# LAE-ST-MOE: BOOSTED LANGUAGE-AWARE ENCODER USING SPEECH TRANSLATION AUXILIARY TASK FOR E2E CODE-SWITCHING ASR

*Guodong Ma[†], Wenxuan Wang[†], Yuke Li[†*], Yuting Yang[†], Binbin Du[†], Haoran Fu[‡]*

[†]NetEase Yidun AI Lab, Hangzhou, China
[‡]Department of Civil Engineering, Zhejiang University

## ABSTRACT

Recently, to mitigate the confusion between different languages in code-switching (CS) automatic speech recognition (ASR), the conditionally factorized models, such as the language-aware encoder (LAE), explicitly disregard the contextual information between different languages. However, this information may be helpful for ASR modeling. To alleviate this issue, we propose the LAE-ST-MoE framework. It incorporates speech translation (ST) tasks into LAE and utilizes ST to learn the contextual information between different languages. It introduces a task-based mixture of expert modules, employing separate feed-forward networks for the ASR and ST tasks. Experimental results on the ASRU 2019 Mandarin-English CS challenge dataset demonstrate that, compared to the LAE-based CTC, the LAE-ST-MoE model achieves a 9.26% mix error reduction on the CS test with the same decoding parameter. Moreover, the well-trained LAE-ST-MoE model can perform ST tasks from CS speech to Mandarin or English text.

***Index Terms***— Automatic speech recognition, Mandarin-English code-switching, speech translation, mixture of expert

## 1. INTRODUCTION

With the rise of end-to-end (E2E) automatic speech recognition (ASR), researchers [1–25] explore different E2E ASR scenarios. An utterance that includes two or more languages is known as a code-switching (CS) scenario, which is generally divided into occurring at an utterance level (extra-sentential CS) or within an utterance (intra-sentential CS). It is still a challenging ASR scenario.

Several challenges are conventionally encountered in modeling CS speech: firstly, the real paired CS audio is data-scarce, and secondly, the conventional models are not good at modeling CS speech due to the confusion between different languages. To alleviate the first issue, researchers propose technical methods to study the rules of CS occurrence and synthesize CS paired data [11–15] or explore the affection of monolingual data [26–28]. As for the second issue, the

structures like Connectionist Temporal Classification (CTC)-, attention-, and transducer-based E2E models have been investigated for CS ASR [13–22]. Recently, to mitigate the second issue, the conditionally factorized frameworks [29–32] are proposed to decompose the CS task (e.g., Mandarin-English CS) into two modeling steps: 1) recognizing Mandarin and English part, respectively, and 2) composing processed monolingual segments into a CS sequence. However, in modeling step 1) for these methods, the model only utilizes the information of the monolingual part. We know that, when modeling the non-streaming E2E ASR task, the prediction of each unit generally relies on overall audio contextual information.

To solve the issues of the conditionally factorized models [29–32] (e.g., LAE [32]), we propose the LAE-ST-MoE framework. It incorporates speech translation (ST) tasks into LAE [32] and utilizes ST to facilitate the learning of contextual information between Mandarin and English, thereby impacting the model's encoder through joint learning. In addition, inspired by [22–24], the LAE-ST-MoE introduces a task-based mixture of expert (MoE) approach, employing separate feed-forward networks (FFNs) for the ASR and ST tasks.

Our experiment is conducted on the classic CS benchmark, i.e., ASRU 2019 Mandarin-English CS challenge dataset [33]. Since the data does not have ST labels, we use the large machine translation (MT) model from ModeScope to label the data, which is based on the CSANMT algorithm [34]. In the experiments, compared to the LAE-based system, the LAE-ST-MoE model achieves a relative performance improvement of about 6%-9% in ASR tasks on all test sets. Moreover, our model does not introduce extra decoding computational complexity. In addition, the trained LAE-ST-MoE model can perform ST tasks from CS speech to Mandarin or English text and has achieved good BLEU. Then, it is easy to extend our model to one-to-many ST tasks.

Our main contributions are as follows: (1) To our best knowledge, we are the first to propose using the ST task to introduce richer cross-lingual contextual information to boost the monolingual modeling stage of LAE; (2) We introduce an MoE between ASR and ST tasks to make each task more focused, thereby improving the overall recognition performance of the model without extra decoding computational complexity; (3) The well-trained LAE-ST-MoE model can perform

---

* Corresponding author

ST tasks from CS speech to Mandarin or English text, and the structure is easy to extend to one-to-many ST tasks.

## 2. PROBLEM FORMULATIONS AND MOTIVATION

In the Mandarin-English CS ASR system [29–32], we know that the basis is to model the label-to-frame alignments. For each T-length speech feature sequence $X = \{x_t | t = 1, ..., T\}$ and L-length CS label sequence $Y = \{y_\ell \in (\mathbb{V}^{\mathrm{Man}} \bigcup \mathbb{V}^{\mathrm{En}} | \ell = 1, ..., L)\}$, there are several possible T-length label-to-frame sequences $Z = \{z_t \in (\mathbb{V}^{\mathrm{Man}} \bigcup \mathbb{V}^{\mathrm{En}} \bigcup \{\emptyset\}) | t = 1, ..., T\}$, where $\emptyset$ denotes a blank symbol in CTC [1] based system, $\mathbb{V}^{\mathrm{Man}}$, and $\mathbb{V}^{\mathrm{En}}$ respectively represents to the Mandarin and English part in CS. However, for each CS Z, there always are two corresponding monolingual label-to-frame sequences $Z^{\mathrm{Man}} = \{z_t^{\mathrm{Man}} \in \{\mathbb{Z}^{\mathrm{Man}} \bigcup \{\emptyset\}\} | t = 1, ..., T\}$ and $Z^{\mathrm{En}} = \{z_t^{\mathrm{En}} \in \mathbb{Z}^{\mathrm{En}} \bigcup \{\emptyset\}\} | t = 1, ..., T\}$. Therefore, the label-to-frame posterior $P(Y|X)$ can thus be represented in terms of CS, $P(Z|X)$, and monolingual, $P(Z^{\mathrm{Man}}|X)$ and $P(Z^{\mathrm{En}}|X)$, label-to-frame posteriors:

$$P(Y|X) = \sum_{Z \in \mathbb{Z}} \sum_{Z^{\mathrm{Man}} \in \mathbb{Z}^{\mathrm{Man}}} \sum_{Z^{\mathrm{En}} \in \mathbb{Z}^{\mathrm{En}}} P(Z, Z^{\mathrm{Man}}, Z^{\mathrm{En}}|X) \quad (1)$$

where $\mathbb{Z}$ and $\mathbb{Z}^{\mathrm{Man/En}}$ denote sets of all possible CS and monolingual label-to-frame alignments for a given Y. By applying Bayes' formula, the $P(Z, Z^{\mathrm{Man}}, Z^{\mathrm{En}}|X)$ in Eq.(1) can be transformed into the following expression:

$$P(Z, Z^{\mathrm{Man}}, Z^{\mathrm{En}}|X) = P(Z|Z^{\mathrm{Man}}, Z^{\mathrm{En}}, X) \\ \times P(Z^{\mathrm{Man}}, Z^{\mathrm{En}}|X) \quad (2)$$

and

$$P(Z^{\mathrm{Man}}, Z^{\mathrm{En}}|X) = P(Z^{\mathrm{Man}}|Z^{\mathrm{En}}, X) \times P(Z^{\mathrm{En}}|X). \quad (3)$$

Two assumptions are made. The first assumption is that once $Z^{\mathrm{Man}}$ and $Z^{\mathrm{En}}$ are given, no additional information from observation X is needed to determine Z. The second assumption is that $Z^{\mathrm{Man}}$ and $Z^{\mathrm{En}}$ are independent, given X. Therefore, combined with Eq. (1-3), the eq. (1) can be shown:

$$P(Y|X) \approx \sum_{Z \in \mathbb{Z}} P(Z|Z^{\mathrm{Man}}, Z^{\mathrm{En}}) \times \sum_{Z^{\mathrm{Man}} \in \mathbb{Z}^{\mathrm{Man}}} P(Z^{\mathrm{Man}}|X) \\ \times \sum_{Z^{\mathrm{En}} \in \mathbb{Z}^{\mathrm{En}}} P(Z^{\mathrm{En}}|X). \quad (4)$$

To achieve the transformation from Eq. (1) to Eq. (4), the monolingual-specific encoder is introduced by the conditionally factorized structures [29–32] to optimize the representation of each language separately. For example, the token sequence of the CS audio is like " 真 正 做 到 _happy _every day". When forwarding Mandarin-specific encoder, the reference text will be replaced with " 真 正 做 到 <En_tok> <En_tok> <En_tok>" and ignore the English part, where <En_tok> can refer to <UNK> [31] or <Eng> [32]. As for the English-specific encoder, as shown in Figure 1, it is the same as the Mandarin-specific encoder. Further consider the

modeling process, e.g., Mandarin-specific encoder, the model will not learn English contextual information in the CS audio, which could potentially improve its performance on the Mandarin part. However, the ST model is capable of converting contextual information from various languages into one language. Therefore, applying ST tasks to enrich the contextual information between the two languages in CS ASR can be reasonable and feasible. Based on the LAE architecture [32] and joint learning mechanism, we propose LAE-ST-MoE architecture, which uses ST as an auxiliary task to bring more contextual information for ASR. The details of our proposed LAE-ST-MoE will be presented in the next section.

## 3. PROPOSED FRAMEWORKS

### 3.1. LAE-ST-MoE architecture

The LAE structure [32] has a shared encoder module, two language-specific encoders for Mandarin and English, and a global ASR decoder. The monolingual-specific encoder is imposed by a corresponding monolingual-specific CTC loss. To alleviate the issues of LAE discussed in section 2, we propose the LAE-ST-MoE model architecture, as shown in Figure 1, which introduces two LAE-ST-MoE encoders and two ST decoders based on LAE. If $N_{\mathbf{Share}}$ represents the number of layers in the shared encoder, and $N_{\mathbf{Mono}}$ represents the number of layers in the monolingual-specific encoder. Then, the LAE-ST-MoE encoder has $N$ layers, where $N$ is equal to ($N_{\mathbf{Encoder}} - N_{\mathbf{Share}} - N_{\mathbf{Mono}}$) and $N_{\mathbf{Encoder}}$ refers to the overall encoder layers. A common ST cross-entropy loss imposes the ST decoder, which consists of 6 Transformer-based blocks. In addition, the ST and ASR tasks are jointly trained using FFN-based MoE. A detailed explanation of the proposed LAE-ST-MoE model is presented as follows.

If given the input feature sequence X, the shared Transformer encoder will transform it to representation $\mathbf{H_{share}}$:

$$\mathbf{H_{share}} = \mathrm{Encoder}_{\mathrm{share}}(X). \quad (5)$$

Furthermore, the $\mathbf{H_{share}}$ will be forwarded to the LAE-ST-MoE encoder, which replaces the FFN of the Transformer encoder with the FFN-MoE module. It produces the hybrid ASR-ST representation $\mathbf{H_{Man\_ASR\_ST}}$ and $\mathbf{H_{En\_ASR\_ST}}$ using multi-head self-attention (MHSA):

$$\mathbf{H_{Man\_ASR\_ST}} = \mathrm{MHSA}(\mathrm{LNorm}(\mathbf{H_{share}})) \quad (6)$$

$$\mathbf{H_{En\_ASR\_ST}} = \mathrm{MHSA}(\mathrm{LNorm}(\mathbf{H_{share}})) \quad (7)$$

where LNorm denotes LayerNorm [35]. Based on $\mathbf{H_{En\_ASR\_ST}}$ and $\mathbf{H_{Man\_ASR\_ST}}$, the FFN-MoE is forward to get ASR representation $\mathbf{H^0_{En\_ASR}}$, $\mathbf{H^0_{Man\_ASR}}$, and ST representation $\mathbf{H_{Man2En\_ST}}$, $\mathbf{H_{En2Man\_ST}}$, respectively:

$$\mathbf{H^0_{Man\_ASR}} = \mathrm{LNorm}(\mathrm{FFN\_MoE}(\mathbf{H_{Man\_ASR\_ST}})) \quad (8)$$

$$\mathbf{H_{En2Man\_ST}} = \mathrm{LNorm}(\mathrm{FFN\_MoE}(\mathbf{H_{Man\_ASR\_ST}})) \quad (9)$$

$$\mathbf{H^0_{En\_ASR}} = \mathrm{LNorm}(\mathrm{FFN\_MoE}(\mathbf{H_{En\_ASR\_ST}})) \quad (10)$$

