

A Zhuang Speech-to-Text Translation Method with Source Language-Aware Conditional Attention Mechanism

Abstract

End-to-end speech-to-text translation aims to directly convert spoken input into text in the target language. As a minority language in China, Zhuang faces challenges such as data scarcity and limited technological support in the field of speech translation. To support research on Zhuang speech translation, we construct a recording platform that is used to compile speech-text parallel corpus. We introduce a Source Language-Aware Conditional Attention Mechanism (LACA), which incorporates source language information into the Transformer to bias attention toward Zhuang-specific linguistic features. Additionally, an Implicit Connectionist Temporal Classification Auxiliary Mechanism (ICAM) is employed during training to provide auxiliary supervision for alignment learning between speech and text representations. Experimental results demonstrate that our model achieves a BLEU score of 32.46 on Zhuang speech-to-text translation, outperforming the baseline by 4.12 BLEU points.

1 Introduction

In recent years, end-to-end models (Ren et al., 2020; Prabhavalkar et al., 2023; Chen et al., 2024; Zhang et al., 2024a) have become mainstream in speech translation, driving the development of large systems like GPT-4o (OpenAI, 2024), Gemini (Google, 2024), and SeamlessM4T (Communication, 2023).

Previous research on Zhuang text processing has investigated a range of technical approaches. Xian Wu (2024) employed word embeddings, relative position encoding, and a Transformer architecture for Chinese-Zhuang translation. Zhang et al. (2024b) proposed Dipmy++,

incorporating a bilingual dictionary and 5,000 parallel sentences. Weiquan Zhang (2022) developed an F-BiLSTM-CNN-CRF model for named entity recognition, while Kui Ning (2019) designed a system for translating and proofreading Zhuang folk culture texts. Compared to text-based research, research on Zhuang speech remains relatively scarce. Huang et al. (2025) incorporated the SimAM attention mechanism into a Bidirectional Encoder Representations from Transformers model to perform speech classification across different Zhuang dialects. Jie Wang (2024) developed a Zhuang speech synthesis system supporting text normalization, phonetic conversion, and waveform generation.

However, research on end-to-end Zhuang speech-to-text translation (S2TT) remains limited. Most current end-to-end S2TT models are developed based on the Fairseq framework (Wang et al., 2020), such as s2t_transformer and s2t_conformer. Despite their success in high-resource settings, these models struggle to align speech and text representations under low-resource conditions due to limited supervision and alignment optimization. Moreover, the differences in speech modality between Zhuang and the languages used during pretraining limit their ability to generalize to unseen languages.

To address the above challenges, we designed a Source Language-Aware Conditional Attention mechanism (LACA) and introduced Connectionist Temporal Classification (CTC) as an implicit auxiliary supervision signal, referred to as the Implicit CTC Auxiliary Mechanism (ICAM), thereby enabling end-to-end S2TT for Zhuang. First, we introduced CTC loss during the training phase to assist the pre-trained speech encoder model wav2vec_small (Baevski et al., 2020) in learning feature alignment. Then, a one-dimensional convolutional adaptor (referred to as

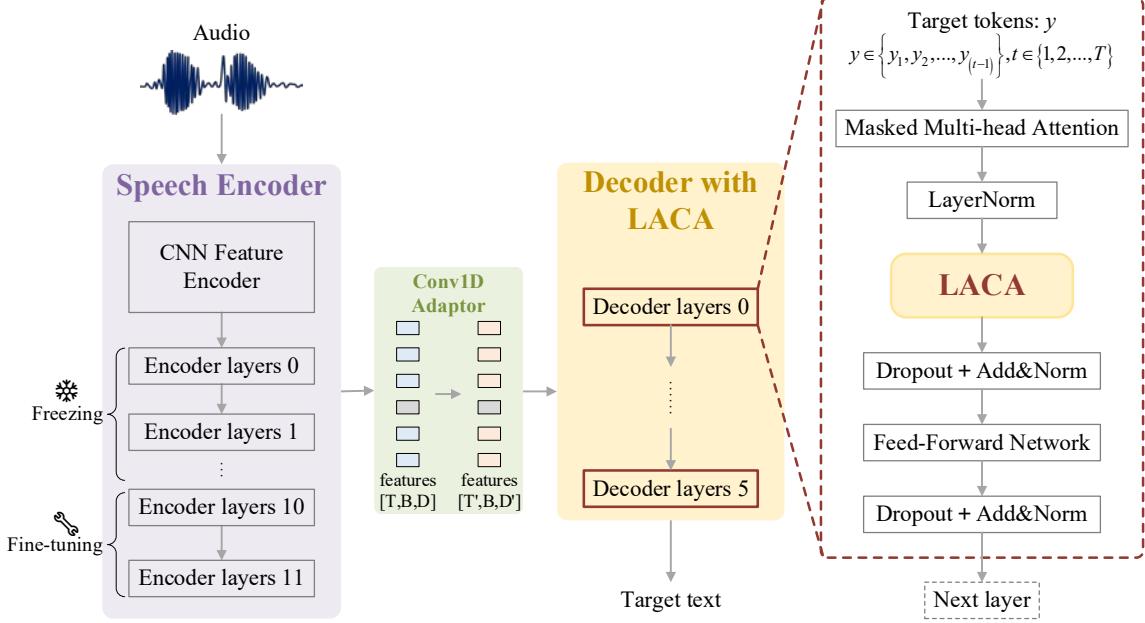


Figure 1: The overall architecture of our model. Raw speech is first encoded with the assistance of ICAM to enhance feature alignment. The resulting representations are then passed to the decoder, where LACA dynamically adjusts attention weights based on source language features.

