

Boosting Translation Capabilities of Large Language Models with Code-Switching Pretraining

Anonymous ACL submission

Abstract

001 Recently, there has been significant attention
002 on adapting the translation capabilities of Large
003 Language Models. Represented by ALMA (Xu
004 et al., 2023), a two-stage training recipe has
005 been developed: first, utilizing a large amount
006 of monolingual data for pretraining to enhance
007 proficiency in non-English languages, followed
008 by fine-tuning with a small amount of high-
009 quality bilingual data. However, in the pretrain-
010 ing process, explicit cross-lingual alignment
011 information is not provided, and excessive use
012 of bilingual data can lead to catastrophic for-
013 getting issues, both of which hinder the further
014 advancement of the model’s translation abili-
015 ties. In this article, we address this issue by
016 introducing a new pretraining process based
017 on Code-Switching pretraining data. In this
018 stage of pretraining, we can provide rich cross-
019 lingual alignment information while ensuring
020 that the training data is semantically coherent
021 documents, which helps alleviate catastrophic
022 forgetting. Moreover, the training process re-
023 lies solely on monolingual data and a pair of
024 traditional machine translation models, making
025 it highly versatile. Experimental results show
026 that our method has improved the translation
027 quality, achieving state-of-the-art results in sim-
028 ilar works.

029 1 Introduction

030 The rapid development of large language models
031 (LLMs) (Brown et al., 2020b; Chowdhery et al.,
032 2023; Touvron et al., 2023), represented by the
033 GPT series (Brown et al., 2020b), has brought ex-
034 citing progress to the field of Natural Language
035 Processing (NLP). The powerful language under-
036 standing, abstract summarization, and conversa-
037 tional generation capabilities of large models are
038 revolutionizing numerous NLP tasks (Shao et al.,
039 2023; Singhal et al., 2023; Zhang et al., 2024; Min
040 et al., 2023), and the field of machine translation is
041 no exception.

042 Extensive work has verified that large models
043 can achieve zero-shot and few-shot translation
044 through their powerful in-context learning (Hendy
045 et al., 2023; Zhang et al., 2023a) capabilities, with-
046 out the need for specific adaptations for translation
047 tasks. However, since large language models are
048 often trained on English as the primary language,
049 the insufficient data in other languages has resulted
050 in most LLMs’ translation capabilities falling short
051 compared to commercial traditional models or top
052 commercial LLMs (Jiao et al., 2023b).

053 ALMA (Xu et al., 2023) has already proven that
054 we can enhance the translation capabilities of LLM
055 through continual training. They first enhance the
056 proficiency of LLM in these non-English languages
057 by adding monolingual data in those languages for
058 Continual Pretraining (CPT), and then stimulate the
059 translation capabilities of LLM by utilizing small
060 amount high-quality bilingual data for Supervised
061 Finetune (SFT). However, in the pretraining phase
062 they proposed, there was no explicit modeling of
063 cross-lingual alignment, which may hinder further
064 quality improvement. In contrast, Guo et al. (2024)
065 attempted to adding sentence-level parallel data
066 during the pretraining phase using an Interlinear
067 text format, but this method has two drawbacks: a)
068 Extensive sentence-level bilingual data has been
069 demonstrated to induce catastrophic forgetting (Xu
070 et al., 2023) and erase previously acquired knowl-
071 edge. b) The pattern of sentence-level parallel data
072 diverges from that of standard pretraining data, ne-
073 cessitating high-quality, semantically coherent doc-
074 ument data.

075 In studies of traditional machine translation
076 (MT) models, the code-switching strategy (Lin
077 et al., 2020; Yang et al., 2020) (i.e., replacing
078 words or phrases in the current sentence with ex-
079 pressions from another language) has been shown
080 to be effective in aligning multilingual representa-
081 tion spaces. Drawing inspiration from this, we at-
082 tempted to modify the pretraining corpus of LLMs

083 using a sentence-level code-switching strategy and
084 obtained semantically coherent document data com-
085 posed of sentences from two languages. Then LLM
086 can learn cross-lingual contextual dependencies
087 and alignment information embedded in the such
088 data through standard pretraining.

089 To achieve our goals, the most ideal training
090 data is document-level parallel corpora, but such
091 data only exists between high-resource languages
092 and in limited quantities. Nevertheless, leveraging
093 the strong fundamental capabilities of LLMs along
094 with specific markers enables us to utilize tradi-
095 tional MT models for generating document-level
096 back translation data as an alternative.

097 More specifically, we use monolingual data in
098 English and the target language, along with a pair
099 of traditional MT models, to generate two types of
100 code-switching pretraining data: from English to
101 the target language and from the target language to
102 English. Subsequently, a novel pretraining phase,
103 denote as Code-Switching Continual PreTraining
104 (standard pretraining on the code-switching data),
105 is integrated into the two-stage training recipe sug-
106 gested by ALMA. In the end, experiments show
107 that our improved training recipe significantly en-
108 hances LLM’s cross-lingual alignment capability.
109 The translation quality from the target language to
110 English and from English to the target language has
111 both improved. At the same time, we found that in
112 the new pretraining stage, the contribution of code-
113 switching pretraining data in the same direction is
114 greater than in the opposite direction, and we pre-
115 liminarily analyze that such data may help improve
116 the model’s automatic post-editing capabilities.

