608

Xuelong Geng, Qijie Shao, Hongfei Xue, Shuiyuan  
Wang, Hanke Xie, Zhao Guo, Yi Zhao, Guojian  
Li, Wenjie Tian, Chengyou Wang, Zhixian Zhao,  
Kangxiang Xia, Ziyu Zhang, Zhennan Lin, Tian-  
lun Zuo, Mingchen Shao, Yuang Cao, Guobin Ma,  
Longhao Li, and 4 others. 2025a. OSUM-EChat:  
Enhancing End-to-End Empathetic Spoken Chatbot  
via Understanding-Driven Spoken Dialogue. *CoRR*,  
abs/2508.09600. 609  
610  
611  
612  
613  
614  
615  
616  
617

Xuelong Geng, Kun Wei, Qijie Shao, Shuiyun Liu,  
Zhennan Lin, Zhixian Zhao, Guojian Li, Wenjie Tian,  
Peikun Chen, Yangze Li, Pengcheng Guo, Mingchen  
Shao, Shuiyuan Wang, Yuang Cao, Chengyou Wang,  
Tianyi Xu, Yuhang Dai, Xinfu Zhu, Yue Li, and 2  
others. 2025b. OSUM: Advancing Open Speech  
Understanding Models with Limited Resources in  
Academia. *CoRR*, abs/2501.13306. 618  
619  
620  
621  
622  
623  
624  
625

Xuelong Geng, Tianyi Xu, Kun Wei, Bingshen Mu,  
Hongfei Xue, He Wang, Yangze Li, Pengcheng Guo,  
Yuhang Dai, Longhao Li, Mingchen Shao, and Lei  
Xie. 2024. Unveiling the Potential of LLM-Based  
ASR on Chinese Open-Source Datasets. In *Proc.  
ISCSLP*, pages 26–30. 626  
627  
628  
629  
630  
631

Haorui He, Zengqiang Shang, Chaoren Wang, Xuyuan  
Li, Yicheng Gu, Hua Hua, Liwei Liu, Chen Yang,  
Jiaqi Li, Peiyang Shi, Yuancheng Wang, Kai Chen,  
Pengyuan Zhang, and Zhizheng Wu. 2024. Emilia:  
An Extensive, Multilingual, and Diverse Speech  
Dataset For Large-Scale Speech Generation. In *Proc.  
SLT*, pages 885–890. 632  
633  
634  
635  
636  
637  
638

Ludwig Kürzinger, Dominik Winkelbauer, Lujun Li,  
Tobias Watzel, and Gerhard Rigoll. 2020. CTC-  
Segmentation of Large Corpora for German End-  
to-End Speech Recognition. In *Proc. Speech and  
Computer*, pages 267–278. 639  
640  
641  
642  
643

Hexin Liu, Xiangyu Zhang, Haoyang Zhang,  
Leibny Paola Garcia-Perera, Andy W. H. Khong,  
Eng Siong Chng, and Shinji Watanabe. 2025. Align-  
ing Speech to Languages to Enhance Code-Switching  
Speech Recognition. *IEEE Transactions on Audio,  
Speech and Language Processing*, 33:4712–4725. 644  
645  
646  
647  
648  
649

