

Enhanced Training Methods for Multiple Languages

Anonymous ACL submission

Abstract

Document-grounded dialogue generation based on multilingual is a challenging and realistic task. Unlike previous tasks, it need to tackle with multiple high-resource languages facilitating low-resource languages. This paper summarizes our research based on a three-stage pipeline that includes retrieval, re-rank and generation where each component is individually optimized. In different languages with limited data scenarios, we mainly improve the robustness of the pipeline through data augmentation and embedding perturbation with purpose of improving the performance designing three training methods: cross-language enhancement training, weighted training with neighborhood distribution augmentation, and ensemble adversarial training, all of that can be used as plug and play modules. Through experiments with different settings, it has been shown that our methods can effectively improve the generalization performance of pipeline with score ranking 6th among the public submissions on leaderboards.

1 Introduction

Question Answering (QA) system has received extensive attention in recent researches. The QA system aims to provide precise answers in response to the user’s questions in natural language. An essential task in the QA system is conversational question answering and document-grounded dialogue modeling. Lack of data is one of the main challenges (Zhang et al., 2020).

Retrieval-augmented Generation (RAG) (Lewis et al., 2020) proposes a two-stage generation method with retriever extracting multiple documents related to the query and feeding them into answer generator. A survey of document-grounded dialogue systems (Ma et al., 2020) points

that it is a mainstream method to indirectly search for key text before directly generating replies. There have been various works for knowledge-grounded dialogue systems (Zhan et al., 2021; Wen et al., 2022; Ma et al., 2020) to address this problem. A new framework UniGDD (Gao et al., 2022) use prompt learning for context guidance and design multitask learning. PPTOD (Su et al., 2022) proposes a dialogue pre-trained model that implements the current SOTA.

As a more realistic task, MultiDoc2Dial (Feng et al., 2021) faces challenges of identifying useful pieces of text from documents and generating response simultaneously which is goal-oriented dialogues generation based on multiple documents. Unlike former task, Doc2dial (Zhang et al., 2023) upgrades the difficulty level by introducing multiple languages.

To alleviate the problem of limited datasets in low-resource languages, on the one hand, it is necessary to effectively utilize datasets in the other high-resource languages. On the other hand, we design three training methods. These designs are all aimed at enhancing the generalization ability of the model. Our model is based on a three-stage framework: retriever, re-ranker and generator, the aims of first and second step are obtaining the most relevant paragraphs to the question, and then generating answer text. The first stage is responsible for the coverage of relevant texts that is the comprehensiveness of input texts; in the second stage, it is necessary to filter out the most relevant text that is the accuracy of the input text; the third stage generates answers based on the input text, which is clearly the most important part. Our contributions are as follows:

- a cross language enhancement training method is designed which can effectively improve generalization ability by replacing the high-frequency tokens of

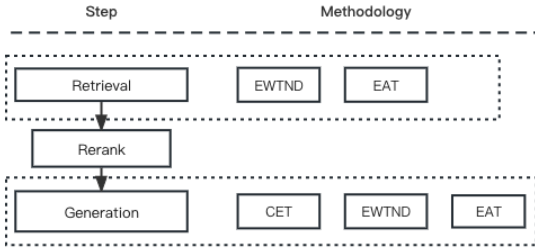


Figure 1: Training process of our pipeline.

high-resource languages with that of low-resource languages in pre-trained model.

- enhanced weighted training approach based on neighborhood distribution is presented, the diversity of input texts can be increased through data augmentation, and the problem of semantic inaccuracy can be alleviated through weight.
- ensemble adversarial training method is proposed including two classic adversarial training methods to improve the model's anti-interference ability and reduce text generation bias.

The above three enhancement training methods can be easily applied to other languages models as plug and play modules. Based on the published dataset, sufficient experiments are conducted confirming the method can effectively improve the generalization performance of the model.

2 Task Definition

Given dialogue history $\{q_1, \dots, q_{t-1}\}$ and current user's query q_t , DialDoc task need to produce the response based on knowledge from a set of relevant documents $D_0 \subseteq D$, where D denotes all knowledge documents. Besides, the task provides similar format dataset of four languages including two high-resource languages (English and Chinese) and two low-resource languages (French and Vietnamese), and the latter one is evaluated.