Conv1D adaptor) was employed to transform the encoder output, serving as an interface to ensure compatibility with the decoder. Finally, LACA integrates language embeddings to enable the decoder to dynamically adjust attention weights based on the source language features during cross-attention. In addition, we developed a Zhuang translation and recording platform, and released several demos showcasing Zhuang S2TT¹.

2 Methods

To realize Zhuang S2TT, we designed LACA based on the xm_transformer² model and incorporated the ICAM. This section presents the overall framework as illustrated in Figure 1. The model consists of a speech encoder, a Conv1D adaptor, and a Transformer decoder.

2.1 Speech Encoder

Following prior approaches (Le et al., 2023; Xu et al., 2024; Yan et al., 2023), we incorporated CTC loss during training to facilitate alignment between speech and textual representations. The auxiliary CTC loss is computed as:

$$\mathcal{L}_{CTC} = CTC_{loss}(\hat{Y}, Y_{target}), \quad (1)$$

where $Encoder(X)$ denotes the output sequence of the speech encoder given the input audio X , and $\hat{Y} = Linear_{CTC}(Encoder(X))$ is the predicted probability sequence obtained by applying a linear projection to the encoder output. Y_{target} is the ground-truth sequence. Unlike conventional CTC-based models, this auxiliary CTC branch is excluded during inference and does not affect the main model architecture. To train the decoder, we adopt the standard cross-entropy loss between the predicted sequence Y_{pred} and the reference sequence Y_{target} :

$$\mathcal{L}_{Decoder} = CrossEntropy(Y_{pred}, Y_{target}). \quad (2)$$

The overall training objective combines the decoder loss and the auxiliary CTC loss, weighted by λ_{CTC} :

$$\mathcal{L}_{total} = \mathcal{L}_{Decoder} + \lambda_{CTC} \cdot \mathcal{L}_{CTC}. \quad (3)$$

To mitigate overfitting in low-resource settings, as suggested by Yang et al. (2022; 2023), we fine-tuned only the last two layers of the encoder while keeping the first ten layers frozen.

2.2 Conv1D Adaptor

The Conv1D Adaptor is a lightweight feature transformation module that applies one-dimensional convolution to the output of the speech encoder, adjusting the feature dimension to

¹ <http://510.english-gxu.net/>

² https://github.com/facebookresearch/fairseq/tree/main/fairseq/models/speech_to_text

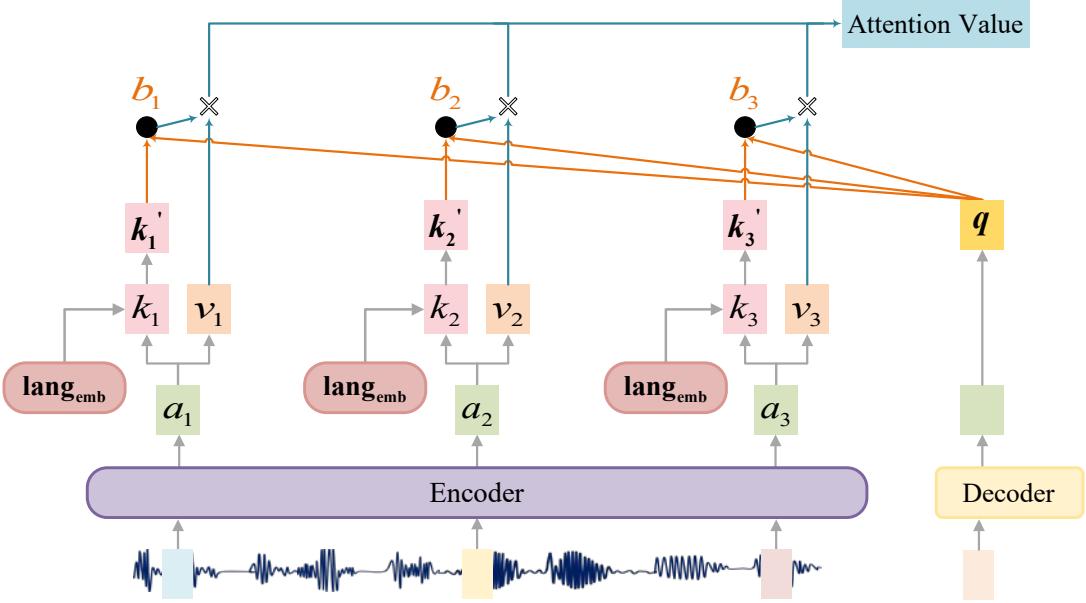


Figure 2: Source Language-Aware Conditional Attention Mechanism. lang_{emb} : The expanded source language embedding. ●: Inner product. ×: Weighted operation.