|     |  |   |     |
|-----|--|---|-----|
| 650 | Michael McAuliffe, Michaela Socolof, Sarah Mihuc,          | Xinsheng Wang, Mingqi Jiang, Ziyang Ma, Ziyu Zhang,       | 706 |
| 651 | Michael Wagner, and Morgan Sonderegger. 2017.              | Songxiang Liu, Linqin Li, Zheng Liang, Qixi Zheng,        | 707 |
| 652 | Montreal Forced Aligner: Trainable Text-Speech             | Rui Wang, Xiaoqin Feng, Weizhen Bian, Zhen Ye,            | 708 |
| 653 | Alignment Using Kaldi. In <i>Proc. Interspeech</i> , pages | Sitong Cheng, Ruibin Yuan, Zhixian Zhao, Xinfu            | 709 |
| 654 | 498–502.   | Zhu, Jiahao Pan, Liumeng Xue, Pengcheng Zhu,              | 710 |
| 655 | Bingshen Mu, Hexin Liu, Hongfei Xue, Kun Wei, and          | and 6 others. 2025a. Spark-TTS: An Efficient LLM-         | 711 |
| 656 | Lei Xie. 2026. Hearing More with Less: Multi-              | Based Text-to-Speech Model with Single-Stream De-         | 712 |
| 657 | Modal Retrieval-and-Selection Augmented Conversa-          | coupled Speech Tokens. <i>CoRR</i> , abs/2503.01710.      | 713 |
| 658 | tional LLM-Based ASR. In <i>Proc. AAAI</i> .               | Xiong Wang, Yangze Li, Chaoyou Fu, Yike Zhang, Yun-       | 714 |
| 659 | Vassil Panayotov, Guoguo Chen, Daniel Povey, and           | hang Shen, Lei Xie, Ke Li, Xing Sun, and Long Ma.         | 715 |
| 660 | Sanjeev Khudanpur. 2015. Librispeech: An ASR               | 2025b. Freeze-Omni: A Smart and Low Latency               | 716 |
| 661 | corpus based on public domain audio books. In <i>Proc.</i> | Speech-to-speech Dialogue Model with Frozen LLM.          | 717 |
| 662 | <i>ICASSP</i> , pages 5206–5210.                           | In <i>Proc. ICML</i> .                                    | 718 |
| 663 | Daniel Povey, Arnab Ghoshal, Gilles Boulianne, Lukas       | Jin Xu, Zhifang Guo, Hangrui Hu, Yunfei Chu, Xiong        | 719 |
| 664 | Burget, Ondrej Glembek, Nagendra Goel, Mirko               | Wang, Jinzheng He, Yuxuan Wang, Xian Shi, Ting            | 720 |
| 665 | Hannemann, Petr Motlicek, Yanmin Qian, Petr                | He, Xinfu Zhu, Yuanjun Lv, Yongqi Wang, Dake              | 721 |
| 666 | Schwarz, and 1 others. 2011. The Kaldi speech recog-       | Guo, He Wang, Linhan Ma, Pei Zhang, Xinyu                 | 722 |
| 667 | nition toolkit. In <i>Proc. ASRU</i> , pages 1–5.          | Zhang, Hongkun Hao, Zishan Guo, and 19 oth-               | 723 |
| 668 | Vineel Pratap, Qiantong Xu, Anuroop Sriram, Gabriel        | ers. 2025. Qwen3-Omni Technical Report. <i>CoRR</i> ,     | 724 |
| 669 | Synnaeve, and Ronan Collobert. 2020. MLS: A                | abs/2509.17765.   | 725 |
| 670 | Large-Scale Multilingual Dataset for Speech Re-            | An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,          | 726 |
| 671 | search. In <i>Proc. Interspeech</i> , pages 2757–2761.     | Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao,               | 727 |
| 672 | Elena Rastorgueva, Vitaly Lavrukhin, and Boris Gins-       | Chengen Huang, Chenxu Lv, Chujie Zheng, Day-              | 728 |
| 673 | burg. 2023. NeMo Forced Aligner and its application        | iheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao              | 729 |
| 674 | to word alignment for subtitle generation. In <i>Proc.</i> | Ge, Haoran Wei, Huan Lin, Jialong Tang, and 40            | 730 |
| 675 | <i>Interspeech</i> , pages 5257–5258.                      | others. 2025. Qwen3 technical report. <i>CoRR</i> ,       | 731 |
| 676 | Xian Shi, Yanni Chen, Shiliang Zhang, and Zhijie Yan.      | abs/2505.09388.   | 732 |
| 677 | 2023. Achieving Timestamp Prediction While Recog-          | An Yang, Baosong Yang, Beichen Zhang, Binyuan             | 733 |