3 Methodology

To start with design, our pipeline is based on the three-stage baseline (Zhang et al., 2023). The three training augmentation methods that we propose can be applied to retrieval and generation. The specific framework process is as Figure 1.

3.1 Cross-Language Enhancement Training (CET)

From perspective of tokenizer, we designed a enhancement training method with token exchange between various languages. In different languages pairs, words with high frequency may have similar semantics, so that transfer learning can be used to facilitate low-resource languages training with embedding layers of high-resource languages. The basic idea is that as for pre-training model's tokenizer, replace high-resource languages' tokens with that of low-resource languages according to the rank of tokens' frequency which should follow four principles: (i) the total number of tokens of the high-resource languages need to be larger than that of the low-resource languages. (ii) select every similar language pairs, replace the high-resource tokens with low-resource tokens according to the rank order of frequency separately. In this paper, it should replace Chinese with Vietnamese and English with French. (iii) if the tokens of a language pair are insufficient, they can be mapped to the remaining unaligned tokens of another language. In this paper, there does not need to do it as the number of tokens in English higher than that of French, so do Chinese and Vietnamese. (iv) punctuation marks, [UNK] and other special marks remain unchanged.

After obtaining the mapping relationship of the tokenizer, we replace low-resource languages' datasets into high-resource languages' datasets as additional data, setting training weight w for the new one.

3.2 Enhanced Weighted Training of Neighborhood Distribution (EWTND)

To alleviate the limited datasets about low-resource languages, we propose enhanced weighted training of neighborhood distribution method. By enhancing the texts from semantic neighborhood distribution, the diversity of input text increases, and the problem of semantic inaccuracy of neighborhood distribution is alleviated through weighted training. The steps of the method are as follows: (i) in top n words $\{w_1, \dots, w_n\}$ with the highest frequency, using the last layer of pre-trained mT5 (Xue et al., 2021; Raffel et al., 2020; Zhang et al., 2020) encoder to produce 512 dimensional vectors $\{v_1, \dots, v_n\}$ for each token (except for punctuation mark). (ii) for every v , find the k words with the largest similarity through vector retrieval by Faiss (Johnson et al., 2019) vector retrieval library, and record their similarities. So we get the text neighborhood

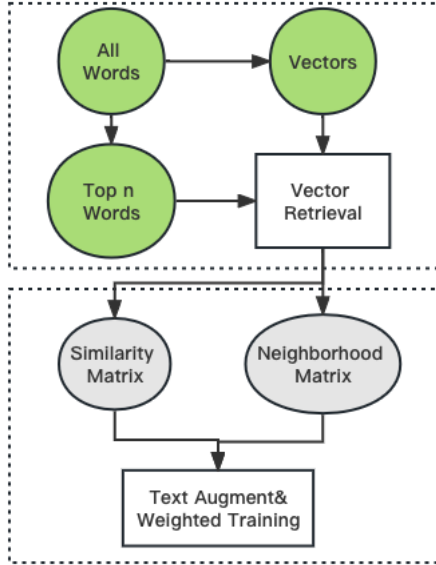


Figure 2: The key parts of EWTND.

matrix t_{ij} and similarity matrix s_{ij} , where $1 \leq i \leq n, 1 \leq j \leq k$. (iii) during training, each sentence has a $p\%$ probability to apply replacing that is words in w are replaced by one of its neighborhood from t with equal probability, and the calculation weight of sample loss is updated to the mean of similarity from s in every sentence.

3.3 Ensemble Adversarial Training (EAT)

As a regularization method, adversarial training can improve the robustness of the model by introducing perturbations in embedding (Tramèr et al., 2020; Miyato et al., 2021). We propose an ensemble adversarial training method that blend two classic adversarial training methods to improve the model's anti-interference ability and reduce text generation bias. Adversarial training can be described by a general formula as follows: (Madry et al., 2019)

$$\min_{\theta} \mathbb{E}_{(x,y) \sim D} [\max_{\Delta x \in \Omega} L(x + \Delta x, y; \theta)]$$

where D is training dataset, x is input, y is target, θ is model parameter, $L(x + \Delta x, y; \theta)$ is loss of single sample, Ω is disturbance space, Δx is perturbation. What's more, the main changes in different adversarial training methods are Δx and Ω . FGM method (Jan et al., 2015; Wong et al., 2020) raise the gradient with parameter ϵ and standardize it getting new Δx :

$$\Delta x = \epsilon \frac{\nabla L(x, y; \theta)}{\|\nabla L(x, y; \theta)\|}$$