match the input requirements of the Transformer decoder. It also performs temporal downsampling¹⁶² via convolutional stride. To mitigate the high information density of speech features, a Gated¹⁶³ Linear Unit (GLU) is applied following the¹⁶⁴ convolution, facilitating the selective transmission¹⁶⁵ of the most relevant temporal features to the¹⁶⁶ decoder.¹⁶⁷

2.3 Decoder with LACA Mechanism

We proposed the LACA mechanism as illustrated¹⁷⁰ in Figure 2. Let the decoder-derived query matrix¹⁷¹ be defined as:¹⁷²

$$Q = [q_1, q_2, \dots, q_n] \in R^{n \times d_k}, \quad (4)$$

where n is the number of decoding steps and d_k is¹⁷⁴ the attention dimension. The encoder output is¹⁷⁵ used to construct the key and value matrices:¹⁷⁶

$$\begin{aligned} K &= [k_1, k_2, \dots, k_m] \in R^{m \times d_k}, \\ V &= [v_1, v_2, \dots, v_m] \in R^{m \times d_k}, \end{aligned} \quad (5)$$

where m is the number of encoder steps. To inject¹⁷⁸ language-specific information, we define a¹⁷⁹ language embedding vector $e_L \in R^{B \times d_k}$, where B ¹⁸⁰ is the batch size. This vector is expanded along¹⁸¹ the temporal dimension to match the encoder¹⁸² output:¹⁸³

$$\text{lang}_{\text{emb}} = \text{Expand}(e_L) \in R^{m \times B \times d_k}. \quad (6)$$

The expanded bias is then added element-wise to¹⁸⁴ the original key to obtain the modified key:¹⁸⁵

$$K' = K + \text{lang}_{\text{emb}}. \quad (7)$$

Finally, the attention is calculated using the standard scaled dot-product formulation:¹⁸⁶

$$\text{Attention}(Q, K', V) = \text{softmax}\left(\frac{QK'^T}{\sqrt{d_k}}\right)V. \quad (8)$$

Compared with existing language-aware methods (Xu et al., 2023; Tian et al., 2022), the proposed LACA differs significantly in both the position and manner of language information injection. By introducing the language embedding as a bias term into the key vectors of the cross-attention module, LACA stably modulates the attention distribution. This lightweight design enhances training robustness under low-resource conditions and enables seamless extension to multilingual speech translation tasks.

3 Experimental Setup

This section introduces the experimental setup for the Zhuang end-to-end S2TT task, including the dataset, baseline system, evaluation metrics, and other implementation details.

3.1 Dataset

This paper focuses on the Dejing dialect of Zhuang (Lv, 2019). We developed a Zhuang speech recording platform³ on a private server and collected speech from native speakers. All recordings were manually verified by linguistic experts. In addition, part of the dataset was sourced from news broadcasts on Jingxi TV, and

³ <http://510.english-gxu.net/record>

187 The Data Collection, Recording, and Display²³²
 188 Platform for the Chinese Language Resources²³³
 189 Protection Project⁴. During data preprocessing, all²³⁴
 190 audio recordings were standardized to a sampling²³⁵
 191 rate of 16 kHz and converted to mono (Huang et²³⁶
 192 al., 2025), followed by data augmentation. The²³⁷
 193 dataset consists of both isolated vocabulary items²³⁸
 194 and complete sentences, with a total duration of²³⁹
 195 19.23 hours. To evaluate the generalization
 196 capability of the model, we conducted validation
 197 on the Thai and English subsets of the Common
 198 Voice (Ardila et al., 2020).

199 3.2 Baseline systems and metrics

200 In this paper, s2t_transformer (Wang et al., 2020)²⁴⁰
 201 is used as the baseline system. It is a Transformer
 202 based sequence-to-sequence model for S2TT. In
 203 the experiment, accuracy on the validation set
 204 during the training and BLEU score during
 205 inference are selected as evaluation metrics.

206 3.3 Implementation details

207 The model was trained on a single NVIDIA²⁴¹
 208 GeForce RTX 4090 using the Adam optimizer
 209 with a learning rate of 2e-4. The weight of the²⁴³
 210 CTC loss was set to 0.3.²⁴⁴

211 4 Results and Discussion

212 With a limited Zhuang speech-text parallel corpus,²⁴⁸
 213 our model still achieved a BLEU score of 32.46²⁴⁹
 214 on the test set. In the following sections, we²⁵⁰
 215 compare our model with the baseline, and analyze²⁵¹
 216 the contribution of each component.²⁵²

217 4.1 Comparison with the Baseline System

218 We compared the performance of different models
 219 on the Zhuang dataset for S2TT task. As shown in
 220 Table 1, although the validation accuracy of our
 221 model was comparable to that of the baseline
 222 system, it achieved a BLEU score that was 4.12²⁵³
 223 points higher than the baseline.