117 Our core contributions are as follows:

- 118 • The Code-Switching Continual PreTraining
119 stage we proposed can enhance the cross-
120 lingual alignment capability of LLM, address-
121 ing the shortcomings of previous work.
- 122 • We introduced traditional MT models into the
123 optimization process of LLM’s translation ca-
124 pabilities in the form of back translation.
- 125 • The final optimized model achieved State-of-
126 the-Art performance in some translation direc-
127 tions compared to similar works.

128 2 Related Work

129 2.1 Large Language Models

130 Large language models generally refer to
131 transformer-based (Vaswani et al., 2017) neural

132 models with billions of parameters. Both open-
133 source models like Llama (Touvron et al., 2023),
134 Mistral (Jiang et al., 2023) and GLM (Zeng et al.,
135 2022) and closed-source models like GPT-3.5/4
136 (Brown et al., 2020a), Claude (Anthropic) demon-
137 strate enhanced language comprehension and gener-
138 ation capabilities. Mainstream LLMs follow a
139 Decoder-only architecture, expanding their param-
140 eter size by layering Transformer decoder units.
141 During training, LLMs initially undergo pretrain-
142 ing on a diverse range of document-level monolin-
143 gual data (such as internet data, books, code, etc.)
144 to establish a foundational model. Subsequently,
145 they undergo training using algorithms like Super-
146 vised Finetune and RLHF (Ouyang et al., 2022)
147 to align with human preferences and ultimately
148 achieve a robust multi-turn Instruct/Chat model for
149 diverse tasks.

150 When adapting LLM for downstream tasks,
151 there are two common strategies: the Incontext-
152 Learning (ICL) strategy based on prompt (Zhu
153 et al., 2023) technology and various evolving tech-
154 niques represented by COT (Wei et al., 2022). An-
155 other strategy involves fine-tuning (Ding et al.,
156 2023) the model using downstream data, which
157 often achieves higher performance. Technologies
158 like Low-rank Adaptation (LoRA) (Hu et al., 2021),
159 which solve the training cost issue, significantly en-
160 hance the applicability of this strategy.

161 2.2 Machine Translation Task

162 **Traditional Methods** The traditional machine
163 translation models, represented by transformers
164 (Vaswani et al., 2017), utilize an Encoder-decoder
165 architecture to autoregressively decode the target
166 language. Among various optimization methods,
167 data augmentation (Burlot and Yvon, 2018) tech-
168 niques like Back Translation (BT) (Edunov et al.,
169 2018; Hoang et al., 2018; Pham et al., 2021) has
170 been proven to be more effective. BT comprises
171 different variations such as sampling BT, Noise
172 BT, Tag BT (Caswell et al., 2019), and so on. In
173 the training phase, BT incorporates a variety of
174 monolingual data in the target language to boost
175 the language model’s capabilities, aiding in pro-
176 ducing more natural and accurate outputs (Edunov
177 et al., 2020). Additionally, Forward Translation
178 (FT), which translates the source text into the tar-
179 get language, is frequently paired with BT data.

180 **LLM-based Methods** As we mentioned before,
181 when adapting the translation capability of LLMs,
182 there are two types of strategies. The first type

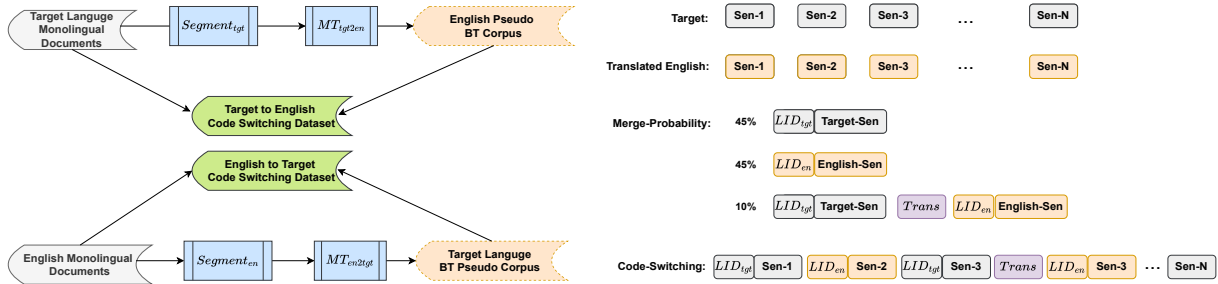


Figure 1: Construction process of Code-Switching pretraining data. The left side displays the key flow nodes involved in data construction, while the right side elaborates on the specifics of constructing Code-Switching data using original monolingual and BT pseudo-corpus. LID and $Trans$ are special tokens.

