

DialectMoE: An End-to-End Multi-Dialect Speech Recognition Model with Mixture-of-Experts

Anonymous ACL submission

Abstract

Dialect speech recognition has always been one of the challenges in Automatic Speech Recognition (ASR) systems. While lots of ASR systems perform well in Mandarin, their performance significantly drops when handling dialect speech. This is mainly due to the obvious differences between dialects and Mandarin in pronunciation and the data scarcity of dialect speech. In this paper, we propose DialectMoE, a Chinese multi-dialects speech recognition model based on Mixture-of-Experts (MoE) in a low-resource conditions. Specifically, DialectMoE assigns input sequences to a set of experts using a dynamic routing algorithm, with each expert potentially trained for a specific dialect. Subsequently, the outputs of these experts are combined to derive the final output. Due to the similarities among dialects, distinct experts may offer assistance in recognizing other dialects as well. Experimental results on the Ai-Datatang dialect public dataset show that, compared with the baseline model, DialectMoE reduces Character Error Rate(CER) for Sichuan, Yunnan, Hubei and Henan dialects by 23.6%, 32.6%, 39.2% and 35.09% respectively. The proposed DialectMoE model demonstrates outstanding performance in multi-dialects speech recognition.

1 Introduction

The application domains of speech recognition technology are extensive, encompassing diverse fields such as voice assistants, smart homes, and automotive voice interaction, among others. Thanks to the advancements in deep learning, Automatic Speech Recognition(ASR) systems have made remarkable strides in recognizing Mandarin speech(Malik et al., 2021; Wang et al., 2019; Alharbi et al., 2021). Dialect serves as a prevalent mode of everyday communication among the Chinese populace. However, the performance of ASR systems remains limited in dialect speech, posing a

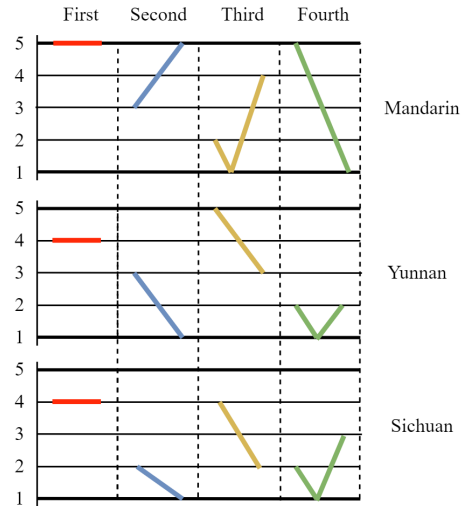


Figure 1: The tonal distinctions among Standard Mandarin, Yunnan dialect, and Sichuan dialect.

significant challenge in the field of speech recognition technology(Hinsvark et al., 2021; Alsharhan and Ramsay, 2020) due to the inherent variations and distinct characteristics in pronunciation among dialects and Mandarin. Therefore, improving the accuracy and adaptability of Chinese ASR systems is significant and meaningful for multi-dialect. Our study mainly focuses on Chinese dialects, the proposed method can also be generalized to other dialects.

Chinese dialects are typically classified into ten main categories, each exhibiting notable differences in pronunciation, tone, vocabulary, and grammar(Ho, 2015). Chinese is a tonal language, where each character corresponds to a specific tone, a feature that is prevalent in most of its dialects as well. The pronunciation of a given Chinese character with different tones imparts markedly distinct meanings. This underscores the profound significance of tones in the comprehension of Chinese phonetics. (Ho, 2015; Sproat et al., 2004). Figure 1 depicts the tonal distinctions among Standard Mandarin, Yunnan dialect, and Sichuan di-

065 alect. It is evident that notable differences exist
066 in tonal between Standard Mandarin and Yunnan
067 dialect as well as Sichuan dialect in the second,
068 third, and fourth tones. However, Yunnan dialect
069 and Sichuan dialect exhibit a pronunciation sim-
070 ilarity in specific tones. The change in tone re-
071 veal differences and similarities between Standard
072 Mandarin and various dialects. Hence, considering
073 both the differences and similarities in pronun-
074 ciation among various dialects alongside Standard
075 Mandarin becomes crucial for the advancement of
076 Chinese speech recognition systems.

077 In recent years, numerous researches have fo-
078 cused on tackling the challenge of poor perfor-
079 mance in dialect speech recognition models(Li
080 et al., 2018; Ren et al., 2019; Zhang et al., 2022b).
081 The traditional way is based on different modeling
082 methods to improve the effect of dialect speech
083 recognition. Humphries et al. (1996) employed an
084 adaptive method that utilizes a pronunciation vo-
085 cabulary with dialect data to capture differences
086 between standard and dialect pronunciations. Li
087 et al. (2019) proposed a novel method for modeling
088 Chinese characters based on radicals, effectively
089 addressing the issue of dialect modeling difficulty.
090 This method significantly reduces the required size
091 of radical dictionaries compared to ordinary char-
092 acter dictionaries. Recently, multitask-based meth-
093 ods have been widely used in the task of dialect
094 speech recognition. Compared with the traditional
095 method, the multi-task learning method is more
096 efficient. Elfeky et al. (2016) proposed construct-
097 ing a dialect classification model and a separate
098 speech recognition model for each dialect. The
099 dialect classification model is used to select the
100 corresponding dialect speech recognition model.
101 Dan et al. (2022) proposed a multi-task training
102 strategy that combines dialect classification with
103 dialect speech recognition, bridging the substantial
104 gap between Mandarin and dialect acoustic prop-
105 erties. However,, these investigations are contingent
106 upon extensive dialectal datasets and do not exam-
107 ine the potential influence of commonalities among
108 various dialects on model performance.