224 To assess the generalization performance of the²⁵⁴
 225 model on S2TT tasks, we conducted validation
 226 using the Common Voice dataset. As indicated in²⁵⁵
 227 Table 2, the model achieved a validation accuracy²⁵⁶
 228 of 86.89% and a BLEU score of 24.71 on the Thai²⁵⁷
 229 dataset, demonstrating its cross-linguistic²⁵⁸
 230 generalization capability. On the English dataset,²⁵⁹
 231 since English and Zhuang belong to different²⁶⁰
 232 language families and exhibit differences in
 233 speech feature patterns, our approach was
 234 sensitive to these differences, resulting in a lower
 235 validation accuracy. Nevertheless, the BLEU
 236 score of 25.05 suggests that the model-generated
 237 sentences exhibited local word order similarity to
 238 the reference translations and preserved a coherent
 239 syntactic structure at the sentence level.

language families and exhibit differences in speech feature patterns, our approach was sensitive to these differences, resulting in a lower validation accuracy. Nevertheless, the BLEU score of 25.05 suggests that the model-generated sentences exhibited local word order similarity to the reference translations and preserved a coherent syntactic structure at the sentence level.

| Model | accuracy(dev) | BLEU |
|-----------------|---------------|-------|
| s2t_transformer | 96.73% | 28.34 |
| s2t_conformer | 85.20% | 14.70 |
| Ours | 96.35% | 32.46 |

Table 1: S2TT Performance of different models on the Zhuang dataset.

| Dataset | accuracy(dev) | BLEU |
|---------|---------------|-------|
| Thai | 86.89% | 24.71 |
| English | 31.17% | 25.05 |
| Zhuang | 96.35% | 32.46 |

Table 2: S2TT Performance of our model on Zhuang, Thai, and English datasets.

4.2 Ablation study

To further evaluate the effectiveness of ICAM and LACA, we conducted ablation experiments, with the results summarized in Table 3. The integration of ICAM led to improvements in both validation accuracy and BLEU score. Upon incorporating LACA, despite a 0.42% drop in validation accuracy, the model attained a BLEU score of 32.46, marking a 2.42-point improvement and reflecting a favorable trade-off between accuracy and translation quality.

| Model | accuracy(dev) | BLEU |
|---------------------------|---------------|-------|
| xm_transformer | 96.59% | 30.04 |
| xm_transformer+ICAM | 96.77% | 30.42 |
| xm_transformer+ICAM +LACA | 96.35% | 32.46 |

Table3: The results of the ablation study.

5 Conclusion

In this study, we developed a Zhuang translation and recording platform and introduced the LACA mechanism and the ICAM training strategy to enhance Zhuang speech-to-text translation (S2TT) task. Experimental results show that our method achieves accurate and effective performance in Zhuang S2TT. This work contributes to both low-resource speech translation research and the preservation of minority language heritage.