| 678 | nizing with Non-Autoregressive End-to-End ASR              | Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayi-              | 734 |
| 679 | Model. <i>CoRR</i> , arXiv:2301.12343.                     | heng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian           | 735 |
| 680 | Zhiyuan Tang, Dong Wang, Yanguang Xu, Jianwei Sun,         | Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Ji-       | 736 |
| 681 | Xiaoning Lei, Shuaijiang Zhao, Cheng Wen, Xingjun          | axi Yang, Jingren Zhou, Junyang Lin, Kai Dang, and        | 737 |
| 682 | Tan, Chuandong Xie, Shuran Zhou, Rui Yan, Chen-            | 22 others. 2024. Qwen2.5 technical report. <i>CoRR</i> ,  | 738 |
| 683 | jia Lv, Yang Han, Wei Zou, and Xiangang Li. 2021.          | abs/2412.15115.   | 739 |
| 684 | KeSpeech: An Open Source Speech Dataset of Man-            | Heiga Zen, Viet Dang, Rob Clark, Yu Zhang, Ron J.         | 740 |
| 685 | darin and Its Eight Subdialects. In <i>Proc. NeurIPS</i> . | Weiss, Ye Jia, Zhifeng Chen, and Yonghui Wu. 2019.        | 741 |
| 686 | Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier      | LibriTTS: A Corpus Derived from LibriSpeech for           | 742 |
| 687 | Martinet, Marie-Anne Lachaux, Timothée Lacroix,            | Text-to-Speech. In <i>Proc. Interspeech</i> , pages 1526– | 743 |
| 688 | Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal         | 1530.   | 744 |
| 689 | Azhar, Aurélien Rodriguez, Armand Joulin, Edouard          | Binbin Zhang, Hang Lv, Pengcheng Guo, Qijie Shao,         | 745 |
| 690 | Grave, and Guillaume Lample. 2023a. LLaMA:                 | Chao Yang, Lei Xie, Xin Xu, Hui Bu, Xiaoyu Chen,          | 746 |
| 691 | Open and Efficient Foundation Language Models.             | Chenchen Zeng, Di Wu, and Zhendong Peng. 2022.            | 747 |
| 692 | <i>CoRR</i> , abs/2302.13971.                              | WENETSPEECH: A 10000+ Hours Multi-Domain                  | 748 |
| 693 | Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-         | Mandarin Corpus for Speech Recognition. In <i>Proc.</i>   | 749 |
| 694 | bert, Amjad Almahairi, Yasmine Babaei, Nikolay             | <i>ICASSP</i> , pages 6182–6186.                          | 750 |
| 695 | Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti          | <b>A Appendix</b>   | 751 |
| 696 | Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-        | <b>A.1 Data Statistics</b>                                | 752 |
| 697 | Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,         | Table 7 summarizes the sources and durations of           | 753 |
| 698 | Jude Fernandes, Jeremy Fu, Wenyin Fu, and 49 oth-          | the overall datasets used for LLM-ForcedAligner.          | 754 |
| 699 | ers. 2023b. Llama 2: Open Foundation and Fine-             | Consequently, we provide a brief introduction to          | 755 |
| 700 | Tuned Chat Models. <i>CoRR</i> , abs/2307.09288.           | the open-source datasets included.                        | 756 |
| 701 | Liang-Hsuan Tseng, Yu-Kuan Fu, Heng-Jui Chang,             | <b>AISHELL-1</b> (Bu et al., 2017): The dataset con-      | 757 |
| 702 | and Hung-yi Lee. 2021. Mandarin-English                    | tains 400 speakers and over 170 hours of Mandarin         | 758 |
| 703 | code-switching speech recognition with self-               | speech data, covering 5 topics: “Finance”, “Sci-          | 759 |
| 704 | supervised speech representation models. <i>CoRR</i> ,     | ence and Technology”, “Sports”, “Entertainments”,         | 760 |
| 705 | abs/2110.03504.  |   |     |