109 To construct a reliable dialect speech recognition
110 model in low-resource conditions, Jiang (2023) in-
111 troduced a transfer learning-based approach, it in-
112 volves a model trained on Mandarin and fine-tunes
113 with small-scale dialect data. However, relying
114 solely on transfer learning may not adequately cap-
115 ture the distinctions between dialects and Mandarin.

116 Wang et al. (2023) proposed Aformer model with
117 multi-stage training strategy, which can capture
118 diverse acoustic information in different training
119 stage, enabling the model to effectively adapt to
120 dialect data. The aforementioned studies focus on
121 the training strategy and model expansion, they do
122 not fully consider the differences and similarities
123 between dialects and Mandarin.

124 In this paper, we present DialectMoE, a multi-
125 dialect speech recognition model based on Mixture-
126 of-Experts (MoE), aimed at improving the perfor-
127 mance of multi-dialects speech recognition in low-
128 resource conditions. DialectMoE is architecturally
129 structured with dual encoders: a dialect encoder
130 and a general encoder. The main contributions of
131 the paper include:

- 132 • We propose a three-stage training methodol-
133 ogy designed to enhance the model’s adapt-
134 ability in addressing low-resource multi-
135 dialect scenarios through different stages. De-
136 tailed specifics will be expounded upon in
137 Section 3.3.
- 138 • We introduce MoE layers and enhance the
139 dynamic routing algorithm to enable the com-
140 bination of acoustic features from both the
141 input sequence and the dialect encoder during
142 the expert selection process.
- 143 • The experiment results show that DialectMoE
144 reduces Character Error Rate(CER) compared
145 to the baseline model for Sichuan, Yunnan,
146 Hubei and Henan dialects by 23.6%, 32.6%,
147 39.2% and 35.09%, respectively.

148 2 Related Work

149 2.1 Conformer-based ASR

150 The Conformer, a convolution-augmented Trans-
151 former introduced in (Gulati et al., 2020), has been
152 widely acknowledged as the state-of-the-art end-
153 to-end ASR technology owing to its exceptional
154 performance in ASR tasks. In recent years, sev-
155 eral researchers have proposed Conformer variants
156 (Peng et al., 2022; Sehoon et al., 2023) to further
157 enhance the capabilities of speech recognition. The
158 Conformer module comprises two feed-forward
159 modules, a multi-heads self-attention module, and
160 a convolution module. The output y of one Con-
161 former block for a given input x can be defined as

162 follows:

$$163 \quad \hat{x} = x + \frac{1}{2} \text{FFN}_1(x) \quad (1)$$

$$164 \quad \tilde{x} = \hat{x} + \text{MHSA}(\hat{x}) \quad (2)$$

$$165 \quad \bar{x} = \tilde{x} + \text{Conv}(\tilde{x}) \quad (3)$$

$$166 \quad y = \text{LN}(\bar{x} + \frac{1}{2} \text{FNN}_2(\bar{x})) \quad (4)$$

167 where FNN_1 denotes the first feedforward module,
168 FNN_2 denotes the second feedforward network,
169 MHSA denotes the multi-head self-attention mod-
170 ule, Conv denotes the convolution module, and
171 LN denotes layer normalization. For additional
172 information regarding the Conformer ASR model,
173 please refer to (Gulati et al., 2020).

174 During Conformer training, the Joint CTC-
175 Attention loss function (Hori et al., 2017) is utilized.
176 This loss function is commonly used in present-day
177 speech recognition technology. In this paper, the
178 joint CTC-Attention loss is incorporated into the
179 total loss function. The loss function is outlined as
180 follows:

$$181 \quad \mathcal{L}_{all} = (1 - \lambda)\mathcal{L}_{att} + \lambda\mathcal{L}_{ctc} \quad (5)$$

182 where \mathcal{L}_{att} denotes the decoding loss of the Atten-
183 tion decoder, and \mathcal{L}_{ctc} denotes the CTC loss., λ is a
184 hyper parameter which denotes the weight of these
185 two loss.

186 2.2 Mixture-of-Experts Based Speech 187 Recognition

188 The MoE based methods offer a solution for more
189 efficient training and inference by selectively acti-
190 vating different experts in the model based on differ-
191 ent inputs (Jacobs et al., 1991; Shazeer et al., 2017).
192 This enables the model to adapt to a wide range of
193 inputs and scale to more parameters while main-
194 taining a consistent computational cost. The MoE
195 based models have demonstrated their effectiveness
196 in natural language processing (Fedus et al., 2022;
197 Du et al., 2022) and computer vision (Riquelme
198 et al., 2021; Fan et al., 2022).

199 In real-world application, speech recognition sys-
200 tems need to handle various input conditions, in-
201 cluding speaker variation, accent variation, and
202 acoustic environment (Zilvan et al., 2021). How-
203 ever, conventional speech recognition models have
204 a fixed computational cost and cannot adapt to the
205 complexity of input instances. You et al. (2021,
206 2022) explore the MoE based model for speech
207 recognition, named SpeechMoE, and propose a

208 new router architecture which integrates additional
209 global domain and various embedding into router
210 input to promote adaptability. Additionally, a mul-
211 tilingual speech recognition network (MoLE) was
212 introduced (Kwon and Chung, 2023) to analyze au-
213 dio input data from multiple languages and identify
214 expert networks suitable for each language. Simul-
215 taneously, a language-independent expert network
216 was also introduced, and the selected expert net-
217 work and the language-independent expert network
218 collectively fulfill the language requirements nec-
219 essary for effective speech recognition.