While PGD method (Madry et al., 2019) split Δx into multiple steps, set the constraint space to a sphere:

$$\Delta x_{t+1} = \prod_{x \in S} (\Delta x_t + \alpha \frac{\nabla L(x_t, y; \theta)}{\|\nabla L(x_t, y; \theta)\|})$$

where $S = \{r \in \mathbb{R}^d, \|r\|_2 < \epsilon\}$, α is step size.

We add the FGM and PGD into training. For each batch in training process, we set the probabilities of the different training methods, there is $p_1\%$ probability of PGD, $p_2\%$ probability of FGM, and $p_3\%$ probability of not changing. The proportion can be determined by the ordinal of the model's convergence effect. In this paper, the rank of PGD, FGM, and non enhancement are 3:2:1 respectively, which means the probabilities are 50%, 33%, 17%. After multiple experiments, we believe that there is a correlation between the final convergence loss of the method and the dataset, so the all possibilities should cannot be directly set and need to be determined based on the training results.

4 Experiments

We evaluate our methods using datasets provided by shared task which include four languages. As for generator, EWTND uses French and Vietnamese dialogue generation dataset, while CET also requires English and Chinese dialogue dataset. Besides, the score is calculated based on the sum of token-level F1, SacreBleu and Rouge-L metrics.

The experiments are mainly conducted on fine-tuning the retriever and generator based on the open-source baseline in three-stage framework. All the performances of methods can be evaluated by score of generator.

| w | F1 | Sarcebleu | Rouge-L | Score |
|----------------------|-------|-----------|---------|---------------|
| 0 | 58.55 | 42.03 | 55.83 | 156.42 |
| 0.2 | 60.74 | 43.30 | 57.92 | 161.96 |
| 0.25 | 61.85 | 43.72 | 59.21 | 164.78 |
| 0.3 | 61.97 | 44.38 | 59.31 | 165.66 |
| 0.35 | 61.71 | 43.63 | 59.08 | 164.42 |
| 0 ^{half} bz | 61.13 | 43.36 | 58.18 | 162.67 |

Table 1: The results of CET on Doc2dial validation dataset.

Implementation As for CET and EWTND, when they are used in generator, we change the "passages" and "re-rank" corresponding text in dataset; when they are used in retriever, we change the "positive" and "negative" corresponding text in dataset; while "query" text and "target" text won't

240 be changed. As for EWTND, we use the cosine 277
 241 similarity. Faiss vector retrieval use product
 242 quantization to divide vector into 8 sub vectors,
 243 with 100 k-means clustering for each sub vector.
 244 There is no threshold set to limit the number of
 245 synonyms k which facilitates parallelization
 246 acceleration. We also set no limit to training epochs
 247 with early stopping epochs as 5, as EAT will need
 248 at least double training time.

249

250 **Results** Table 1 reports the performance of
 251 generator by using CET. When the weight is small,
 252 there can be a significant improvement. As weight
 253 increases to a certain extent, there will be score
 254 jitter. It proves that the CET can utilize the
 255 embedding of high-resource languages to improve
 256 low-resource languages. Meanwhile, this may also
 257 be due to more training batches. By reducing the
 258 batch size to half, it can be observed that score still
 259 improves, but under nearly equal training time,
 260 CET still achieves better results.

261

| n | k | p | Score |
|------|-----|-----|---------------|
| 500 | 1 | 0.2 | 170.23 |
| 500 | 2 | 0.2 | 172.45 |
| 500 | 3 | 0.2 | 166.38 |
| 500 | 2 | 0.3 | 171.81 |
| 1000 | 2 | 0.2 | 170.75 |

262 Table 2: The results of EWTND on Doc2dial
 263 validation dataset.

264 Table 2 shows the effect of generator by using
 265 EWTND, it still use CET and EWTND but only
 266 strengthen the origin data. When k increases from
 267 2 to 3, the reason why score drops might be
 268 uncertainty of the neighborhood’s semantic
 269 meaning, the same reason can explain the time
 270 when n increases.