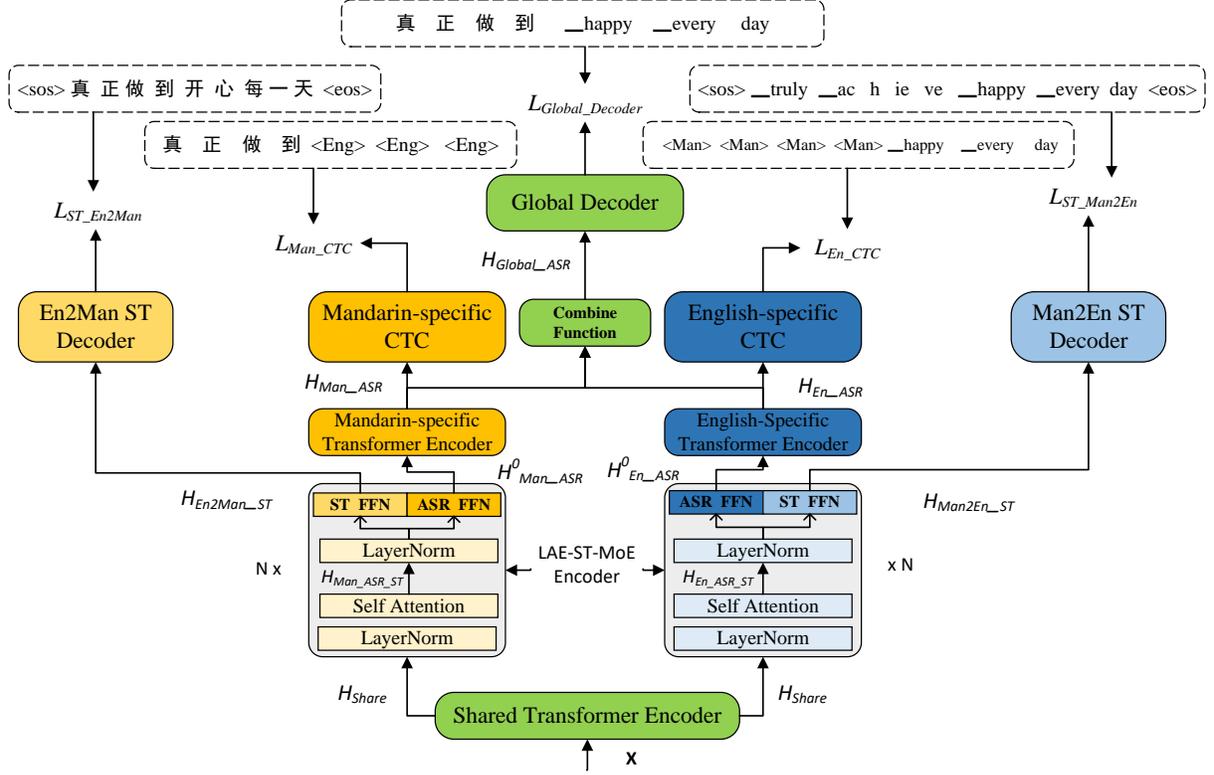$$\mathbf{H_{Man2En\_ST}} = \mathrm{LNorm}(\mathrm{FFN\_MoE}(\mathbf{H_{En\_ASR\_ST}})). \quad (11)$$

**Fig. 1**. The framework of the proposed LAE-ST-MoE.

On the ST task side, $\mathbf{H_{En2Man\_ST}}$ and $\mathbf{H_{Man2En\_ST}}$ will forward to the En2Man and Man2En ST decoder, respectively. In ASR, it is the same as LAE, based on $\mathbf{H^0_{Man\_ASR}}$ and $\mathbf{H^0_{En\_ASR}}$, the Monolingual-specific encoder will produce the monolingual-specific representation $\mathbf{H_{Man\_ASR}}$, $\mathbf{H_{En\_ASR}}$ and combine these to get the global ASR representation $\mathbf{H_{Global\_ASR}}$:

$$\mathbf{H_{Man\_ASR}} = \mathrm{Encoder_{Man\_Spec}}(\mathbf{H^0_{Man\_ASR}}) \quad (12)$$

$$\mathbf{H_{En\_ASR}} = \mathrm{Encoder_{En\_Spec}}(\mathbf{H^0_{En\_ASR}}) \quad (13)$$

$$\mathbf{H_{Global\_ASR}} = \mathbf{H_{Man\_ASR}} + \mathbf{H_{En\_ASR}}. \quad (14)$$

### 3.2. Training and Decoding

In the LAE-ST-MoE model training stage, if the label text sequence for speech feature X is Y, we will apply the $\mathrm{Model_{MT}}$ from ModelScope to translate Y into Mandarin $\mathrm{Y^{Man}}$ and English $\mathrm{Y^{En}}$ text:

$$\mathrm{Y^{Man}} = \mathrm{Model_{MT\_En2Man}}(Y) \quad (15)$$

$$\mathrm{Y^{En}} = \mathrm{Model_{MT\_Man2En}}(Y). \quad (16)$$

Like [31, 32], we replace Y with monolingual-specific label $\mathrm{Y^{Man\_Spec}}$ and $\mathrm{Y^{En\_Spec}}$ using <Eng> and <Man>, respectively. Based on monolingual-specific ASR representation $\mathbf{H_{Man\_ASR}}$ and $\mathbf{H_{En\_ASR}}$, the monolingual-specific

ASR object $\mathcal{L}_{Spec}$ will be shown as follow:

$$\mathcal{L}_{Man\_CTC} = \mathrm{CTC_{Man\_Spec}}(\mathrm{Y^{Man\_Spec}}|\mathbf{H_{Man\_ASR}}) \quad (17)$$

$$\mathcal{L}_{En\_CTC} = \mathrm{CTC_{En\_Spec}}(\mathrm{Y^{En\_Spec}}|\mathbf{H_{En\_ASR}}) \quad (18)$$

$$\mathcal{L}_{Spec} = \frac{(\mathcal{L}_{Man\_CTC} + \mathcal{L}_{En\_CTC})}{2}. \quad (19)$$