183 focuses on harnessing LLMs’ Incontext-learning
 184 feature and employing prompt techniques to en-
 185 hance the model’s translation ability. Many studies
 186 (Hendy et al., 2023; Zhang et al., 2023a; Wang
 187 et al., 2023; Gulcehre et al., 2017) have conducted
 188 detailed explorations in this direction. Another type
 189 involves fine-tuning the model with specific data
 190 from translation tasks to achieve better translation
 191 quality. Different studies may attempt to fine-tune
 192 the model at different training stages. For example,
 193 fine-tuning the model with monolingual data (Tan
 194 et al., 2023; Yang et al., 2023; Wei et al., 2023)
 195 in the target language or domain during the pretrain-
 196 ing phase. Alternatively, using translation-related
 197 instruction (Li et al., 2024; Zhang et al., 2023b)
 198 data during the SFT phase. Xu et al. (2024) aim
 199 to enhance translation quality by fine-tuning the
 200 model using comparison data with varying quality
 201 through reinforcement learning.

202 3 Methodology

203 In this chapter, we will describe the details of the
 204 code-switching strategy we proposed, as well as
 205 the training recipe we optimized for adapting the
 206 translation capabilities of LLMs.

207 3.1 Code-Switching Pretraining Data

208 In the traditional MT and multilingual language
 209 model (MLM) field, the code-switching strategy
 210 has been proven to provide cross-lingual alignment
 211 information (Lin et al., 2020; Yang et al., 2020).
 212 In order to adapt to pretraining tasks for LLMs,
 213 we use a sentence-level code-switching strategy
 214 and obtain semantically coherent document data
 215 composed of sentences from two languages.

216 In Figure 1, we illustrate the specific approach.
 217 We refer to the target language as tgt and English
 218 as en . Utilizing a pair of pre-trained traditional MT
 219 models, we translate monolingual English and tar-

220 get language corpora to generate BT pseudo-corpus
 221 denote as D_{bt} . When constructing Code-Switching
 222 pretraining data (D_{cs}), we randomly select origi-
 223 nal and translated sentences with equal probability,
 224 and with a 10% probability, we allow them to ap-
 225 pear simultaneously. To effectively differentiate
 226 Code-switching data and prevent language confu-
 227 sion during inference, we incorporate some special
 228 tokens. The design of special tokens and an ex-
 229 ample of Code-Switching data are provided in the
 230 Appendix A.

231 3.2 A New Training Recipe

232 We proposed a new training recipe, in which we
 233 added a "Code-Switching Continual Pretraining"
 234 stage to ALMA’s two-stage training recipe, aiming
 235 to more efficiently inject cross-lingual alignment
 236 information. Figure 2 illustrates our training recipe
 237 and the differences between our work and typical
 238 similar works.

239 **Stage-1: Continual Pretraining with Monolin-
 240 gual Data** LLMs like LLaMA are pre-trained on
 241 English-dominated corpora. They may encounter
 242 issues with insufficient comprehension and gener-
 243 ation abilities in the target language. By incorpo-
 244 rating a large amount of monolingual data in the
 245 target language for continual pretraining, we can
 246 alleviate this issue. At this stage, we can train with
 247 the full set of parameters or utilize LoRA technol-
 248 ogy to enhance training efficiency. Training data
 249 often comes from widely available internet sources,
 250 such as Common Crawl (Foundation, 2023), as
 251 well as some cleaned versions like OSCAR (Ortiz
 252 Su’arez et al., 2019; Kreutzer et al., 2022). It is
 253 worth noting that the outcome of this stage is to
 254 obtain a foundational LLM with multilingual capa-
 255 bilities, where we can conduct the training process
 256 ourselves or obtain pre-trained models from the
 257 open-source community.

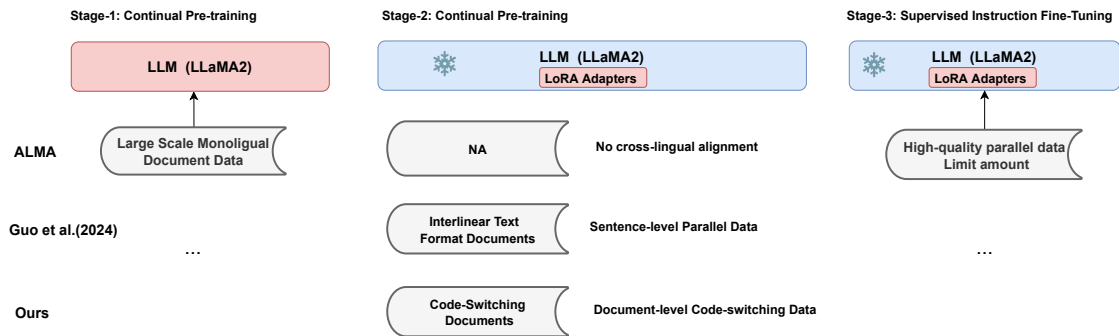


Figure 2: Training process for our and similar works. Overall, we use a three-stage training recipe. And by using Code-Switching strategy, we provide rich cross-lingual alignment information to solve the problems faced in previous works.