220 Employing the MoE mechanism to determine
221 expert activation during the forward propagation
222 process manifests a notable capacity for accom-
223 modating the inherent variability in multi-dialectal
224 speech across different input sequences. Neverthe-
225 less, the conventional MoE paradigm relies solely
226 on the input sequence for expert selection, and the
227 information in the present input does not inher-
228 ently ensure the optimal suitability of the selected
229 experts. Therefore, the incorporation of supple-
230 mentary dialectal information to facilitate expert
231 selection stands forth as a judicious resolution, en-
232 hancing the precision and adaptability of the chosen
233 experts to the distinctive intricacies characterizing
234 the multi-dialectal speech context. Furthermore,
235 the exploration of MOE-based methods in the do-
236 main of multi-dialect speech recognition remains
237 limited.

238 3 DialectMoE

239 3.1 Overall Architecture of DialectMoE

240 The overall architecture of DialectMoE is shown in
241 Figure 2a. The original audio sequence undergoes
242 preprocessing by the frontend module to extract
243 FBank features. Subsequently, the convolutional
244 downsampling is applied to temporally downsam-
245 ple the audio feature sequence. The dialect encoder,
246 consisting of 6 layers of vanilla Conformer encoder,
247 captures dialect information from the feature se-
248 quences. On the other hand, the general encoder
249 comprises 12 layers of DialectMoE encoder blocks,
250 are responsible for capturing speech information
251 in a normal manner. Both encoders share the same
252 input but focus on different aspects of information.

253 The detailed structure of the DialectMoE block
254 is presented in Figure 2b. With the DialectMoE
255 block, the input sequence firstly passes through
256 the Feed-Forward Network (FFN) layer, followed
257 by Attention and Convolutional Neural Network

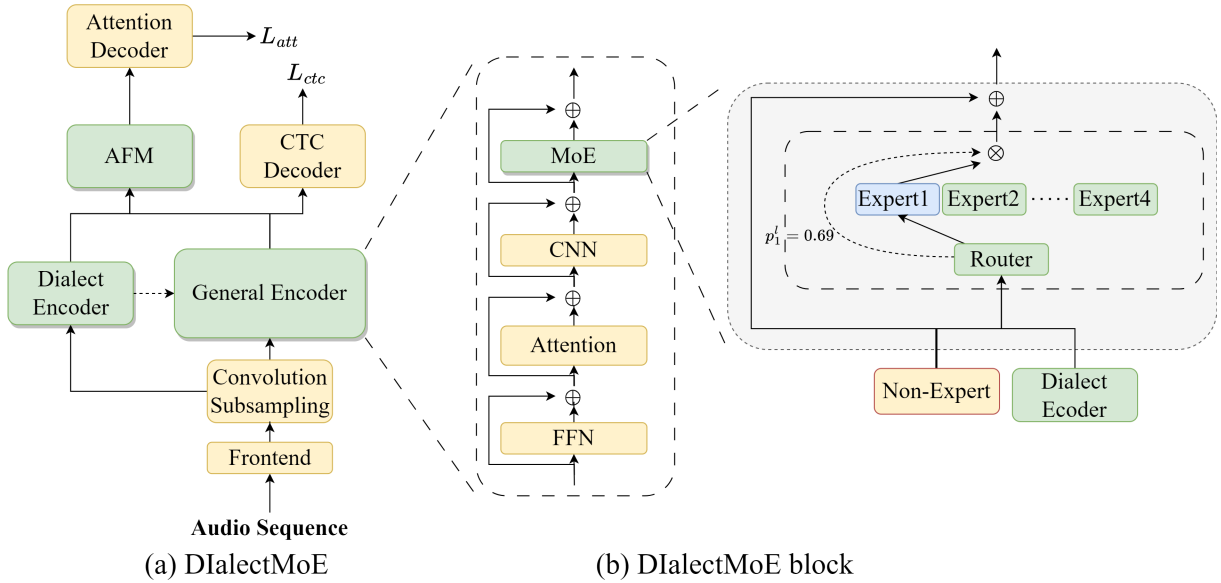


Figure 2: (a) DialectMoE overall architecture, where the general encoder consists of N DialectMoE blocks. (b) Architecture of the DialectMoE encoder block module.

(CNN) layer to extract global and local information, respectively. Then the appropriate expert within the MoE layer is selected based on the dynamic routing. The output of experts are multiplied by the weight assigned by the router layer.

Compared to the widely used vanilla Conformer block (Gulati et al., 2020), our DialectMoE block incorporates MoE layers to address complex and variable scenarios encountered in real-world situations. The dialect information captured by the dialect encoder is weighted by the router layer, which enables the router layer to choose more appropriate experts based on both dialect features and general features obtained from two encoders. This dynamic routing mechanism proves more effective in intricate speech scenarios, especially those involving multiple dialects.

3.2 Dialect Adaptive Dynamic Expert Routing

In the context of multi-dialect speech recognition, effectively addressing the diversity of dialectal variations is crucial. We present a novel dynamic routing algorithm aimed at enhancing the adaptability and generalization of the model to diverse dialects speech input sequences. The proposed algorithm leverages the input sequences from the current MoE layer and the dialect information provided by the dialect encoder to select appropriate experts. To evaluate the impact of different dialect embedding on routing, we consider the following three strategies: input sequence **concat** dialect embedding,