⁴ <https://zhongguoyuyan.cn/index>

- 264 **Limitations** 318
 265 Although our method yields a satisfactory high 319
 266 BLEU score on the Zhuang dataset, performance 320
 267 on other languages remains moderate, indicating 321
 268 that the generalization capability still needs to be 322
 269 improved. In future work, we aim to improve 324
 270 performance in speech-to-text translation tasks 325
 271 across a wider range of languages. Furthermore, 326
 272 we plan to explore speech-to-speech translation 327
 273 both among Zhuang dialects and between Zhuang 329
 274 and other languages as an extension of this 330
 275 research. To facilitate these advancements, we 331
 276 will collect and construct a larger-scale Zhuang 332
 277 speech-text parallel corpus to mitigate the 333
 278 limitations posed by data scarcity. 334
 279 **References** 335
 280 Rosana Ardila, Megan Branson, Kelly Davis, Michael 337
 281 Kohler, Josh Meyer, Michael Henretty, Reuben 338
 282 Morais, Lindsay Saunders, Francis Tyers, and 339
 283 Gregor Weber. 2020. *Common Voice: A Massively- 340
 284 Multilingual Speech Corpus*. In *Proceedings of the 341
 285 Twelfth Language Resources and Evaluation 342
 286 Conference*, pages 4218–4222, Marseille, France. 343
 287 European Language Resources Association. 344
 288 <https://aclanthology.org/2020.lrec-1.520/>. 345
 289 Alexei Baevski, Yuhao Zhou, Abdelrahman 347
 290 Mohamed, Michael Auli. 2020. wav2vec 2.0: 348
 291 a framework for self-supervised learning of speech 349
 292 representations. In *Proceedings of the 34th 350
 293 International Conference on Neural Information 351
 294 Processing Systems*, pages 12449–12460, 352
 295 Vancouver, Canada. Curran Associates, Inc. 353
 296 <https://dl.acm.org/doi/abs/10.5555/3495724.34967> 354
 297 68. 355
 298 Xi Chen, Songyang Zhang, Qibing Bai, Kai Chen, 356
 299 and Satoshi Nakamura. 2024. *LLaST: Improved 357
 300 End-to-end Speech Translation System Leveraged 358
 301 by Large Language Models*. In *Findings of the 359
 302 Association for Computational Linguistics: ACL 360
 303 2024*, pages 6976–6987, Bangkok, Thailand. 361
 304 Association for Computational Linguistics. doi: 362
 305 [10.18653/v1/2024.findings-acl.416](https://doi.org/10.18653/v1/2024.findings-acl.416). 363
 306 Seamless Communication, Loïc Barrault, Yu-An 365
 307 Chung, Mariano Cora Meglioli, David Dale, Ning 366
 308 Dong, Paul-Ambroise Duquenne, Hady 367
 309 Elsahar, Hongyu Gong, Kevin Heffernan, John 368
 310 Hoffman, Christopher Klaiber, Pengwei Li, Daniel 369
 311 Licht, Jean Maillard, Alice Rakotoarison, Kaushik 370
 312 Ram Sadagopan, Guillaume Wenzek, Ethan 371
 313 Ye, Bapi Akula, Peng-Jen Chen, Naji El 372
 314 Hachem, Brian Ellis, Gabriel Mejia 373
 315 Gonzalez, Justin Haaheim, Prangthip 374
 316 Hansanti, Russ Howes, Bernie Huang, Min-Jae 375
 317 Hwang, Hirofumi Inaguma, Somya Jain, Elahe
 318 Kalbassi, Amanda Kallet, Ilia Kulikov, Janice 319
 319 Lam, Daniel Li, Xutai Ma, Ruslan 320
 320 Mavlyutov, Benjamin Peloquin, Mohamed 321
 321 Ramadan, Abinesh Ramakrishnan, Anna 322
 322 Sun, Kevin Tran, Tuan Tran, Igor Tufanov, Vish 323
 323 Vogeti, Carleigh Wood, Yilin Yang, Bokai 324
 324 Yu, Pierre Andrews, Can Balioglu, Marta R. Costa- 325
 325 jussà, Onur Celebi, Maha Elbayad, Cynthia 326
 326 Gao, Francisco Guzmán, Justine Kao, Ann 327
 327 Lee, Alexandre Mourachko, Juan Pino, Sravya 328
 328 Popuri, Christophe Ropers, Safiyyah 329
 329 Saleem, Holger Schwenk, Paden 330
 330 Tomasello, Changhan Wang, Jeff Wang, and Skyler 331
 331 Wang. 2023. *SeamlessM4T: Massively 332
 332 Multilingual & Multimodal Machine Translation*. 333
 333 *arXiv preprint arXiv:2308.11596*.
 334 Gemini Team Google: Petko Georgiev, Ving Ian 335
 335 Lei, Ryan Burnell, Libin Bai, Anmol 336
 336 Gulati, Garrett Tanzer, Damien Vincent, Zhufeng 337
 337 Pan, Shibo Wang, Soroosh Mariooryad, Yifan 338
 338 Ding, Xinyang Geng, Fred Alcober, Roy 339
 339 Frostig, Mark Omernick, Lexi Walker, Cosmin 340
 340 Paduraru, Christina Sorokin, Andrea 341
 341 Tacchetti, Colin Gaffney, Samira Daruki, Olcan 342
 342 Sercinoglu, Zach Gleicher, Juliette Love, Paul 343
 343 Voigtlaender, Rohan Jain, Gabriela Surita, Kareem 344
 344 Mohamed, Rory Blevins, Junwhan Ahn, Tao 345
 345 Zhu, Kornraphop Kawintiranon, Orhan 346
 346 Firat, Yiming Gu, Yujing Zhang, Matthew 347
 347 Rahtz, Manaal Faruqui, Natalie Clay, Justin 348
 348 Gilmer, JD Co-Reyes, Ivo Penchev, Rui 349
 349 Zhu, Nobuyuki Morioka, Kevin Hui, Krishna 350
 350 Haridasan, Victor Campos, Mahdis 351
 351 Mahdieh, Mandy Guo, Samer Hassan, Kevin 352
 352 Kilgour, Arpi Vezer, Heng-Tze Cheng, Raoul de 353
 353 Liedekerke, Siddharth Goyal, Paul Barham, DJ 354
 354 Strouse, Seb Noury, Jonas Adler, Mukund 355
 355 Sundararajan, Sharad Vikram, Dmitry 356
 356 Lepikhin, Michela Paganini, Xavier Garcia, Fan 357
 357 Yang, Dasha Valter, Maja Trebacz, Kiran 358
 358 Vodrahalli, Chulayuth Asawaroengchai, Roman 359
 359 Ring, Norbert Kalb, Livio Baldini 360
 360 Soares, Siddhartha Brahma, David Steiner, Tianhe 361
 361 Yu, Fabian Mentzer, Antoine He, Lucas 362
 362 Gonzalez, Bibo Xu, Raphael Lopez 363
 363 Kaufman, Laurent El Shafey, Junhyuk Oh, Tom 364
 364 Hennigan, George van den Driessche, Seth 365
 365 Odoom, Mario Lucic, Becca Roelofs, Sid 366
 366 Lall, Amit Marathe, Betty Chan, Santiago 367
 367 Ontanon, Luheng He, Denis Teplyashin, Jonathan 368
 368 Lai, Phil Crone, Bogdan Damoc, Lewis 369
 369 Ho, Sebastian Riedel, Karel Lenc, Chih-Kuan 370
 370 Yeh et al. (1037 additional authors not shown). 371
 371 2024. *Gemini 1.5: Unlocking multimodal 372
 372 understanding across millions of tokens of context*. 373
 373 *arXiv preprint arXiv:2403.05530*.
 374 Min Huang, Xuejun Zhang and Wenkang Chen. 2025. 375
 375 Classification of Zhuang Dialect combined with