Stage-2: Code-Switching Continual Pretraining (CS-CPT) Since our main objective is to enhance the translation capabilities of LLM, cross-lingual alignment information holds significant importance. During the initial training phase, the absence of explicit cross-lingual alignment information necessitates the LLM to learn implicitly, which may not be the most efficient method. We mitigate this issue by performing Continual Pretraining on Code-switching data (presented in 3.1). And CS-CPT offering three key advantages:

- Code-Switching data explicitly provides cross-lingual contextual dependencies, which can compel the model to learn semantic-level alignment relationships.
- Code-Switching data is essentially semantically coherent document data, which maintains consistency with the standard pretraining data format and can alleviate catastrophic forgetting.
- It only requires an additional pair of traditional MT models, making resource consumption and complexity controllable.

We use LoRA technology to carry out the training in this stage, but set the *embed_tokens* and *lm_head* modules to be updatable so that the model can learn token-level alignment information. It is worth mentioning that, as the data pattern is consistent with the first stage, we can even merge them together for training. We also validated this point in the subsequent experimental section.

Stage-3: High-Quality Data Fine-tuning In previous research on adapting LLMs to downstream tasks, it has been confirmed that the quality of data during the SFT phase is more important than the quantity (Zhou et al., 2023; Maillard et al.,

2023; Gunasekar et al., 2023) of data. Following the settings of previous works ALMA and Guo et al. (2024), we use a small amount of high-quality bilingual data to fine-tune the model in order to enhance its translation capabilities. To ensure data quality, we collect human-written datasets from WMT development and test sets. We also employ LoRA for training.

4 Experiments

We mainly tested our algorithm on translation tasks in four directions in two language pairs: English-Chinese and English-German. Our experiment design closely follows ALMA to ensure a fair comparison.

4.1 Datasets and Evaluation Metrics

The monolingual dataset we used is sourced from OSCAR. Since the base model we chose (Chinese-LLaMA-2 (Cui et al., 2023)) has already undergone the first stage of pretraining in Chinese, we selected only 0.5B of Chinese and English data from the OSCAR dataset for the second stage of training. For the English-German translation task, we opted to pretrain with 1.5B of German and English monolingual data (the average number in the ALMA’s experiments) and similarly used 0.5B for the second stage of training.

For our parallel training data, we collect human-written test datasets from WMT’17 to WMT’20 for EN \leftrightarrow ZH and EN \leftrightarrow DE resulting in a total of 37.6K training examples across all languages.

Furthermore, we include the test sets from the WMT22 competition, which are thoughtfully curated to encompass recent content from various domains like news, social media, e-commerce, and conversations.

Models	De⇒En		En⇒De		Zh⇒En		En⇒Zh	
	BLEU	COMET	BLEU	COMET	BLEU	COMET	BLEU	COMET
SoTA models								
NLLB-54B(Team et al., 2022b)	26.89	78.94	34.50	86.45	16.56	70.70	27.38	78.91
GPT-3.5-D, zero-shot	30.90	84.79	31.80	85.61	25.00	81.60	38.30	85.76
GPT-3.5-T, zero-shot	33.10	85.50	34.40	87.00	26.60	82.90	44.90	87.00
GPT-4	33.87	85.62	35.38	87.44	27.20	82.79	43.98	87.49
Prior Similar Studies								
TIM-7B(Zeng et al., 2023)	27.91	82.80	25.59	82.56	19.33	75.46	19.33	75.46
Parrot-7B(Jiao et al., 2023a)	29.80	83.00	26.10	81.60	20.20	75.90	30.30	80.30
SWIE-7B(Chen et al., 2023)	30.48	82.97	27.21	82.36	21.30	76.48	31.24	80.63
ALMA-7B(Xu et al., 2023)	29.56	83.95	30.31	85.59	23.64	79.78	36.48	85.05
Guo et al. (2024)	31.14	84.70	30.50	85.62	22.20	79.88	41.10	86.37
Parrot-13B(Jiao et al., 2023a)	31.10	83.60	28.10	82.60	21.70	76.70	31.70	81.00
BigTranslate-13B(Yang et al., 2023)	23.35	80.68	21.48	78.81	14.16	74.26	28.56	81.31
Bayling-13B(Zhang et al., 2023b)	27.34	83.02	25.62	82.69	20.12	77.72	37.92	84.62
ALMA-13B(Xu et al., 2023)	31.14	84.56	31.47	85.62	25.46	80.21	39.84	85.96
Guo et al. (2024)	32.24	85.17	32.53	86.14	23.10	80.53	42.30	86.65
Traditional Back Translation Model								
NLLB-distilled-600M-Finetune	26.80	78.53	30.01	85.07	19.72	74.89	33.24	80.76
Ours	Our Recipe with Backbone Model: LLaMA2(Touvron et al., 2023)							
7B Stage1,3	30.05	84.07	30.21	85.55	23.96	79.62	35.31	84.74
7B Stage1,2,3	31.64	85.01	31.20	85.71	26.87	80.44	41.81	86.12
13B Stage1,3	31.20	84.43	31.30	85.77	24.31	80.01	37.34	85.27
13B Stage1,2,3	32.74	85.48	32.49	86.20	27.16	81.06	42.84	86.63