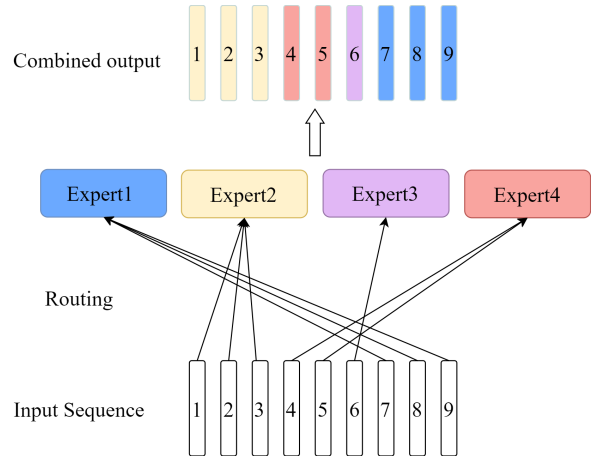


Figure 3: Illustration of dynamic routing algorithm.

input sequence **add** dialect embedding, and only use embedding (**embed**) of dialects. The output of the dialect encoder is denoted as $\mathbf{X}_{encoder}^D \subseteq \mathbb{R}^{T \times d}$, where T represents the sequence length and d denotes the feature dimension. Assuming that there are N experts, the output $r \subseteq \mathbb{R}^{T \times N}$ of the routing layer can be defined as follows. :

$$r = W_r \cdot \text{Concat}(\bar{x}, \mathbf{X}_{encoder}^D) \quad (6)$$

$$r = W_r \cdot \text{Add}(\bar{x}, \mathbf{X}_{encoder}^D) \quad (7)$$

$$r = W_r \cdot \mathbf{X}_{encoder}^D \quad (8)$$

where W_r represents the weight parameter of the router layer, and \bar{x} denotes the output of the convolution module. These three dynamic routing strategies are the ones we consider employing. It is worth

noting that while the general router layer selects experts based on the input sequence \bar{x} , the algorithm we designed intuitively makes more sense as it incorporates the output of the dialect encoder to select the most suitable expert.

The router layer selects the expert with the highest probability through dynamic routing, which is based on the input sequence r . The dynamic routing probability is then defined as follows:

$$p_i = \frac{\exp^{r_i}}{\sum_{j=1}^N \exp^{r_j}} \quad (9)$$

where $p_i \subseteq \mathbb{R}^{T \times N}$ is the probability that the i expert is selected, the output $\mathbf{O}_{moe} \subseteq \mathbb{R}^{T \times d}$ of the MoE layer can be formally defined as follows:

$$\mathbf{O}_{moe} = p_i \cdot E_i(\bar{x}) \quad (10)$$

where E_i is the output of the i expert selected. Figure 3 illustrates the process of dynamic routing.

In order to incorporate dialect information into the decoder, DialectMoE incorporates an information fusion step by combining the outputs of two separate encoders. This fusion process, illustrated as the Acoustic Fusion Module (AFM) in Figure 2(a), occurs prior to transmitting the results to the decoder. The fusion process is defined as follows.

$$\mathbf{X}_{encoder}^A = \text{Concat}(\mathbf{X}_{encoder}^G, \mathbf{X}_{encoder}^D) \quad (11)$$

where $\mathbf{X}_{encoder}^A$ denotes the result of fusion of information output by two different encoders, and $\mathbf{X}_{encoder}^G$ denotes the result output by a general encoder.

The comprehensive loss function for speech recognition comprises the combined CTC-Attention loss (Hori et al., 2017), as explained previously, along with the supplementary balance loss (Fedus et al., 2022). The complete formulation of the loss function is as follows:

$$\mathcal{L}_{all} = \lambda \mathcal{L}_{ctc} + (1 - \lambda) \mathcal{L}_{att} + \alpha \mathcal{L}_b \quad (12)$$

where α is the weight of the balance loss ($\alpha = 0.1$) and λ is the weight of the speech recognition loss ($\lambda = 0.3$), \mathcal{L}_b denotes the balance loss.

3.3 Training Strategies

Considering the significant disparity in the quantities of Mandarin and dialect data, low-resource dialect speech recognition scenarios commonly exhibit limited labeled dialect speech data, typically

ranging from a few to tens of hours. This insufficiency hampers the development of a reliable speech recognition model. To address this issue, this study introduces a multi-stage training strategy. The training process encompasses the following sequential steps:

1. Pre-training: The Conformer model is used as a general encoder for DialectMoE to implement pre-training on Mandarin datasets. The pre-training step allows the model to capture various common speech features, thus reducing the complexity of learning for the dialect recognition task.
2. Training Dialect Encoder: A Conformer Encoder is initialized as a dialect encoder and is trained on the dialect classification task using both dialect and Mandarin data. The objective of this step is to enable the dialect encoder to learn the acoustic differences between Mandarin and dialects, assisting the general encoder in dialect speech recognition tasks.
3. Training DialectMoE: The parameters of the dialect encoder are frozen, and the second feedforward network layer in the pre-trained Conformer model is initialized with N experts. Use only low-resource dialect training data to train the final DialectMoE model.

By pre-training phase, the initial model acquires a substantial set of effective parameters, thereby conferring notable advantages for last training stages. In the second phase, the dialect encoder is trained on a dialect identification task, enabling it to focus on differences between multi-dialects and Mandarin. In the last stage, only multi-dialect data is used for training, This strategic approach enhances DialectMoE’s capability to adeptly capture shared acoustic characteristics across diverse dialects. This approach enhances DialectMoE’s capability to adeptly capture shared acoustic characteristics across multi-dialects. The proposed method is evaluated based on extensive comparison and ablation experiments, which are comprehensively detailed in Section 4.