271

| p_1 | p_2 | p_3 | Score |
|-------|-------|-------|---------------|
| 100% | 0% | 0% | 175.05 |
| 0% | 100% | 0% | 172.45 |
| 50% | 33% | 17% | 175.39 |
| 60% | 25% | 15% | 174.48 |
| 45% | 35% | 20% | 173.60 |

272 Table 3: The results of EAT on Doc2dial validation
 273 dataset.

274 Table 3 shows the ensemble effect of adversarial
 275 training, it proves that such training method will
 276 provide stable improving although not much.

| Method | EWTND | EAT | CET | Score |
|-----------|-------|-----|-----|--------|
| Retriever | ✓ | | | 181.57 |
| Retriever | ✓ | ✓ | | 181.60 |
| mT5 | | | | 173.42 |
| mT5 | | | ✓ | 183.05 |
| mT5 | ✓ | | ✓ | 186.71 |
| mT5 | ✓ | ✓ | ✓ | 188.62 |

278 Table 4: The results of adding training methods into
 279 other models on Doc2dial validation dataset.

280 Table 4 shows effectiveness of three training
 281 methods as plug and play modules. By enhancing
 282 the retriever, the generator still improves but
 283 disadvantage is that it increases training time
 284 around 1.5 times. Besides, the improved
 285 performance is not as good as methods applied to
 286 the generator. With the best retriever and origin re-
 287 ranker, we replace the generator with origin mT5
 288 (Xue et al., 2021) model which shows that it is
 289 better than generator in baseline. Finally, we
 290 achieve best performance by adding three
 291 enhanced training methods into mT5.

292 The above experiments have shown that our
 293 methods have significant advantages: (i) three
 294 training methods can effectively increase model’s
 295 performance without affecting prediction speed. (ii)
 296 almost all language models with token as input can
 297 apply these methods. (iii) the methods can have
 298 more potentials in future work, especially in cross
 299 language scenarios, EWTND can be extended to
 300 more similar language pairs; EAT can use more
 301 complex sampling methods based on the neighbor-
 302 hood distribution of different languages.

303 5 Conclusion

304 In this paper, we propose three training methods to
 305 improve model’s performance from perspective of
 306 embedding enhancement and data augmentation.
 307 CET Introduces cross language learning through
 308 high-frequency words; EWTND use weighted
 309 augmentation from the neighborhood distribution
 310 of high-frequency words; EAT strengthen the
 311 robustness of the model through embedding
 312 perturbation. Compared to the baseline mode, our
 313 methods achieve the stable rise in score.