Table 1: **The main results.** Bold numbers represent the best scores among prior similar studies. After integrating CS-CPT, the translation quality of the model has been significantly improved. Our 7B and 13B models have achieved top performance in most evaluation metrics compare to similar studies. Even the BLEU score for the Zh⇒En direction is on par with that of GPT-4.

For automatic evaluation, we utilize SacreBLEU, which implements BLEU(Papineni et al., 2002), and COMET(Rei et al., 2020) from Unbabel/wmt22-comet-da. SacreBLEU calculates similarity based on n-gram matching, while COMET leverages cross-lingual pretrained models for evaluation. We rely more on COMET than BLEU due to its better alignment with human evaluations (Freitag et al., 2022).

4.2 Training Setup

Our experiments were carried out using HuggingFace Transformers¹ with open-source LLaMA (Touvron et al., 2023) family as our foundation model. Most of our verification experiments were conducted on 7B model, but we will also report the results of the 13B model to explore the impact of model size.

Specifically, we chose to use Chinese-LLaMA2 (Cui et al., 2023) as the base model for our training because it handles Chinese more efficiently (using expanded vocabulary) and has already completed the first stage of training in Chinese. Building on this foundation, we can proceed with the second

and third stages of training for Chinese tasks. For German tasks, we will execute the training of the first and second stages together.

In the training of the first and second stages, we use the LoRA approach to adapt the key, query, value, and output layers of the self-attention mechanism, and the LoRA hyperparameters are set to $R = 32$ and $a = 64$. At the same time, the modules *embed_tokens* and *lm_head* are also set as updatable parameters. We fine-tune the foundation model for one epoch using a batch size of 256, a warm-up ratio of 0.01, and sequences with a maximum of 1024 tokens in total.

During the third stage of training, we follow the ALMA’s approach by updating only 0.1% of the parameters using LoRA. We train the model for 2 epochs and select the best model based on the lowest validation loss. For both stages, we adopt deepspeed (Rasley et al., 2020) to accelerate our training.

We employ the NLLB-600M-distil² as our traditional MT model for BT pseudo-corpus. Additionally, we leveraged training data from WMT21 to improve the translation quality for the target lan-

¹<https://huggingface.co/docs/transformers/en/index>

²<https://github.com/facebookresearch/fairseq/tree/nllb>

376 guages German and Chinese, thereby ensuring the
377 fundamental quality of the BT pseudo-corpus.

378 4.3 Baselines

379 We compare our method against two baseline cat-
380 egories. Firstly, we examine previous studies that
381 share our objective of utilizing LLMs for transla-
382 tion. Secondly, we assess against the latest state-of-
383 the-art translation models.

384 **Prior Similar Work** We compare our model
385 with BigTranslate (Yang et al., 2023), which ex-
386 tends LLaMA-1-13B to over 100 translation di-
387 rections; TIM (Zeng et al., 2023), which uses cor-
388 rect and incorrect examples to help LLM to learn
389 translation; SWIE (Chen et al., 2023), which im-
390 proves LLM in translation via instruction augmen-
391 tation; Parrot(Jiao et al., 2023a), through three
392 types of instructions improves the translation per-
393 formance of LLM after SFT; and BayLing (Zhang
394 et al., 2023b), which uses interactive translation
395 instructions; and ALMA (Xu et al., 2023), a two-
396 stage finetuning method that initially fine-tunes on
397 monolingual data and subsequently on a small set
398 of high-quality parallel data; and Guo et al. (2024),
399 expand on ALMA’s approach by introducing an
400 additional stage for fine-tuning with parallel sen-
401 tences with Interlinear text format.

402 **SoTA Models** We focus on the NLLB-54B
403 model, the top-tier translation model in the NLLB
404 family (Team et al., 2022a), as well as the zero-shot
405 capabilities of GPT3.5-text-davinci-003 (**GPT-3.5-
406 D**) and GPT-3.5-turbo-0301 (**GPT-3.5-T**), along
407 with GPT-4³.

408 5 Results

409 **Main Results** Table 1 summarizes the main re-
410 sults of our experiments. In summary, our final
411 optimized model has shown consistent improve-
412 ment in translation quality, surpassing ALMA in
413 both BLEU and COMET metrics. The improve-
414 ment in the Chinese translation task is greater than
415 that in the German task, and the BLEU metric for
416 ZH⇒EN task even on par with GPT-4. Compared
417 to similar works, the 13B model has achieved a
418 leading position in most metrics.