4 Experiments

4.1 Datasets

The Aishell dataset (Bu et al., 2017) serves as the Mandarin speech corpus in this study. This extensive collection of Mandarin Chinese speech data,

| Dataset | Train(h) | Test(h) |
|-------------|----------|---------|
| Aishell | 164 | 10 |
| Sichuan(SC) | 28.5 | 1.5 |
| Yunnan(YN) | 28.5 | 1.5 |
| Henan(HN) | 0 | 1.5 |
| Hubei(HB) | 0 | 1.5 |

Table 1: Details of both Dialect and Mandarin datasets.

encompasses diverse acoustic scenarios such as reading and dialogue.

For the Chinese dialect dataset, an open-source dataset provided by AiDatatang¹ is utilized in this study. It comprises a training set of 30 hours of Sichuan and Yunnan dialects and a test set of 1.5 hours featuring Henan and Hubei dialects. Within this study, Sichuan(SC) and Yunnan(YN) dialects are used to test the adaptability of the model to multi-dialect data, and Henan(HN) and Hubei(HB) dialects are used to test the generalization of the model to multi-dialect data. More details are shown in Table 1.

4.2 Experiment Setup

All experiments were conducted using the Wenet(Zhang et al., 2022a) end-to-end speech toolkit. Our methodology involved extracting an 80-dimensional log-Mel filter bank (Fbank) as the acoustic input feature, with a window size of 25 ms and a step size of 10 ms. To ensure feature normalization, we applied cepstrum mean variance normalization (CMVN) calculated from the training set on Fbank. To augment the low-resource dialect data, we employed speed perturbation and SpecAugment(Park et al., 2019) techniques. No additional language models were incorporated into the experiments.

For the pre-training model, we utilized a Conformer encoder trained on the Mandarin dataset. The general encoder of DialectMoE consists of 12 Conformer encoder layers with a feed-forward dimension of 2048 and an attention dimension of 256, employing 4 self-attention heads. This model was trained using the Adam optimizer (Kingma and Ba, 2014). Furthermore, we adopted the warmup learning schedule (Gotmare et al., 2018) for the initial 25K training iterations and set the label smoothing (Szegedy et al., 2016) weight and dropout to 0.1 for model regularization. The decoder consists of a 6-layer Transformer, while the dialect encoder

¹<https://www.datatang.com>

comprises a 6-layer Conformer encoder. The loss function for the dialect classification task applies cross-entropy loss to all training datasets.

The proposed DialectMoE is initialized with pre-trained general encoder, dialect encoder, and decoder. The second feedforward layer in each Conformer layer of the general encoder is initialized as an N expert ($N = 4$), with the expert parameters being the pre-trained feedforward network parameters. Training employed the same Adam optimizer, and the number of warm-up steps in the pre-training learning plan was adjusted to 10000, with an initial learning rate of 0.001.

4.3 Main Results

In this paper, we meticulously design comparison experiments with other speech recognition models to showcase the effectiveness of our proposed method. The experimental results presented in this study were reproduced using the open-source speech processing toolkit Wenet (Zhang et al., 2022a). Table 2 illustrates the performance of each ASR model in dialect speech recognition under low-resource conditions. The evaluation metric employed is the Character Error Rate (CER).

M1 represents the Conformer model that was exclusively pre-trained on the Aishell Mandarin dataset, consisting of 178 hours of data.

M2 denotes the model fine-tuned from M1 using the low-resource dialect dataset.

M3 corresponds to the model trained directly on the combined dataset of both dialect and Mandarin speech.

M4 and **M5** refer to the multi-tasking models trained on the combined dataset, with a distinction that M4 predicts the dialect category in the encoder while the decoder focuses on recognizing the speech text, whereas M5 predicts both the dialect category and the speech text in the decoder.

M6 represents the multi-pass model proposed in (Wang et al., 2023) for the training of the Aformer.

M7 signifies the DialectMoE model proposed in this paper.

The results obtained from the M1 model demonstrate a notably poor performance in recognizing dialectal speech within the Mandarin speech recognition model. However, by fine-tuning the M1 model with dialect data, the CER of the M2 model for Sichuan and Yunnan dialects is significantly reduced, although further optimization is still required. To address this, our paper proposes the Di-

| ID | Model | Params(M) | SC | YN | HN | HB |
|----|------------------------------|-----------|--------------|-------------|--------------|--------------|
| M1 | Conformer (Mandarin) | 46.1 M | 82.75 | 81.42 | 82.26 | 87.47 |
| M2 | FT-Conformer | 46.1 M | 15.60 | 14.06 | 54.91 | 57.18 |
| M3 | Conformer (Mandarin+Dialect) | 46.1 M | 13.86 | 12.02 | 41.28 | 48.37 |
| M4 | MT-Conformer(DID+ASR) | 47.2 M | 17.79 | 15.64 | 47.78 | 56.55 |
| M5 | MT-Conformer(DID&ASR) | 46.2 M | 16.09 | 14.05 | 41.28 | 48.37 |
| M6 | Aformer | 68.3 M | 13.21 | 12.76 | 35.32 | 39.89 |
| M7 | DialectMoE | 93.8 M | 11.91 | 9.84 | 33.38 | 37.11 |

Table 2: CER(%) on Chinese Dialect ASR task. **FT** represents the fine-tuning step and **MT** represents the multitask-based approach.

| Strategy | SC | YN | HN | HB |
|---------------|--------------|-------------|--------------|--------------|
| normal | 13.43 | 11.83 | 35.22 | 38.67 |
| embed | 12.53 | 10.20 | 34.95 | 38.65 |
| concat | 12.18 | 9.97 | 33.55 | 37.09 |
| add | 12.41 | 10.36 | 34.19 | 37.23 |
| normal+fusion | 13.92 | 12.19 | 35.27 | 38.72 |
| embed+fusion | 12.23 | 10.21 | 33.50 | 37.55 |
| concat+fusion | 11.91 | 9.84 | 33.38 | 37.11 |
| add+fusion | 12.49 | 10.13 | 33.66 | 37.65 |

Table 3: Ablation of different routing strategy.