314 References

315 Song Feng, Siva Sankalp Patel, Hui Wan, and
 316 Sachindra Joshi. 2021. Multidoc2dial: Modeling

- 317 dialogues grounded in multiple documents. 370 Longxuan Ma, Wei-Nan Zhang, Mingda Li and Ting
318 *In EMNLP*. 371 Liu. 2020. A Survey of Document Grounded
319 Zhang, Yeqin and Fu, Haomin and Fu, Cheng and Yu, 372 Dialogue Systems (DGDS). *arXiv preprint arXiv:*
320 Haiyang and Li, Yongbin and Nguyen, Cam-Tu. 373 *2004.13818*.
- 321 2023. Coarse-to-Fine Knowledge Selection for 374 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine
322 Document Grounded Dialogs. *2023 IEEE* 375 Lee, Sharan Narang, Michael Matena, Yanqi Zhou,
323 *International Conference on Acoustics, Speech and* 376 Wei Li, Peter J. Liu. 2020. Exploring the Limits of
324 *Signal Processing*. 377 Transfer Learning with a Unified Text-to-Text
325 Patrick S. H. Lewis, Ethan Perez, Aleksandra Pik- 378 Transformer. *arXiv preprint arXiv: 1910.10683*.
- 326 tus, 379 Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Ian
327 Fabio Petroni, Vladimir Karpukhin, Naman Goyal, 380 Goodfellow, Dan Boneh and Patrick McDaniel.
328 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim 381 2020. Ensemble Adversarial Training: Attacks and
329 Rocktäschel, Sebastian Riedel, and Douwe Kiela. 382 Defenses. *arXiv preprint arXiv: 1705.07204*.
- 330 2020. Retrieval-augmented generation for 383 Takeru Miyato, Andrew M. Dai, Ian Goodfellow. 2021.
331 knowledge-intensive NLP tasks. *In Advances in* 384 Adversarial Training Methods for Semi-Supervised
332 *Neural Information Processing Systems 33:* 385 Text Classification. *arXiv preprint arXiv:*
333 *Annual Conference on Neural Information* 386 *1605.07725*.
- 334 *Processing Systems 2020, NeurIPS 2020, December*
335 Haolan Zhan, Lei Shen, Hongshen Chen, and Hainan 387 Longxuan Ma, Weinan Zhang, Runxin Sun, Ting Liu.
336 Zhang. 2021. CoLV: A collaborative latent variable 388 2020. A Compare Aggregate Transformer for
337 model for knowledge-grounded dialogue 389 Understanding Document-grounded Dialogue.
338 generation. *In Proceedings of the 2021 Conference* 390 *arXiv preprint arXiv: 2010.00190*.
- 339 *on Empirical Methods in Natural Language* 391 Zheng Zhang, Ryuichi Takanobu, Qi Zhu, Minlie
340 *Processing, pages 2250–2261, Online and Punta* 392 Huang and Xiaoyan Zhu. 2020. Recent Advances
341 *Cana, Dominican Republic. Association for*
- 342 *Computational Linguistics*. 393 and Challenges in Task-oriented Dialog System.
343 Xiaofei Wen, Wei Wei and Xian-Ling Mao. 2022. 394 *arXiv preprint arXiv: 2003.07490*.
- 344 Sequential Topic Selection Model with Latent 395 Jingqing Zhang, Yao Zhao, Mohammad Saleh, Peter J.
345 Variable for Topic-Grounded Dialogue. *In* 396 Liu. 2020. PEGASUS: Pre-training with Extracted
346 *Proceedings of the 2022 Conference on Empirical* 397 Gap-sentences for Abstractive Summarization.
347 *Methods in Natural Language Processing (Findings* 398 *arXiv preprint arXiv: 1912.08777*.
- 348 *of EMNLP'2022), Abu Dhabi*. 399 Eric Wong, Leslie Rice, J. Zico Kolter. 2020. Fast is
349 Ian J. Goodfellow, Jonathon Shlens and Christian 400 better than free: Revisiting adversarial training.
350 Szegedy. 2015. Explaining and Harnessing 401 *ICLR 2020*.
- 351 Adversarial Examples. *arXiv preprint arXiv:* 402 Yixuan Su, Lei Shu, Elman Mansimov, Arshit Gupta,
352 *1412.6572*. 403 Deng Cai, YiAn Lai, Yi Zhang. 2022. Multi-Task
353 Aleksander Madry, Aleksandar Makelov, Ludwig 404 Pre-Training for Plug-and-Play Task-Oriented
354 Schmidt, Dimitris Tsipras and Adrian Vladu. 2019. 405 Dialogue System. *Proceedings of the 60th Annual*
355 Towards Deep Learning Models Resistant to 406 *Meeting of the Association for Computational*
356 Adversarial Attacks. *arXiv preprint arXiv:* 407 *Linguistics*.
357 *1706.06083*. 408
- 358 Linting Xue, Noah Constant, Adam Roberts, Mihir 359 Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua
359 and Colin Raffel. 2021. mT5: A massively
360 multilingual pre-trained text-to-text transformer.
361 *arXiv preprint arXiv: 2010.11934*.
- 362 Johnson, Jeff Douze, Matthijs Jegou, Herve. 2019.
363 Billion-scale similarity search with {GPUs}. *IEEE*
364 *Transactions on Big Data*.
- 365 Chang Gao, Wenxuan Zhang and Wai Lam. 2022.
366 UniGDD: A Unified Generative Framework for
367 Goal-Oriented Document-Grounded Dialogue.
368 *arXiv preprint arXiv: 2204.07770*.
- 369