419 **Effectiveness of Code-Switching Continual
420 Pretraining** The training in the second stage in-
421 deed improved the model’s translation ability. Tak-
422 ing Chinese tasks as an example, the COMET

423 scores for ZH⇒EN and EN⇒ZH improved by
424 0.82 and 1.38, while BLEU scores improved by
425 2.91 and 6.5, respectively. For the German task,
426 the overall trend is consistent with Chinese, but the
427 improvement is slightly smaller relative to Chinese.
428 This may be because the alignment information in
429 the foundation models for German and English is
430 richer compared to Chinese (with a higher charac-
431 ter overlap rate).

432 **Compared with Prior Similar Studies** Com-
433 pared to the strong baseline ALMA, our 7B model
434 achieved an average BLEU improvement of 2.88
435 and a COMET improvement of 0.73. Our results
436 exceed those of Guo et al. (2024) in the tasks for
437 ZH⇒EN and DE⇒EN, but are on par with theirs
438 in the EN⇒ZH and EN⇒DE directions. But it is
439 important to note that we did not use parallel cor-
440 pora in our training process. Moreover, unlike the
441 OSCAR data that we employed, they utilized the
442 WMT bilingual training data, which is more closer
443 to the domain of the current test set.

444 6 Analysis

445 In this chapter, we will analyze the key points of
446 the model. Some analysis experiments will be con-
447 ducted on Chinese tasks because Chinese and En-
448 glish have relatively greater linguistic distances.

449 6.1 Cross-lingual alignment analysis

450 To verify whether our model in the second stage
451 has achieved the goal of cross-lingual alignment,
452 we referenced relevant works (Lin et al., 2020)
453 and conducted quantitative analyses in two dimen-
454 sions. Firstly, we calculated the similarity of word
455 embeddings for words with the same meanings in
456 different languages. We selected the top 1000 most
457 frequent words according to the MUSE⁴ dictionary.
458 We averaged the sub-word sequences of words to
459 obtain word embeddings and calculated the cosine
460 similarity between the two languages. Additionally,
461 we analyzed representations at the sentence level
462 for sentences with the same meanings. We used
463 the Flores (NLLB Team, 2022) test set to calculate
464 sentence-level embeddings using the same method
465 and computed the corresponding similarities. The
466 results of stage-1 and stage-2 pretraining models
467 are summarized in Figure 3. From the figure, it is
468 evident that, both word and sentence-level similar-
469 ities have significantly improved after our CS-CPT,
470 regardless of whether it is language pairs with rela-

³GPT-3.5-D, GPT-3.5-T and GPT-4 results are sourced from Xu et al., 2023

⁴<https://github.com/facebookresearch/MUSE>

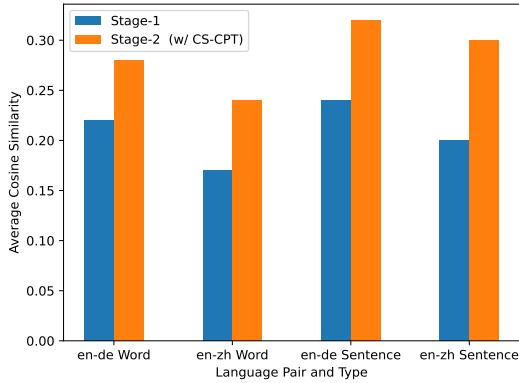


Figure 3: The average cosine similarity results of models from various stages are sourced from the 7B version. We observe an increase in similarity after the second training stage, affirming the effectiveness of our training approach.

tively close distances like EN-DE or distant pairs like Chinese-English. This once again proves that CS-CPT can indeed serve the intended purpose, aligning the model’s cross-lingual representations to some extent.

6.2 Using of traditional MT models

When creating Code-Switching data, we introduced a of traditional sentence-level MT model to ensure the method’s versatility and overcome challenges in obtaining document-level parallel corpora or document-level MT models. The results in Table 1 indicate that they did not achieve higher translation quality in terms of BLEU and COMET scores compared to the first-stage model. This finding dismisses the idea of LLM gaining knowledge via distillation from pseudo-corpus affirms that the model acquired alignment information beneficial for translation from the Code-Switching data after training in the second stage.

6.3 FT is more effective than BT?

Back translation is more effective than forward translation during the optimization of traditional machine translation models. For instance, when optimizing the ZH⇒EN model, the pseudo-corpus in the EN⇒ZH direction is typically more effective. This is because back translation introduces a large amount of monolingual data for the target language side, enhancing the generation capability of the target language (Edunov et al., 2018). With LLMs having already learned a significant amount of monolingual data during the pre-training phase, the target language’s generation ability is already

Models	ZH⇒EN		EN⇒ZH	
	BLEU	CMT	BLEU	CMT
7B Stage-1,3	23.96	79.62	35.31	84.74
7B Stage-1,2,3	26.87	80.44	41.81	86.12
Only D_{cs}^{zh2en}	26.30	80.10	38.01	85.32
Only D_{cs}^{en2zh}	24.70	79.80	39.17	85.60

Table 2: Comparative experimental results of code-switching data between BT and FT. "Only D_{cs}^{zh2en} " means using only ZH⇒EN Code-Switching data for training stage-2.

strong. Does this conclusion still hold when adapting LLM to translation tasks?