483 alectMoE model, which surpasses existing studies
484 and baselines in terms of performance. In compari-
485 son to the fine-tuned model of M2, the DialectMoE
486 model exhibits a reduction in CER of 23.6% and
487 32.6% for the Sichuan and Yunnan dialects, re-
488 spective. Additionally, it achieves a reduction
489 of 39.2% and 35.09% for the Henan and Hubei
490 dialects, respectively.

4.4 Ablation Studies

4.4.1 Ablation of dynamic routing strategy

493 This paper incorporates ablation experiments to
494 investigate the effectiveness of the proposed dy-
495 namic routing algorithm and model design. Table
496 3 presents the impact of utilizing different dynamic
497 routing algorithms and the merging of two encoder
498 outputs before the decoder. In "normal", the dy-
499 namic routing algorithm proposed in this paper is
500 not employed, and the experts are directly selected
501 based on the input sequence, similar to the ap-
502 proach in (Fedus et al., 2022). The strategy column
503 in Table3 indicates the usage of different dynamic
504 routing algorithms: "**embed**" signifies the utiliza-
505 tion of only the dialect encoder outputs, "**concat**"
506 denotes the concatenation of the dialect encoder
507 outputs with the input sequence, and "**add**" indi-
508 cates the summation of the dialect encoder outputs
509 with the input sequence. The "**fusion**" entry indi-

| Model | Params(M) | SC | YN | HN | HB |
|--------|-----------|--------------|-------------|--------------|--------------|
| MoE-2e | 68.5M | 12.39 | 10.21 | 33.84 | 38.37 |
| MoE-4e | 93.8M | 11.91 | 9.84 | 33.38 | 37.11 |
| MoE-8e | 134.5M | 12.94 | 10.44 | 34.09 | 38.26 |

Table 4: Ablation of experts number.

510 cates whether or not the two encoder outputs should
511 be fused before reaching the decoder., whether they
512 go through AFM. The experiments employing the
513 "**concat+fusion**" strategy along with the fusion of
514 the two encoder outputs demonstrate optimal re-
515 sults across the four different dialect test sets.

4.4.2 Ablation of experts number

517 To investigate the impact of initializing a different
518 number of experts in DialectMoE on the overall
519 model performance, we conducted an experiment
520 with varying numbers of experts, specifically 2, 4,
521 and 8. The experimental results, as shown in Ta-
522 ble4, highlight that the model size increases with an
523 increasing number of experts. However, the com-
524 mon notion that a larger number of model partici-
525 pants leads to improved performance does not hold
526 true under low-resource conditions. The results
527 indicate that, for low-resource dialect data, an ex-
528 cessive number of experts does not enhance model
529 performance; in fact, it diminishes it. Experimental
530 evidence supports the conclusion that setting the
531 number of experts to 4 is more appropriate in this
532 context. It is noteworthy that when the number of
533 experts is set to 2, the number of model param-
534 eters matches the number of Aformer(Wang et al.,
535 2023) parameters. However, despite this similarity,
536 our results outperform the baseline. This finding
537 further validates the efficacy and correctness of our
538 proposed method.

| Top-k | Time(s) | SC | YN | HN | HB |
|-------|--------------|--------------|-------------|--------------|--------------|
| 4 | 1.46s | 12.01 | 10.03 | 33.31 | 37.08 |
| 2 | 0.94s | 11.93 | 9.81 | 33.57 | 37.24 |
| 1 | 0.68s | 11.91 | 9.84 | 33.38 | 37.11 |

Table 5: Ablation of the number of experts selected.

4.4.3 Ablation of the number of experts selected

In the vanilla MoE, a top-k approach is employed to select a combination of k experts for routing the input sequence. However, in this paper, a Softmax approach, specifically top-1, is utilized. To further investigate the effectiveness of the proposed dynamic routing algorithm, experiments were conducted to explore the impact of the number of selected experts. As presented in Table 5, when the number of selected experts is set to 4, there is an improvement in performance for dialects that are not part of the training dataset (Henan and Hubei dialects). This suggests that increasing the number of selected experts can enhance the model’s generalization to external data. The model’s performance remains similar when the number of selected experts is 2 or 1. However, it is worth noting that the decoding time for a single speech increases by approximately 53% when the number of selected experts is 4 compared to when it is 1. This indicates that the number of selected experts has a minimal impact on the model’s performance but significantly affects decoding efficiency, which is crucial for a robust speech recognition system.

4.5 Layer-Wise Analysis of Experts

In Figure 4, we present a visualization of the expert weights applied to the test sets corresponding to Sichuan and Kunming dialects. We can observe certain patterns in the weights. Across the initial three layers of the model, both dialects manifest a heightened degree of distinctiveness in expert selection, indicative of specific groups of experts concentrating exclusively on dialect-specific information. Within the intermediate layers of the model, expert weights display a diminished prominence, yet discernible differences persist in the expert weights associated with the two dialects. This observation suggests that varying combinations of experts implicitly encapsulate distinctive information pertaining to dialectal variations. In the concluding three layers of the model, the deployed experts exhibit near-identical characteristics, thereby indirectly affirming the model’s proficiency in capturing shared

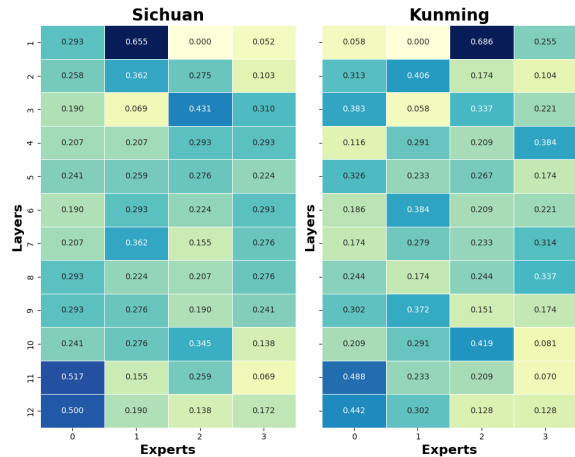


Figure 4: The expert weights are visualized on Sichuan dialect and Kunming dialect.

features between Sichuan and Kunming dialects.