Moreover, given the global ASR representation $\mathbf{H_{Global\_ASR}}$, the global ASR decoder object $\mathcal{L}_{Global\_Decoder}$ is:

$$\mathcal{L}_{Global\_Decoder} = \mathrm{Decoder_{Global}}(Y|\mathbf{H_{Global\_ASR}}). \quad (20)$$

Following [31], we also use $\lambda_{Spec}$ (we set it to 0.3 in the experiments) to combine $\mathcal{L}_{Spec}$ and $\mathcal{L}_{Global\_Decoder}$ to produce the overall ASR loss $\mathcal{L}_{ASR}$:

$$\mathcal{L}_{ASR}=\lambda_{Spec}\times \mathcal{L}_{Spec} + (1 - \lambda_{Spec}) \times \mathcal{L}_{Global\_Decoder}. \quad (21)$$

In the CTC-based ASR system, $\mathcal{L}_{Global\_Decoder}$ only represents the CTC loss. Otherwise, in hybrid CTC attention-based ASR [4], $\mathcal{L}_{Global\_Decoder}$ is the combination between CTC $\mathcal{L}_{Global\_CTC}$ and attention $\mathcal{L}_{Global\_Att}$ loss using $\lambda_{CTC}$:

$$\mathcal{L}_{Global\_Decoder} = \begin{aligned} & \mathcal{L}_{Global\_CTC} \times \lambda_{CTC} \\ & +(1 - \lambda_{CTC}) \times \mathcal{L}_{Global\_Att}. \end{aligned} \quad (22)$$

On the ST task, given ST representation ($\mathbf{H_{En2Man\_ST}}$ and $\mathbf{H_{Man2En\_ST}}$) and ST label sequence ($\mathrm{Y^{Man}}$ and $\mathrm{Y^{En}}$),

the overall ST loss $\mathcal{L}_{\text{ST}}$ is shown as follows:

$$\mathcal{L}_{\text{ST\_Man2En}} = \text{Decoder}_{\text{Man2En}}(Y^{\text{En}}|\mathbf{H_{Man2En\_ST}}) \quad (23)$$

$$\mathcal{L}_{\text{ST\_En2Man}} = \text{Decoder}_{\text{En2Man}}(Y^{\text{Man}}|\mathbf{H_{En2Man\_ST}}) \quad (24)$$

$$\mathcal{L}_{\text{ST}} = \frac{(\mathcal{L}_{\text{ST\_Man2En}} + \mathcal{L}_{\text{ST\_En2Man}})}{2} \quad (25)$$

where we use the cross-entropy loss for the ST tasks.

Based on the overall ASR loss $\mathcal{L}_{\text{ASR}}$ and ST loss $\mathcal{L}_{\text{ST}}$, the final training object $\mathcal{L}_{\text{Final}}$ is:

$$\mathcal{L}_{\text{Final}} = \mathcal{L}_{\text{ASR}} + \beta \times \mathcal{L}_{\text{ST}} \quad (26)$$

where $\beta$ is used to balance and regulate the ST effect.

In the ASR decoding stage, like the LAE structure [32], our model only gets the probabilities from the global ASR decoder. Therefore, compared with [32], our LAE-ST-MoE model has the same decoding computational complexity. In the ST decoding, our model uses the custom auto-regressive manner to forward the corresponding ST branch and get the final ST results. In addition, for monolingual Mandarin input, the En2Man ST decoder is comparable to the Mandarin ASR decoder. Therefore, we can easily fuse it into monolingual Mandarin decoding through rescoring. The same applies to monolingual English decoding.

**Table 1**. The details of the used Datasets

| Lang | Corpora | Dur. (Hrs) | | Utterance(k) | |
|---|---|---|---|---|---|
| | | Train | Eval | Train | Eval |
| CN | ASRU-Man [33] | 482.6 | 14.3 | 545.2 | 16.6 |
| EN | Librispeech [36] | 464.2 | 10.5 | 132.5 | 5.6 |
| CN-EN | ASRU-CS [33] | 199.0 | 20.3 | 186.4 | 16.2 |

## 4. EXPERIMENTS AND RESULTS

### 4.1. Datasets

We experiment on ASRU 2019 Mandarin-English code-switching challenge dataset [33]. Like [22], we split the same Mandarin monolingual subset of the ASRU 2019 dataset as our CN test. Moreover, we use the test-clean and test-other datasets from Librispeech [36] to create our monolingual English test EN. Then, the CS test CN-EN is from the official challenge test set. The details are presented in Table 1.

The 80-dimensional log filter-bank energy is our input acoustic features, which are extracted with a stride size 10ms and a window size 25ms. The cepstral mean and variance normalization (CMVN), and SpecAugment [37] is applied. The vocabulary consists of 7075 unique characters and 4989 BPE [38] tokens. In addition, as for the training and testing ST label, the EN2CN[1] and CN2EN[2] translation model, which is based on the CSANMT algorithm [34], both from ModelScope[3], is used to get the pseudo labels. Then, we

---

[1]https://www.modelscope.cn/models/damo/nlp_csanmt_translation_en2 zh/summary

[2]https://www.modelscope.cn/models/damo/nlp_csanmt_translation_zh2 en/summary

[3]https://github.com/modelscope/modelscope

use WeNet's [39] metrics calculation script[4] for ASR scoring, which includes word (WER), character (CER), mix (MER) error rate, and the sacrebleu [40] tool for ST scoring, which includes BLEU and translation error rate (TER).

For simpler expression, in Table 2, Table 3, Table 4, Table 5, and Table 6, we will use CN, EN, and ALL to represent the CER of monolingual Mandarin, the WER of monolingual English, and the total MER of the CS test set respectively.

### 4.2. Experimental setup

The experiments are both conducted on the ESPnet toolkit [41]. We use the hybrid CTC/Attention [4] model with a $\mathbf{N_{Encoder}}$=12 encoder, $\mathbf{N_{Decoder}}$=6 decoder, and the CTC-only model with a $\mathbf{N_{Encoder}}$=12, called the Vallina model. In the hybrid CTC/Attention model, $\lambda_{\text{CTC}}$ set to 0.3. In our implementation, following [32], the LAE-based baseline model contains a shared encoder block $\mathbf{N_{Share}}$=9 and a language-specific encoder block $\mathbf{N_{Mono}}$=3 for each language. As mentioned in section 3.1, the layers of the LAE-ST-MoE encoder $\mathbf{N}$ are equal to ( $\mathbf{N_{Encoder}}$ - $\mathbf{N_{Share}}$ - $\mathbf{N_{Mono}}$ ), and the number of layers will be given in the result section. In our models, all encoders and decoders are stacked Transformer-based blocks [5, 42] with an attention dimension of 256, 4 attention heads, and a feed-forward dimension of 2048.