Table 7: Overall dataset statistics.

| Language   | Sources  | Hours |
|------------|--|-------|
| Chinese    | AISHELL-1, AISHELL-2,  | 31536 |
|            | WenetSpeech, Aidatatang_200zh, Magicdata, KeSpeech, Private data |       |
| English    | LibriSpeech, GigaSpeech, VCTK, LibriTTS, Private data            | 20322 |
| French     | Emilia, MLS  | 611   |
| German     | Emilia, MLS  | 517   |
| Italian    | MLS, Private data  | 602   |
| Japanese   | Emilia   | 731   |
| Korean     | Emilia, Private data   | 662   |
| Portuguese | MLS, Private data  | 316   |
| Russian    | Private data   | 268   |
| Spanish    | MLS, Private data  | 579   |

and “News”. Transcripts are manually filtered to eliminate improper contents.

**AISHELL-2** (Du et al., 2018): The dataset contains 1991 speakers and over 1000 hours of clean reading speech data. The content of the recording covers 8 major topics: voice commands such as IoT device control and digital sequential input, places of interest, entertainment, finance, technology, sports, English spellings and free speaking without specific topic.

**WenetSpeech** (Zhang et al., 2022): The dataset is a multi-domain Mandarin corpus consisting of 10000+ hours high-quality labeled speech, 2400+ hours weakly labeled speech, and about 10000 hours unlabeled speech, with 22400+ hours in total. The data from YouTube and Podcast, which covers a variety of speaking styles, scenarios, domains, topics and noisy conditions.

**Aidatatang\_200zh**<sup>8</sup>: The dataset is a Chinese Mandarin speech corpus, containing 200 hours of speech data from 600 speakers. The transcription accuracy for each sentence is larger than 98%.

**Magicdata**<sup>9</sup>: The dataset contains 755 hours of scripted read speech data from 1080 native speakers of the Mandarin Chinese spoken in mainland China. The sentence transcription accuracy is higher than 98%. The domain of recording texts is diversified, including interactive Q&A, music search, SNS messages, home command and control, etc.

**KeSpeech** (Tang et al., 2021): The dataset comprises 1,542 hours of speech signals recorded by 27,237 speakers across 34 cities in China, with pronunciation in standard Mandarin and 8 subdialects. Two professional data companies manually label

<sup>8</sup><http://openslr.magicdatatech.com/62/>

<sup>9</sup><https://www.openslr.org/68/>

the dataset with three steps.

**LibriSpeech** (Panayotov et al., 2015): The dataset is derived from audiobooks that are part of the LibriVox project, contains 1000 hours of speech sampled at 16kHz, and has been carefully segmented and aligned.

**GigaSpeech** (Chen et al., 2021): The dataset is a multi-domain English speech recognition corpus with 10,000 hours of high quality labeled audio suitable for supervised training, and 40,000 hours of total audio suitable for semi-supervised and unsupervised training.

**LibriTTS** (Zen et al., 2019): The dataset is a multi-speaker English corpus consisting of approximately 585 hours of read English speech at a 24kHz sampling rate from 2,456 speakers, along with the corresponding texts.

**MLS** (Pratap et al., 2020): The dataset is derived from read audiobooks from LibriVox and consists of 8 languages, including about 44.5K hours of English and about 6K hours of other languages.

**Emilia** (He et al., 2024): The dataset contains over 101k hours of speech data at 24 kHz and covers six languages. It comprises mostly spontaneous speech, covering a wide range of speaking styles.

**Private Data**: The datasets are primarily sourced from professional data providers and encompass multilingual read and conversational speech, and all annotations are manually reviewed.

**Construction of Training Dataset**: All data in Table 7 are first annotated with word-level or character-level start and end timestamp pseudo-labels obtained using MFA. For each language and each data source, a small portion is reserved as the test dataset. From the remaining data, 30% is used as raw speech for training, while the remaining 70% is used to construct long-form speech through mixing. Specifically, a target duration is randomly sampled from a uniform distribution between 30s and 500s, and raw speech from different languages is randomly concatenated until the target duration is reached. To further enhance data diversity, various types of noise are randomly mixed into the mixed speech.