- 376 Bert and SimAM. In *ICASSP 2025 - 2025 IEEE*⁴³³
 377 *International Conference on Acoustics, Speech and*⁴³⁴
 378 *Signal Processing*, pages 1–5, Hyderabad, India⁴³⁵
 379 IEEE. doi: [10.1109/ICASSP49660.2025.10890693](https://doi.org/10.1109/ICASSP49660.2025.10890693).⁴³⁶
- 380 Phuong-Hang Le, Hongyu Gong, Changhan Wang⁴³⁷
 381 Juan Pino, Benjamin Lecouteux, Didier Schwab⁴³⁸
 382 2023. Pre-training for speech translation: CTC⁴³⁹
 383 meets optimal transport. In *Proceedings of the 40th*⁴⁴⁰
 384 *International Conference on Machine Learning*⁴⁴¹
 385 pages 18667–18685, Honolulu, Hawaii USA⁴⁴²
 386 Proceedings of Machine Learning Research.⁴⁴³
 387 <https://proceedings.mlr.press/v202/le23a.html>.⁴⁴⁴
- 388 Songsong Lv. 2019. Evolution of Proximal⁴⁴⁵
 389 Demonstratives in Dejing Dialect of Zhuang⁴⁴⁶
 390 Language——A Perspective Based on Linguistic⁴⁴⁷
 391 Contact. *Journal of Guangxi Normal University*⁴⁴⁸
 392 (*Philosophy and Social Sciences Edition*)⁴⁴⁹
 393 55(04):108–118. doi: [10.16088/j.issn.1001-450](https://doi.org/10.16088/j.issn.1001-450)
 394 6597.2019.04.014.⁴⁵¹
- 395 Kui Ning. 2019. Zhuang Intelligent Translation⁴⁵²
 396 System 1.0 Based on a Parallel Corpus. Master's⁴⁵³
 397 thesis, Guangxi University, Nanning, China, April.⁴⁵⁴
- 398 OpenAI: Aaron Hurst, Adam Lerer, Adam P.⁴⁵⁵
 399 Goucher, Adam Perelman, Aditya Ramesh, Aidan⁴⁵⁶
 400 Clark, AJ Ostrow, Akila Welihinda, Alan⁴⁵⁷
 401 Hayes, Alec Radford, Aleksander Mądry, Alex⁴⁵⁸
 402 Baker-Whitcomb, Alex Beutel, Alex⁴⁵⁹
 403 Borzunov, Alex Carney, Alex Chow, Alex⁴⁶⁰
 404 Kirillov, Alex Nichol, Alex Paino, Alex⁴⁶¹
 405 Renzin, Alex Tachard Passos, Alexander⁴⁶²
 406 Kirillov, Alexi Christakis, Alexis Conneau, Ali⁴⁶³
 407 Kamali, Allan Jabri, Allison Moyer, Allison⁴⁶⁴
 408 Tam, Amadou Crookes, Amin Tootoochian, Amin⁴⁶⁵
 409 Tootoonchian, Ananya Kumar, Andrea⁴⁶⁶
 410 Vallone, Andrej Karpathy, Andrew⁴⁶⁷
 411 Braunstein, Andrew Cann, Andrew⁴⁶⁸
 412 Codispoti, Andrew Galu, Andrew⁴⁶⁹
 413 Kondrich, Andrew Tulloch, Andrey⁴⁷⁰
 414 Mishchenko, Angela Baek, Angela Jiang, Antoine⁴⁷¹
 415 Pelisse, Antonia Woodford, Anuj Gosalia, Arka⁴⁷²
 416 Dhar, Ashley Pantuliano, Avi Nayak, Avital⁴⁷³
 417 Oliver, Barret Zoph, Behrooz Ghorbani, Ben⁴⁷⁴
 418 Leimberger, Ben Rossen, Ben Sokolowsky, Ben⁴⁷⁵
 419 Wang, Benjamin Zweig, Beth Hoover, Blake⁴⁷⁶
 420 Samic, Bob McGrew, Bobby Spero, Bogo⁴⁷⁷
 421 Giertler, Bowen Cheng, Brad Lightcap, Brandon⁴⁷⁸
 422 Walkin, Brendan Quinn, Brian Guaraci, Brian⁴⁷⁹
 423 Hsu, Bright Kellogg, Brydon Eastman, Camillo⁴⁸⁰
 424 Lugaresi, Carroll Wainwright, Cary Bassin, Cary⁴⁸¹
 425 Hudson, Casey Chu, Chad Nelson, Chak Li, Chan⁴⁸²
 426 Jun Shern, Channing Conger, Charlotte⁴⁸³
 427 Barette, Chelsea Voss, Chen Ding, Cheng⁴⁸⁴
 428 Lu, Chong Zhang, Chris Beaumont, Chris⁴⁸⁵
 429 Hallacy, Chris Koch, Christian Gibson, Christina⁴⁸⁶
 430 Kim, Christine Choi, Christine⁴⁸⁷
 431 McLeavey, Christopher Hesse, Claudia⁴⁸⁸
 432 Fischer, Clemens Winter, Coley Czarnecki, Colin⁴⁸⁹
- 433 Jarvis, Colin Wei, Constantin Koumouzelis, Dane⁴⁹⁰
 434 Sherburn et al. (318 additional authors not shown).⁴⁹¹