We use the Adam optimizer with a Transformer-lr scale of 1 and warmup steps of 25k to train 100 epochs on 8 Tesla V100 GPUs. The dropout rate is 0.1 to prevent the model from over-fitting. In the training stage, we adopt a dynamic batch size strategy with a maximum batch size of 128. Moreover, we use Kenlm [43] to train a 4-gram language model with all training transcriptions and adopt the CTC prefix beam search for ST decoder rescore with a fixed beam size 10.

### 4.3. Experimental Results

#### 4.3.1. Main results

To show the effectiveness of our proposed LAE-ST-MoE framework, we compare it with LAE-based CTC and attention-based (AED) ASR models. We set the $\mathbf{N_{Mono}}$ to 1 and $\beta$ to 0.6 in these experiments. The ablation on $\beta$ and $\mathbf{N_{Mono}}$ will be shown in section 4.3.4. The results are shown in Table 2.
**CTC System:** Compared with the LAE-CTC ASR system (S2), our proposed LAE-ST-MoE CTC model (S3) achieve 9.26%, 8.57%, and 7.55% relative performance gain over the CS, mono EN, and CN tests, respectively, with the same decoding parameter. Especially in the English part of the CS test, our LAE-ST-MoE CTC (S3) shows a 10.09% WER reduction over the LAE CTC (S2) system. Moreover, it demonstrates a superior performance gain compared to Vanilla CTC (S1), which shows an about 20% error rate reduction in the CS test. Furthermore, the proposed LAE-ST-MoE CTC achieves a comparable performance with Conformer-based LAE [32] and an obvious gain compared to FLR-MoE CTC [22].

---

[4]https://github.com/wenet-e2e/wenet/blob/main/tools/compute-wer.py

**Table 2**. Results of proposed models and the baselines. The numbers in brackets indicates the relative error rate reduction comparing with the corresponding LAE-based model (S2 and S5).

| System | Model | Infer Params | Code-Switch | | | Mono | |
|---|---|---|---|---|---|---|---|
| | | | ALL | CN | EN | EN | CN |
| CTC-based ASR system | | | | | | | |
| Literature | | | | | | | |
| - | Conformer CTC [32] | - | 11.6 | - | - | - | - |
| | + LAE [32] | - | **9.5** | - | - | - | - |
| - | FLR-MoE CTC [22] | 25.8 M | 10.5 | 7.7 | 33.1 | 10.1 | 5.1 |
| Our results | | | | | | | |
| S1 | Vallina CTC | 19.8 M | 12.2 | 9.0 | 38.9 | 12.4 | 7.1 |
| S2 | LAE CTC (baseline) | 26.5 M | 10.8 | 8.0 | 33.7 | 10.5 | 5.3 |
| S3 | LAE-ST-MoE CTC (proposed) | 26.5 M | **9.8 (9.26% ↓)** | **7.3** | **30.3** | **9.6 (8.57% ↓)** | **4.9 (7.55% ↓)** |
| Attention-based ASR system | | | | | | | |
| Literature | | | | | | | |
| - | Hybrid CTC + Attention [21] | 28.8 M | 10.9 | 8.8 | 28.1 | - | - |
| | + Bi-En. (MoE-in-unsup) [21] | 45.6 M | 9.8 | 7.7 | 26.6 | - | - |
| - | FLR-MoE AED [22] | 40.7 M | 9.7 | 7.4 | 28.4 | 9.6 | **4.7** |
| Our results | | | | | | | |
| S4 | Vallina AED | 34.7 M | 11.2 | 8.6 | 32.5 | 11.7 | 6.3 |
| S5 | LAE AED (baseline) | 41.4 M | 10.0 | 7.7 | 29.2 | 9.9 | 5.0 |
| S6 | LAE-ST-MoE AED (proposed) | 41.4 M | **9.3 (7% ↓)** | **7.1** | **27.4** | **9.2 (7.07% ↓)** | **4.7 (6% ↓)** |

**AED System:** The results also show that our LAE-ST-MoE-based system (S6) performs better than the Vallina (S4) and LAE-based (S5) AED ASR. Moreover, the LAE-ST-MoE-based AED system (S6) also shows an obvious MER reduction compared with the Bi-encoder [21] based and FLR-MoE [22] based system on the CS test.

**CTC vs. AED system:** We can find that the proposed LAE-ST-MoE-based CTC (S3) shows a little performance gain to the LAE AED system (S5) and comparable results with Bi-Encoder [21] based and FLR-MoE [22] based AED system.

These results suggest that the ST auxiliary task can improve the ASR performance based on the LAE structure, which is consistent with our motivation.

### 4.3.2. Results of the w/ or w/o MoE in LAE-ST-MoE model

**Table 3**. Performance of the w/ or w/o MoE.

| Model | Code-Switch | | | Mono | | CS → EN | CS → CN |
|---|---|---|---|---|---|---|---|
| | ALL | CN | EN | EN | CN | BLEU | BLEU |
| LAE-ST CTC | 10.0 | 7.4 | 31.6 | 9.8 | 5.2 | 16.2 | 65.8 |
| + MoE | **9.8** | **7.3** | **30.3** | **9.6** | **4.9** | **17.7** | **66.6** |

Table 3's LAE-ST CTC model replaces the MoE layer in LAE-ST-MoE with a regular FFN. From the results, we can see that due to the introduction of the MoE module, the performance of ASR and ST is both improved obviously, which further confirms our motivation that introducing the MoE module will make ASR and ST tasks more focused.