**Construction of Test Datasets**: For each language and data source, 30% of the reserved raw speech is used as the MFA-labeled Raw test datasets, while the remaining 70% is used to create long-form speech through mixing. The mixing procedure is the same as for the training dataset, but divided into monolingual and crosslingual mixing. The resulting mixed test datasets serve as the MFA-labeled

Table 8: Configuration of Compared FA methods. “-” indicates that no model is available for this language, all model names are identical to the corresponding official open-source model names.

| Language   | Monotonic-Aligner | NFA                                  | WhisperX                  |
|------------|-------------------|--------------------------------------|---------------------------|
| Chinese    | Paraformer        | stt_zh_citrinet_1024_gamma_0_25      | -                         |
| English    | -                 | stt_en_fastconformer_hybrid_large_pc | WAV2VEC2_ASR_BASE_960H    |
| French     | -                 | stt_fr_conformer_ctc_large           | VOXPOPULI_ASR_BASE_10K_FR |
| German     | -                 | stt_de_fastconformer_hybrid_large_pc | VOXPOPULI_ASR_BASE_10K_DE |
| Italian    | -                 | stt_it_fastconformer_hybrid_large_pc | VOXPOPULI_ASR_BASE_10K_IT |
| Japanese   | -                 | -                                    | -                         |
| Korean     | -                 | -                                    | -                         |
| Portuguese | -                 | -                                    | -                         |
| Russian    | -                 | stt_ru_quartznet15x5                 | -                         |
| Spanish    | -                 | stt_es_fastconformer_hybrid_large_pc | VOXPOPULI_ASR_BASE_10K_ES |

Mixed test datasets. In addition, we evaluate the performance of LLM-ForcedAligner on an internal, manually annotated Chinese data source. Specifically, this source is first used as the human-labeled Raw test dataset, and background noise is then added to create the human-labeled Raw-Noisy test dataset. It is also concatenated to form the human-labeled Mixed-60s and human-labeled Mixed-300s test datasets, each with a maximum duration of 60s and 300s, respectively. Finally, it is combined with the MFA-labeled Raw test dataset to create the human-labeled Mixed-Crosslingual test dataset.

phoneme recognition models specific to different languages. Its standard evaluation procedures are on its homepage<sup>13</sup>, and the available checkpoint list is on the inference script file<sup>14</sup>.

## A.2 Details of Compared FA Methods

Since the compared FA methods require switching backbone models across languages, Table 8 lists the backbone model used by each method for each language in the comparison experiments. We strictly reproduce the results on our test datasets following the standard evaluation procedures of these FA methods.

**Monotonic-Aligner** (Shi et al., 2023): It supports only Chinese and uses Paraformer as its backbone model. Its standard evaluation procedures are on its homepage<sup>10</sup>.

**NFA** (Rastorgueva et al., 2023): It is a tool for generating token-, word-, and segment-level timestamps of speech in audio using NeMo’s CTC-based ASR models. Its standard evaluation procedures are on its homepage<sup>11</sup>, and the available checkpoint list is on the NeMo ASR collection page<sup>12</sup>.

**WhisperX** (Bain et al., 2023): It performs FA using

<sup>10</sup>[https://modelscope.cn/models/iic/speech\\_timestamp\\_prediction-v1-16k-offline](https://modelscope.cn/models/iic/speech_timestamp_prediction-v1-16k-offline)

<sup>11</sup>[https://github.com/NVIDIA-NeMo/NeMo/tree/main/tools/nemo\\_forced\\_aligner](https://github.com/NVIDIA-NeMo/NeMo/tree/main/tools/nemo_forced_aligner)

<sup>12</sup>[https://catalog.ngc.nvidia.com/orgs/nvidia/collections/nemo\\_asr](https://catalog.ngc.nvidia.com/orgs/nvidia/collections/nemo_asr)

<sup>13</sup><https://github.com/m-bain/whisperX>

<sup>14</sup><https://github.com/m-bain/whisperX/blob/main/whisperx/alignment.py>