5 Conclusion

In this manuscript, we present a multi-dialectal speech recognition model based on MoE termed DialectMoE. Structurally, it incorporates a dual-encoder architecture, wherein the general encoder is dedicated to acquiring general acoustic representations, and the dialect encoder is specialized for acquiring acoustic representations across various dialects. A refinement in the dynamic routing strategy within the MoE layer of the universal encoder has been introduced to enable the selection of appropriate experts based on the acoustic information specific to the dialect in the input sequence. Furthermore, we propose a three-stage training methodology to facilitate DialectMoE in learning distinct tasks at different phases, thereby enhancing its adaptability and performance across varying aspects of the multi-dialectal speech recognition task. Experimental results demonstrate that the proposed DialectMoE model achieves remarkable performance in multi-dialects speech recognition tasks.

6 Limitations

While the MoE-based approach can effectively enhance model performance, it inherently results in an increase in the number of model parameters. This increase in parameters can lead to higher training costs and longer inference times, which are inevitable consequences. Therefore, it is imperative to conduct further research on model compression techniques to mitigate these issues.

615
616
617
618
619
620
621

622
623
624
625

626
627
628
629
630
631
632

633
634
635
636

637
638
639
640
641
642

643
644
645
646
647

648
649
650
651
652
653

654
655
656
657
658

659
660
661
662
663

664
665
666
667
668
669

References

Sadeen Alharbi, Muna Alrazgan, Alanoud Alrashed, Turkiyah Alnomasi, Raghad Almojel, Rimah Alharbi, Saja Alharbi, Sahar Alturki, Fatimah Alshehri, and Maha Almojil. 2021. Automatic speech recognition: Systematic literature review. *IEEE Access*, 9:131858–131876.

Eiman Alsharhan and Allan Ramsay. 2020. Investigating the effects of gender, dialect, and training size on the performance of arabic speech recognition. *Language Resources and Evaluation*, 54:975–998.

Hui Bu, Jiayu Du, Xingyu Na, Bengu Wu, and Hao Zheng. 2017. Aishell-1: An open-source mandarin speech corpus and a speech recognition baseline. In *2017 20th conference of the oriental chapter of the international coordinating committee on speech databases and speech I/O systems and assessment (O-COCOSDA)*, pages 1–5. IEEE.

Zhengjia Dan, Yue Zhao, Xiaojun Bi, Licheng Wu, and Qiang Ji. 2022. Multi-task transformer with adaptive cross-entropy loss for multi-dialect speech recognition. *Entropy*, 24(10):1429.

Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. 2022. Glam: Efficient scaling of language models with mixture-of-experts. In *International Conference on Machine Learning*, pages 5547–5569. PMLR.

Mohamed Elfeky, Meysam Bastani, Xavier Velez, Pedro Moreno, and Austin Waters. 2016. Towards acoustic model unification across dialects. In *2016 IEEE Spoken Language Technology Workshop (SLT)*, pages 624–628. IEEE.

Zhiwen Fan, Rishov Sarkar, Ziyu Jiang, Tianlong Chen, Kai Zou, Yu Cheng, Cong Hao, Zhangyang Wang, et al. 2022. M³vit: Mixture-of-experts vision transformer for efficient multi-task learning with model-accelerator co-design. *Advances in Neural Information Processing Systems*, 35:28441–28457.

William Fedus, Barret Zoph, and Noam Shazeer. 2022. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *The Journal of Machine Learning Research*, 23(1):5232–5270.

Akhilesh Gotmare, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. 2018. A closer look at deep learning heuristics: Learning rate restarts, warmup and distillation. *arXiv preprint arXiv:1810.13243*.

Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al. 2020. Conformer: Convolution-augmented transformer for speech recognition. *arXiv preprint arXiv:2005.08100*.

Arthur Hinsvark, Natalie Delworth, Miguel Del Rio, Quinten McNamara, Joshua Dong, Ryan Westerman, Michelle Huang, Joseph Palakapilly, Jennifer Drexler, Ilya Pirkin, et al. 2021. Accented speech recognition: A survey. *arXiv preprint arXiv:2104.10747*.

Dah-an Ho. 2015. Chinese dialects. *The Oxford handbook of Chinese linguistics*, pages 149–159.

Takaaki Hori, Shinji Watanabe, and John R Hershey. 2017. Joint ctc/attention decoding for end-to-end speech recognition. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 518–529.

Jason J Humphries, Philip C Woodland, and D Pearce. 1996. Using accent-specific pronunciation modelling for robust speech recognition. In *Proceeding of Fourth International Conference on Spoken Language Processing. ICSLP’96*, volume 4, pages 2324–2327. IEEE.

Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. 1991. Adaptive mixtures of local experts. *Neural computation*, 3(1):79–87.

Rui Jiang. 2023. [Chinese dialect recognition based on transfer learning](#). *INTERSPEECH 2023*.

Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Yoohwan Kwon and Soo-Whan Chung. 2023. Mole: Mixture of language experts for multi-lingual automatic speech recognition. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE.

Bo Li, Tara N Sainath, Khe Chai Sim, Michiel Bacchiani, Eugene Weinstein, Patrick Nguyen, Zhifeng Chen, Yanghui Wu, and Kanishka Rao. 2018. Multi-dialect speech recognition with a single sequence-to-sequence model. In *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 4749–4753. IEEE.

Sheng Li, Xugang Lu, Chenchen Ding, Peng Shen, Tatsuya Kawahara, and Hisashi Kawai. 2019. Investigating radical-based end-to-end speech recognition systems for chinese dialects and japanese. In *INTERSPEECH*, pages 2200–2204.

Mishaim Malik, Muhammad Kamran Malik, Khawar Mehmood, and Imran Makhdoom. 2021. Automatic speech recognition: a survey. *Multimedia Tools and Applications*, 80:9411–9457.

Daniel S Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D Cubuk, and Quoc V Le. 2019. Specaugment: A simple data augmentation method for automatic speech recognition. *arXiv preprint arXiv:1904.08779*.

| | | |
|-----|---|-----|
| 722 | Yifan Peng, Siddharth Dalmia, Ian Lane, and Shinji Watanabe. 2022. Branchformer: Parallel mlp-attention architectures to capture local and global context for speech recognition and understanding. In <i>International Conference on Machine Learning</i> , pages 17627–17643. PMLR. | |
| 723 | | |
| 724 | | |
| 725 | | |
| 726 | | |
| 727 | | |
| 728 | Zongze Ren, Guofu Yang, and Shugong Xu. 2019. Two-stage training for chinese dialect recognition. <i>arXiv preprint arXiv:1908.02284</i> . | |
| 729 | | |
| 730 | | |
| 731 | Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André Susano Pinto, Daniel Keysers, and Neil Houlsby. 2021. Scaling vision with sparse mixture of experts. <i>Advances in Neural Information Processing Systems</i> , 34:8583–8595. | |
| 732 | | |
| 733 | | |
| 734 | | |
| 735 | | |
| 736 | | |
| 737 | Kim Sehoon, Gholami Amir, Shaw Albert, Lee Nicholas, Mangalam Karttikeya, Keutzer Kurt, et al. 2023. Squeezeformer: An efficient transformer for automatic speech recognition. | |
| 738 | | |
| 739 | | |
| 740 | | |
| 741 | Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. <i>arXiv preprint arXiv:1701.06538</i> . | |
| 742 | | |
| 743 | | |
| 744 | | |
| 745 | | |
| 746 | Richard Sproat, Liang Gu, Jing Li, Yanli Zheng, Yi Su, Haolang Zhou, Philip Bramsen, MIT David Kirsch, Izhak Shafran, Stavros Tsakalidis, et al. 2004. Dialectal chinese speech recognition. In <i>CLSP Summer Workshop</i> . | |
| 747 | | |
| 748 | | |
| 749 | | |
| 750 | | |
| 751 | Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. 2016. Rethinking the inception architecture for computer vision. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 2818–2826. | |
| 752 | | |
| 753 | | |
| 754 | | |
| 755 | | |
| 756 | Dong Wang, Xiaodong Wang, and Shaohe Lv. 2019. An overview of end-to-end automatic speech recognition. <i>Symmetry</i> , 11(8):1018. | |
| 757 | | |
| 758 | | |
| 759 | Xuefei Wang, Yanhua Long, Yijie Li, and Haoran Wei. 2023. Multi-pass training and cross-information fusion for low-resource end-to-end accented speech recognition. <i>arXiv preprint arXiv:2306.11309</i> . | |
| 760 | | |
| 761 | | |
| 762 | | |
| 763 | Zhao You, Shulin Feng, Dan Su, and Dong Yu. 2021. Speechmoe: Scaling to large acoustic models with dynamic routing mixture of experts. <i>arXiv preprint arXiv:2105.03036</i> . | |
| 764 | | |
| 765 | | |
| 766 | | |
| 767 | Zhao You, Shulin Feng, Dan Su, and Dong Yu. 2022. Speechmoe2: Mixture-of-experts model with improved routing. In <i>ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 7217–7221. IEEE. | |
| 768 | | |
| 769 | | |
| 770 | | |
| 771 | | |
| 772 | Binbin Zhang, Di Wu, Zhendong Peng, Xingchen Song, Zhuoyuan Yao, Hang Lv, Lei Xie, Chao Yang, Fuping Pan, and Jianwei Niu. 2022a. Wenet 2.0: More productive end-to-end speech recognition toolkit. <i>arXiv preprint arXiv:2203.15455</i> . | |
| 773 | | |
| 774 | | |
| 775 | | |
| 776 | | |
| | Fengrun Zhang, Xiang Xie, and Xinyue Quan. 2022b. Chinese dialect speech recognition based on end-to-end machine learning. In <i>2022 International Conference on Machine Learning, Control, and Robotics (MLCR)</i> , pages 14–18. IEEE. | 777 |
| | | 778 |
| | | 779 |
| | | 780 |
| | | 781 |
| | Vicky Zilvan, Ana Heryana, Asri Rizki Yuliani, Dikdik Krisnandi, R Sandra Yuwana, and Hilman F Pardede. 2021. Front-end based robust speech recognition methods: A review. In <i>Proceedings of the 2021 International Conference on Computer, Control, Informatics and Its Applications</i> , pages 136–140. | 782 |
| | | 783 |
| | | 784 |
| | | 785 |
| | | 786 |
| | | 787 |