### 4.3.3. Results of using ST decoder for ASR rescore

**Table 4**. Performance of using ST decoder rescore.

| Model | Code-Switch | | | Mono | |
|---|---|---|---|---|---|
| | ALL | CN | EN | EN | CN |
| Vallina CTC | 12.2 | 9.0 | 38.9 | 12.4 | 7.1 |
| LAE CTC | 10.8 | 8.0 | 33.7 | 10.5 | 5.3 |
| LAE-ST-MoE CTC | 9.8 | 7.3 | 30.3 | 9.6 | 4.9 |
| + En2Man ST Dec. res. | **9.7** | **7.1** | 31.2 | 10.2 | **4.8** |
| + Man2En ST Dec. res. | 10.4 | 8.1 | **29.1** | **9.3** | 5.6 |

The En2Man ST decoder is comparable to the Mandarin ASR decoder for monolingual Mandarin input. Therefore, we can easily fuse it into monolingual Mandarin decoding through rescoring. As shown in Table 4, the En2Man ST decoder improves the LAE-ST-MoE CTC system in the mono CN speech. It achieves comparable results to the LAE-ST-MoE AED system (Table 2's S6) in the monolingual Mandarin test. Especially on the Mandarin part of the CS test, the En2Man ST decoder rescoring performs better than the LAE-based AED system (Table 2's S5), which maybe benefit from the Mandarin-English context representation and the decoder LM-related information. In addition, the same phenomenon also can be observed when applying the Man2En ST decoder rescoring. These results show that the information learned by the ST decoder differs from that of the ASR decoder, improving the ASR performance. To a certain extent, the above results also prove the effectiveness of the LAE-ST-MoE.

### 4.3.4. Results of different $\beta$ and $N_{Mono}$ values in LAE-ST-MoE

As mentioned in section 3.2, $\beta$ is used to balance and regulate the ST effect. Therefore, in Table 5, we conduct experiments with $\beta$ values of 1.0, 0.8, 0.6, and 0.4, where we set $N_{Share}$ = 9 and $N_{Mono}$ = 1. From the results, it can be seen that the performance of CS is basically not affected, and the model has the best overall performance at 0.6.

**Table 5**. Results with different $\beta$ when $N_{Share}$ = 9 and $N_{Mono}$ = 1.

| Model | $\beta$ | Code-Switch | | | Mono | |
|---|---|---|---|---|---|---|
| | | ALL | CN | EN | EN | CN |
| Vallina CTC | - | 12.2 | 9.0 | 38.9 | 12.4 | 7.1 |
| LAE-ST-MoE CTC | 1.0 | **9.8** | **7.3** | **30.3** | 9.7 | 5.0 |
| LAE-ST-MoE CTC | 0.8 | **9.8** | **7.3** | **30.3** | 9.7 | 5.1 |
| LAE-ST-MoE CTC | 0.6 | **9.8** | **7.3** | **30.3** | **9.6** | **4.9** |
| LAE-ST-MoE CTC | 0.4 | 9.9 | 7.4 | 30.5 | 9.7 | 5.0 |

In addition, we set $\beta$ to 0.6 and $N_{Share}$ to 9. Then, the effectiveness of $N_{Mono}$ is investigated in Table 6. When $N_{Mono}$ is 0, the ASR and ST share all encoder layers except FFN-MoE. However, when $N_{Mono}$=2, the LAE-ST-MoE encoder layer will reduce to 1. From Table 6, we can see the model achieves the best in $N_{Mono}$=1, which suggests that the LAE-ST-MoE model needs more layers to perform ST, and it also needs to reserve some layers to learn the language-specific ASR representation.

**Table 6**. Results with different $N_{Mono}$ when $N_{Share}$ = 9 and $\beta$ = 0.6.

| Model | $N_{Mono}$ | Code-Switch | | | Mono | |
|---|---|---|---|---|---|---|
| | | ALL | CN | EN | EN | CN |
| Vallina CTC | - | 12.2 | 9.0 | 38.9 | 12.4 | 7.1 |
| LAE-ST-MoE CTC | 0 | 10.1 | 7.5 | 31.6 | 9.9 | 5.1 |
| LAE-ST-MoE CTC | 1 | **9.8** | **7.3** | **30.3** | **9.6** | **4.9** |
| LAE-ST-MoE CTC | 2 | 9.9 | 7.4 | 30.8 | 9.7 | 5.0 |

### 4.3.5. The results of ST auxiliary task in LAE-ST-MoE models

We use ModelScope's MT model to generate pseudo-labels for the test set. From Tables 7 and 8, which show the BLEU score and translation error rate (TER) of our models, we can see that ST is less affected by $\beta$ but more affected by $N_{Mono}$. Furthermore, by combining Tables 5, 6, 7, and 8, we can observe that when the ST BLEU change, the ASR remain basically unchanged. It may be because there is also some confusion between the information on ASR and ST. However, the helpful and confusing information needs to be balanced. Our experimental CS data is Mandarin-dominant, so we have more Mandarin-to-English ST training data than English-to-Mandarin, which results in better BLEU for Mandarin-to-English ST. Furthermore, we test the best ST model on monolingual data in Table 9, and we can see that our model also has

**Table 7**. ST results on the CS test when $N_{Share}$ = 9 and $N_{Mono}$ = 1.

| Model | $\beta$ | CS → EN | | CS → CN | |
|---|---|---|---|---|---|
| | | BLEU | TER ($\downarrow$) | BLEU | TER ($\downarrow$) |
| LAE-ST-MoE CTC | 1.0 | **18.4** | **69.6** | **67.0** | **21.3** |
| LAE-ST-MoE CTC | 0.8 | 18.1 | 70.0 | 66.8 | 21.5 |
| LAE-ST-MoE CTC | 0.6 | 17.7 | 70.3 | 66.6 | 21.6 |
| LAE-ST-MoE CTC | 0.4 | 17.3 | 70.6 | 66.2 | 21.9 |

**Table 8**. ST results on the CS test when $N_{Share}$ = 9 and $\beta$ = 0.6.

| Model | $N_{Mono}$ | CS → EN | | CS → CN | |
|---|---|---|---|---|---|
| | | BLEU | TER ($\downarrow$) | BLEU | TER ($\downarrow$) |
| LAE-ST-MoE CTC | 0 | **18.6** | **69.0** | **67.0** | **21.3** |
| LAE-ST-MoE CTC | 1 | 17.7 | 70.3 | 66.6 | 21.6 |
| LAE-ST-MoE CTC | 2 | 16.3 | 72.6 | 65.5 | 22.4 |