 435 2024. Gpt-4o system card. *arXiv preprint*
 436 *arXiv:2410.21276*.⁴⁹²
- 437 Rohit Prabhavalkar, Takaaki Hori, Tara N. Sainath,⁴⁹³
 438 Ralf Schlüter and Shinji Watanabe. 2023. End-to-⁴⁹⁴
 439 End Speech Recognition: A Survey. *IEEE/ACM*⁴⁹⁵
 440 *Transactions on Audio, Speech, and Language*⁴⁹⁶
 441 *Processing*, 32:325–351. doi:⁴⁹⁷
 442 [10.1109/TASLP.2023.3328283](https://doi.org/10.1109/TASLP.2023.3328283).⁴⁹⁸
- 443 Yi Ren, Jinglin Liu, Xu Tan, Chen Zhang, Tao Qin,⁴⁹⁹
 444 Zhou Zhao, and Tie-Yan Liu. 2020. *SimulSpeech:*⁵⁰⁰
 445 *End-to-End Simultaneous Speech to Text*⁵⁰¹
 446 *Translation*. In *Proceedings of the 58th Annual*⁵⁰²
 447 *Meeting of the Association for Computational*⁵⁰³
 448 *Linguistics*, pages 3787–3796, Online. Association⁵⁰⁴
 449 for Computational Linguistics. doi:⁵⁰⁵
 450 [10.18653/v1/2020.acl-main.350](https://doi.org/10.18653/v1/2020.acl-main.350).⁵⁰⁶
- 451 Jinchuan Tian, Jianwei Yu, Chunlei Zhang, Yuxian⁵⁰⁷
 452 Zou, Dong Yu. 2022. LAE: Language-Aware⁵⁰⁸
 453 Encoder for Monolingual and Multilingual ASR. In⁵⁰⁹
 454 *23th Annual Conference of the International*⁵¹⁰
 455 *Speech Communication Association*, pages 3178–⁵¹¹
 456 3182, Incheon, Korea. International Speech⁵¹²
 457 Communication Association. doi:⁵¹³
 458 [10.21437/Interspeech.2022-923](https://doi.org/10.21437/Interspeech.2022-923).⁵¹⁴
- 459 Xian Wu. 2024. Research on Chinese–Zhuang⁵¹⁵
 460 Machine Translation Based on Neural Networks.⁵¹⁶
 461 Master's thesis, Guangxi University for⁵¹⁷
 462 Nationalities, Nanning, China, March. doi:⁵¹⁸
 463 [10.27035/d.cnki.ggxmc.2024.000933](https://doi.org/10.27035/d.cnki.ggxmc.2024.000933).⁵¹⁹
- 464 Jie Wang. 2024. Research on Zhuang Language⁵²⁰
 465 Speech Synthesis Technology Based on End-to-⁵²¹
 466 End. Master's thesis, Guangxi University for⁵²²
 467 Nationalities, Nanning, China, March. doi:⁵²³
 468 [10.27035/d.cnki.ggxmc.2024.000424](https://doi.org/10.27035/d.cnki.ggxmc.2024.000424).⁵²⁴
- 469 Changhan Wang, Yun Tang, Xutai Ma, Anne Wu,⁵²⁵
 470 Dmytro Okhonko, and Juan Pino. 2020. *Fairseq*⁵²⁶
 471 *S2T: Fast Speech-to-Text Modeling with Fairseq*.⁵²⁷
 472 In *Proceedings of the 1st Conference of the Asia-⁵²⁸*
 473 *Pacific Chapter of the Association for*⁵²⁹
 474 *Computational Linguistics and the 10th*⁵³⁰
 475 *International Joint Conference on Natural*⁵³¹
 476 *Language Processing: System Demonstrations*,⁵³²
 477 pages 33–39, Suzhou, China. Association for⁵³³
 478 Computational Linguistics. doi:⁵³⁴
 479 [10.18653/v1/2020.aacl-demo.6](https://doi.org/10.18653/v1/2020.aacl-demo.6).⁵³⁵
- 480 Chen Xu, Xiaoqian Liu, Erfeng He, Yuhao⁵³⁶
 481 Zhang, Qianqian Dong, Tong Xiao. 2024. Bridging⁵³⁷
 482 the Gaps of Both Modality and Language:⁵³⁸
 483 Synchronous Bilingual CTC for Speech⁵³⁹
 484 Translation and Speech Recognition. In *ICASSP*⁵⁴⁰
 485 *2024-2024 IEEE International Conference on*⁵⁴¹
 486 *Acoustics, Speech and Signal Processing*, pages⁵⁴²
 487 12176–12180, Seoul, Korea. IEEE. doi:⁵⁴³
 488 [10.1109/ICASSP48485.2024.10445849](https://doi.org/10.1109/ICASSP48485.2024.10445849).⁵⁴⁴