To explore this, we conducted comparative experiments on Chinese tasks. Specifically, in our CS-CPT stage, we only used code-switching data in one direction, then obtained the final translation model through the third stage of SFT. The results are summarized in Table 2. We were surprised to find that the improvement brought by forward translation is significantly better than that of back translation. Taking ZH⇒EN task as an example, using only ZH⇒EN direction code-switching data resulted in an improvement ratio of over 70% compared to using a mixture of data from both directions, while the improvement ratio for the quality of EN⇒ZH task was only around 25%. The overall trend for EN⇒ZH task is similar, just not as pronounced as with ZH⇒EN task.

We speculate that apart from bringing benefits in cross-lingual alignment, the forward translation data has also boosted the Automatic Post-Editing (APE) capability of LLM. During the CS-CPT stage, we used special tokens to mark code-switching pseudo data, guiding the model to differentiate between real and pseudo data. In the final SFT stage, the humans-written parallel data inspired the model to output sentences that lean towards real data during translation. By comparing these two types of data, LLM has improved its ability to rewrite machine-translation results into more natural and fluent translations.

To validate our speculation, we conducted a simple test on the APE capabilities of the models from the first and second stages. Specifically, we used the traditional MT model to translate the test set and obtained machine-translation results, then generated APE results using the 3-shot learning. Evaluation results are summarized in Table 3. The APE ability of the second-stage model is stronger than that of the first-stage model, with an aver-

Models	ZH⇒EN		EN⇒ZH	
	BLEU	CMT	BLEU	CMT
NLLB-distil	19.72	74.89	33.24	80.76
Stage1 + APE	20.20	75.24	33.56	82.11
Stage1,2 + APE	20.41	75.78	33.61	82.37

Table 3: Results of APE ability tests for pre-trained models at different stages. The results are all from the 7B version of the model, and the testing method is 3-shot learning. "NLLB-distil" is our traditional MT model used for translating BT pseudo-corpus.

Models	ZH⇒EN		EN⇒ZH	
	BLEU	CMT	BLEU	CMT
7B Stage1,2,3	26.87	80.44	41.81	86.12
CPT-InterLinear	24.12	79.79	37.87	85.54
+ 5-Epoch	23.87	79.36	37.88	85.46
SFT + BT	23.45	78.86	34.57	83.15
+ Full Data	23.01	78.49	34.55	82.71

Table 4: Results of ablation experiments. "CPT + InterLinear" represents replacing D_{cs} with data in InterLinear text format. "SFT + BT" means using BT translation data to replace human-writing data for the stage-3 training with equal data volum. "Full Data" denote using all the BT data.

age COMET improvement of over 0.4 for the final translation results. Further in-depth exploration will be left for future research.

6.4 Ablation for BT Pseudo-Corpus

If we follow the previous work and directly use BT pseudo-corpus in the CPT or SFT stage, how would it compare to the current Code-Switching strategy? To verify this question, we conducted a series of ablation experiments. Firstly, following Guo et al. (2024), we replaced the Code-Switching data with InnerLinear formatted data for the second-stage pretraining and also extended the training time to explore the issue of catastrophic forgetting. Next, we bypassed the second stage and utilized BT pseudo-corpus in the SFT phase, experimenting with varying amounts of data. The results are summarized in Table 4.

From the experimental results, we can draw the following conclusions:

- It is not wise to introduce BT pseudo-corpus in the SFT stage. The improvement in translation quality is not as good as that of human-written data, which aligns with previous findings.

- Using data in InnerLinear Text Format in the second stage can bring limit improvement, and there is a certain gap compared to the Code-Switching strategy in terms of BLEU and COMET metrics. Moreover, as the training time increases, the model indeed exhibits the issue of catastrophic forgetting, with a significant decline in translation quality in the ZH⇒EN direction.

7 Conclusion

In this paper, we focus on the research of adapting the translation capabilities of large models. We attempt to inject cross-lingual alignment information into LLM during the pretraining phase through a Code-Switching strategy, thereby expanding the classic two-stage training recipe. Experiments show that our Code-Switching data constructed based on the back translation strategy achieves desirable results, enhancing the end-to-end translation quality of LLMs. Additionally, we also find in our new training recipe, forward translation data seems to be more efficient, and the model’s APE capability may also benefit from the new training stage. Our Code-Switching strategy and the introduction of traditional MT models in the form of back translation into the optimization work of LLM-based translation models may inspire future research to some extent.

8 Limitations

The code-switching data format is consistent with the standard pre-training data format. In theory, we can further increase the amount of monolingual data for additional optimization. This aspect of work needs to be further explored in the future.