**Table 9**. ST results on the monolingual test.

| Model | CN → EN | | EN → CN | |
|---|---|---|---|---|
| | BLEU | TER ($\downarrow$) | BLEU | TER ($\downarrow$) |
| LAE-ST-MoE CTC | 33.9 | 44.8 | 31.5 | 59.1 |

CS → CN:
Audio text: 他的diary标题我都很喜欢
Modelscope (MT): 他的日记标题我都很喜欢
Our model (ST): 他的日记标题我都很喜欢
EN → CN:
Audio text: I say I've been wondering about this business
Modelscope (MT): 我说我一直在想这项业务
Our model (ST): 我说我一直在想这项业务

CS → EN:
他的diary标题我都很喜欢
I liked his diary title
His diary title I like it very much
CN → EN:
给我介绍几首好听的歌
Introduce me some nice songs
Introduce me some good songs

**Fig. 2**. The examples translated by ModelScope and our model respectively.

good BLEU. For CS data with limited English, the BLEU of CS speech to Mandarin text shows better than CS to English.

Figure 2 provides examples of the translation performed by ModelScope's MT model and our model's ST branch. Specifically, ModelScope's model translates text, whereas ours conducts the ST task. To a certain extent, these examples demonstrate that our model has learned good ST ability.

## 5. CONCLUSIONS

In this paper, we propose an LAE-ST-MoE framework that incorporates ST tasks into LAE and utilizes ST to learn the contextual information between different languages. The experimental results on the ASRU 2019 Mandarin-English CS challenge dataset demonstrate that, compared to the LAE-based CTC and AED system, the proposed LAE-ST-MoE model achieves about 6%-9% relative error rate reduction. Extensive investigations into the w/ or w/o MoE module, comparison with the literature results, and ablation on different $\beta$ and $N_{Mono}$ values have also been carried out and confirm the effectiveness of the LAE-ST-MoE. Moreover, the well-trained LAE-ST-MoE model can perform ST tasks from CS speech to Mandarin or English, and the structure is easy to extend to one-to-many ST tasks. In the future, we will further explore the LAE-ST-MoE to multilingual ASR and one-to-many ST.

# 6. REFERENCES

[1] Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber, "Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks," in *ICML*, New York, NY, USA, 2006, ICML '06, p. 369–376, Association for Computing Machinery.

[2] Alex Graves, "Sequence Transduction with Recurrent Neural Networks," *arXiv e-prints*, p. arXiv:1211.3711, Nov. 2012.

[3] W. Chan, N. Jaitly, Q. Le, and O. Vinyals, "Listen, attend and spell: A neural network for large vocabulary conversational speech recognition," in *ICASSP*, March 2016, pp. 4960–4964.

[4] Suyoun Kim, Takaaki Hori, and Shinji Watanabe, "Joint ctc-attention based end-to-end speech recognition using multi-task learning," in *ICASSP*, 2017, pp. 4835–4839.

[5] L. Dong, S. Xu, and B. Xu, "Speech-transformer: A no-recurrence sequence-to-sequence model for speech recognition," in *ICASSP*, 2018, pp. 5884–5888.

[6] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, and Ruoming Pang, "Conformer: Convolution-augmented Transformer for Speech Recognition," in *Proc. Interspeech*, 2020, pp. 5036–5040.

[7] Daniel S. Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D. Cubuk, and Quoc V. Le, "SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition," in *Proc. Interspeech*, 2019, pp. 2613–2617.

[8] Yuting Yang, Yuke Li, and Binbin Du, "Improving ctc-based asr models with gated interlayer collaboration," in *ICASSP 2023*, 2023, pp. 1–5.

[9] Guodong Ma, Pengfei Hu, Nurmemet Yolwas, Shen Huang, and Hao Huang, "PM-MMUT: Boosted Phone-mask Data Augmentation using Multi-Modeling Unit Training for Phonetic-Reduction-Robust E2E Speech Recognition," in *Proc. Interspeech 2022*, 2022, pp. 1021–1025.

[10] Rui Li, Guodong Ma, Dexin Zhao, Ranran Zeng, Xiaoyu Li, and Hao Huang, "A policy-based approach to the specaugment method for low resource e2e asr," in *2022 APSIPA ASC*, 2022, pp. 630–635.

[11] Kunal Dhawan, Dima Rekesh, and Boris Ginsburg, "Towards training Bilingual and Code-Switched Speech Recognition models from Monolingual data sources," *arXiv e-prints*, p. arXiv:2306.08753, June 2023.

[12] Zheng Liang, Zheshu Song, Ziyang Ma, Chenpeng Du, Kai Yu, and Xie Chen, "Improving Code-Switching and Name Entity Recognition in ASR with Speech Editing based Data Augmentation," in *Proc. INTERSPEECH 2023*, 2023, pp. 919–923.

[13] Chenpeng Du, Hao Li, Yizhou Lu, Lan Wang, and Yanmin Qian, "Data augmentation for end-to-end code-switching speech recognition," in *2021 IEEE SLT*, 2021, pp. 194–200.

[14] Haibin Yu, Yuxuan Hu, Yao Qian, Ma Jin, Linquan Liu, Shujie Liu, Yu Shi, Yanmin Qian, Edward Lin, and Michael Zeng, "Code-switching text generation and injection in mandarin-english asr," in *ICASSP 2023*, 2023, pp. 1–5.

[15] Thien Nguyen, Nathalie Tran, Liuhui Deng, Thiago Fraga da Silva, Matthew Radzihovsky, Roger Hsiao, Henry Mason, Stefan Braun, Erik McDermott, Dogan Can, et al., "Optimizing bilingual neural transducer with synthetic code-switching text generation," *arXiv preprint arXiv:2210.12214*, 2022.

[16] Shun-Po Chuang, Heng-Jui Chang, Sung-Feng Huang, and Hung-yi Lee, "Non-autoregressive mandarin-english code-switching speech recognition," in *2021 IEEE ASRU*, 2021, pp. 465–472.

[17] Zhiping Zeng, Yerbolat Khassanov, Van Tung Pham, Haihua Xu, Eng Siong Chng, and Haizhou Li, "On the End-to-End Solution to Mandarin-English Code-Switching Speech Recognition," in *Proc. Interspeech 2019*, 2019, pp. 2165–2169.