- 489 Haoran Xu, Jean Maillard, and Vedanuj Goswami.
490 2023. *Language-Aware Multilingual Machine*
491 *Translation with Self-Supervised Learning*.
492 In *Findings of the Association for Computational*
493 *Linguistics: EACL 2023*, pages 526–539,
494 Dubrovnik, Croatia. Association for Computational
495 Linguistics. doi: [10.18653/v1/2023.findings-eacl.38](https://doi.org/10.18653/v1/2023.findings-eacl.38).
- 496
- 497 Brian Yan, Siddharth Dalmia, Yosuke Higuchi,
498 Graham Neubig, Florian Metze, Alan W Black, and
499 Shinji Watanabe. 2023. *CTC Alignments Improve*
500 *Autoregressive Translation*. In *Proceedings of the*
501 *17th Conference of the European Chapter of the*
502 *Association for Computational Linguistics*, pages
503 1623–1639, Dubrovnik, Croatia. Association for
504 Computational Linguistics. doi:
505 [10.18653/v1/2023.eacl-main.119](https://doi.org/10.18653/v1/2023.eacl-main.119).
- 506 Hao Yang, Jinming Zhao, Gholamreza Haffari, and
507 Ehsan Shareghi. 2022. *Self-supervised Rewiring of*
508 *Pre-trained Speech Encoders: Towards Faster Fine-*
509 *tuning with Less Labels in Speech Processing*.
510 In *Findings of the Association for Computational*
511 *Linguistics: EMNLP 2022*, pages 1952–1959, Abu
512 Dhabi, United Arab Emirates. Association for
513 Computational Linguistics. doi:
514 [10.18653/v1/2022.findings-emnlp.141](https://doi.org/10.18653/v1/2022.findings-emnlp.141).
- 515 Hao Yang, Jinming Zhao, Gholamreza Haffari, and
516 Ehsan Shareghi. 2023. Investigating pre-trained
517 audio encoders in the low-resource condition. In
518 *24th Annual Conference of the International*
519 *Speech Communication Association*, pages 1498–
520 1502, Dublin, Ireland. International Speech
521 Communication Association. doi:
522 [10.21437/Interspeech.2023-343](https://doi.org/10.21437/Interspeech.2023-343).
- 523 Shaolei Zhang, Qingkai Fang, Shoutao Guo, Zhengrui
524 Ma, Min Zhang, and Yang Feng.
525 2024a. *StreamSpeech: Simultaneous Speech-to-*
526 *Speech Translation with Multi-task Learning*.
527 In *Proceedings of the 62nd Annual Meeting of the*
528 *Association for Computational Linguistics (Volume*
529 *1: Long Papers)*, pages 8964–8986, Bangkok,
530 Thailand. Association for Computational
531 Linguistics. doi: [10.18653/v1/2024.acl-long.485](https://doi.org/10.18653/v1/2024.acl-long.485).
- 532 Chen Zhang, Xiao Liu, Jiaheng Lin, and Yansong
533 Feng. 2024b. *Teaching Large Language Models an*
534 *Unseen Language on the Fly*. In *Findings of the*
535 *Association for Computational Linguistics: ACL*
536 *2024*, pages 8783–8800, Bangkok, Thailand.
537 Association for Computational Linguistics. doi:
538 [10.18653/v1/2024.findings-acl.519](https://doi.org/10.18653/v1/2024.findings-acl.519).
- 539 Weiquan Zhang. 2022. Zhuang Named Entity
540 Recognition Based on Deep Learning. Master's
541 thesis, Guangxi Normal University, Nanning,
542 China, June. doi:
543 [10.27036/d.cnki.ggxu.2022.001249](https://doi.org/10.27036/d.cnki.ggxu.2022.001249).