[18] Yizhou Peng, Jicheng Zhang, Haihua Xu, Hao Huang, and Eng Siong Chng, "Minimum word error training for non-autoregressive transformer-based code-switching asr," in *ICASSP 2022*. IEEE, 2022, pp. 7807–7811.

[19] Hexin Liu, Haihua Xu, Leibny Paola Garcia, Andy W. H. Khong, Yi He, and Sanjeev Khudanpur, "Reducing language confusion for code-switching speech recognition with token-level language diarization," in *ICASSP 2023*, 2023, pp. 1–5.

[20] Zhiyun Fan, Linhao Dong, Chen Shen, Zhenlin Liang, Jun Zhang, Lu Lu, and Zejun Ma, "Language-specific Boundary Learning for Improving Mandarin-English Code-switching Speech Recognition," in *Proc. INTERSPEECH 2023*, 2023, pp. 3322–3326.

[21] Y. Lu, M. Huang, H. Li, J. Guo, and Y. Qian, "Bi-encoder transformer network for mandarin-english code-switching speech recognition using mixture of experts," in *Interspeech*, 2020.

[22] Wenxuan Wang, Guodong Ma, Yuke Li, and Binbin Du, "Language-Routing Mixture of Experts for Multilingual and Code-Switching Speech Recognition," in *Proc. INTERSPEECH 2023*, 2023, pp. 1389–1393.

[23] Zhao You, Shulin Feng, Dan Su, and Dong Yu, "Speech-MoE: Scaling to Large Acoustic Models with Dynamic Routing Mixture of Experts," in *Proc. Interspeech 2021*, 2021, pp. 2077–2081.

[24] Yoohwan Kwon and Soo-Whan Chung, "Mole : Mixture of language experts for multi-lingual automatic speech recognition," in *ICASSP 2023*, 2023, pp. 1–5.

[25] Guodong Ma, Pengfei Hu, Jian Kang, Shen Huang, and Hao Huang, "Leveraging Phone Mask Training for Phonetic-Reduction-Robust E2E Uyghur Speech Recognition," in *Proc. Interspeech 2021*, 2021, pp. 306–310.

[26] Haobo Zhang, Haihua Xu, Van Tung Pham, Hao Huang, and Eng Siong Chng, "Monolingual Data Selection Analysis for English-Mandarin Hybrid Code-Switching Speech Recognition," in *Proc. Interspeech 2020*, 2020, pp. 2392–2396.

[27] Ahmed Ali, Shammur Chowdhury, Amir Hussein, and Yasser Hifny, "Arabic code-switching speech recognition using monolingual data," *arXiv preprint arXiv:2107.01573*, 2021.

[28] Xianghu Yue, Grandee Lee, Emre Yılmaz, Fang Deng, and Haizhou Li, "End-to-end code-switching asr for low-resourced language pairs," in *2019 IEEE ASRU*, 2019, pp. 972–979.

[29] Brian Yan, Matthew Wiesner, Ondřej Klejch, Preethi Jyothi, and Shinji Watanabe, "Towards zero-shot code-switched speech recognition," in *ICASSP 2023*, 2023, pp. 1–5.

[30] Brian Yan, Chunlei Zhang, Meng Yu, Shi-Xiong Zhang, Siddharth Dalmia, Dan Berrebbi, Chao Weng, Shinji Watanabe, and Dong Yu, "Joint modeling of code-switched and monolingual asr via conditional factorization," in *ICASSP 2022*. IEEE, 2022, pp. 6412–6416.

[31] Tongtong Song, Qiang Xu, Meng Ge, Longbiao Wang, Hao Shi, Yongjie Lv, Yuqin Lin, and Jianwu Dang, "Language-specific Characteristic Assistance for Code-switching Speech Recognition," in *Proc. Interspeech 2022*, 2022, pp. 3924–3928.

[32] Jinchuan Tian, Jianwei Yu, Chunlei Zhang, Yuexian Zou, and Dong Yu, "LAE: Language-Aware Encoder for Monolingual and Multilingual ASR," in *Proc. Interspeech 2022*, 2022, pp. 3178–3182.

[33] Xian Shi, Qiangze Feng, and Lei Xie, "The asru 2019 mandarin-english code-switching speech recognition challenge: Open datasets, tracks, methods and results," *arXiv preprint arXiv:2007.05916*, 2020.

[34] Xiangpeng Wei, Heng Yu, Yue Hu, Rongxiang Weng, Weihua Luo, and Rong Jin, "Learning to generalize to more: Continuous semantic augmentation for neural machine translation," in *ACL 2022*, 2022.

[35] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton, "Layer normalization," *arXiv preprint arXiv:1607.06450*, 2016.

[36] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur, "Librispeech: an asr corpus based on public domain audio books," in *2015 ICASSP*. IEEE, 2015, pp. 5206–5210.

[37] Daniel S Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D Cubuk, and Quoc V Le, "Specaugment: A simple data augmentation method for automatic speech recognition," *Proc. Interspeech 2019*, pp. 2613–2617, 2019.

[38] Rico Sennrich, Barry Haddow, and Alexandra Birch, "Neural machine translation of rare words with subword units," *arXiv preprint arXiv:1508.07909*, 2015.

[39] Zhuoyuan Yao, Di Wu, Xiong Wang, Binbin Zhang, Fan Yu, Chao Yang, Zhendong Peng, Xiaoyu Chen, Lei Xie, and Xin Lei, "WeNet: Production Oriented Streaming and Non-Streaming End-to-End Speech Recognition Toolkit," in *Proc. Interspeech 2021*, 2021, pp. 4054–4058.

[40] Matt Post, "A call for clarity in reporting BLEU scores," in *Proceedings of the Third Conference on Machine Translation: Research Papers*, Belgium, Brussels, Oct. 2018, pp. 186–191, Association for Computational Linguistics.

[41] Shinji Watanabe, Takaaki Hori, Shigeki Karita, Tomoki Hayashi, Jiro Nishitoba, Yuya Unno, Nelson Enrique Yalta Soplin, Jahn Heymann, Matthew Wiesner, Nanxin Chen, et al., "Espnet: End-to-end speech processing toolkit," *arXiv preprint arXiv:1804.00015*, 2018.

[42] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

[43] Kenneth Heafield, "Kenlm: Faster and smaller language model queries," in *Proceedings of the sixth workshop on statistical machine translation*, 2011, pp. 187–197.