Current experiments and analyses are based on translation tasks centered around English. Extending our strategies to non-English translation tasks is also worth further research and optimization.

A more in-depth analysis of the principles behind the effectiveness of Code-Switching data and the internal changes in the model will lead to more meaningful discoveries.

References

- Anthropic. [claude](#).
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda

614	Askeff, et al. 2020a. Language models are few-shot learners. <i>Advances in neural information processing systems</i> , 33:1877–1901.	
615		
616		
617	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askeff, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020b. Language models are few-shot learners . <i>CoRR</i> , abs/2005.14165.	
618		
619		
620		
621		
622		
623		
624		
625		
626		
627		
628		
629	Franck Burlot and François Yvon. 2018. Using monolingual data in neural machine translation: a systematic study . In <i>Proceedings of the Third Conference on Machine Translation: Research Papers</i> , pages 144–155, Brussels, Belgium. Association for Computational Linguistics.	
630		
631		
632		
633		
634		
635	Isaac Caswell, Ciprian Chelba, and David Grangier. 2019. Tagged back-translation . In <i>Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers)</i> , pages 53–63, Florence, Italy. Association for Computational Linguistics.	
636		
637		
638		
639		
640	Yijie Chen, Yijin Liu, Fandong Meng, Yufeng Chen, Jinan Xu, and Jie Zhou. 2023. Improving translation faithfulness of large language models via augmenting instructions . <i>Preprint</i> , arXiv:2308.12674.	
641		
642		
643		
644	Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2023. Palm: Scaling language modeling with pathways . <i>J. Mach. Learn. Res.</i> , 24:240:1–240:113.	
645		
646		
647		
648		
649		
650		
651		
652		
653		
654		
655		
656		
657		
658		
659		
660		
661		
662		
663		
664		
665		
666		
667		
668	Yiming Cui, Ziqing Yang, and Xin Yao. 2023. Efficient and effective text encoding for chinese llama and alpaca . <i>arXiv preprint arXiv:2304.08177</i> .	
669		
670		
671	Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin	
672		
	Chen, Chi-Min Chan, Weize Chen, et al. 2023. Parameter-efficient fine-tuning of large-scale pre-trained language models. <i>Nature Machine Intelligence</i> , 5(3):220–235.	673 674 675 676
	Sergey Edunov, Myle Ott, Michael Auli, and David Grangier. 2018. Understanding back-translation at scale . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 489–500, Brussels, Belgium. Association for Computational Linguistics.	677 678 679 680 681 682
	Sergey Edunov, Myle Ott, Marc’Aurelio Ranzato, and Michael Auli. 2020. On the evaluation of machine translation systems trained with back-translation . <i>Preprint</i> , arXiv:1908.05204.	683 684 685 686
	Common Crawl Foundation. 2023. Statistics of common crawl monthly archives by common-crawl. https://commoncrawl.github.io/cc-crawl-statistics/plots/languages .	687 688 689 690
	Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André F. T. Martins. 2022. Results of WMT22 metrics shared task: Stop using BLEU – neural metrics are better and more robust . In <i>Proceedings of the Seventh Conference on Machine Translation (WMT)</i> , pages 46–68, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.	691 692 693 694 695 696 697 698 699
	Caglar Gulcehre, Orhan Firat, Kelvin Xu, Kyunghyun Cho, and Yoshua Bengio. 2017. On integrating a language model into neural machine translation. <i>Computer Speech & Language</i> , 45:137–148.	700 701 702 703
	Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, et al. 2023. Textbooks are all you need. <i>arXiv preprint arXiv:2306.11644</i> .	704 705 706 707 708
	Jiaxin Guo, Hao Yang, Zongyao Li, Daimeng Wei, Hengchao Shang, and Xiaoyu Chen. 2024. A novel paradigm boosting translation capabilities of large language models. <i>arXiv preprint arXiv:2403.11430</i> .	709 710 711 712
	Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are GPT models at machine translation? A comprehensive evaluation . <i>CoRR</i> , abs/2302.09210.	713 714 715 716 717 718
	Vu Cong Duy Hoang, Philipp Koehn, Gholamreza Haffari, and Trevor Cohn. 2018. Iterative back-translation for neural machine translation . In <i>Proceedings of the 2nd Workshop on Neural Machine Translation and Generation</i> , pages 18–24, Melbourne, Australia. Association for Computational Linguistics.	719 720 721 722 723 724 725
	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,	726 727

728	and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> .			
729				
730				
731	Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. <i>arXiv preprint arXiv:2310.06825</i> .			
732				
733				
734				
735				
736	Wenxiang Jiao, Jen tse Huang, Wenxuan Wang, Zhiwei He, Tian Liang, Xing Wang, Shuming Shi, and Zhaopeng Tu. 2023a. Parrot: Translating during chat using large language models tuned with human translation and feedback. <i>Preprint</i> , arXiv:2304.02426.			
737				
738				
739				
740				
741	Wenxiang Jiao, Wenxuan Wang, Jen-tse Huang, Xing Wang, and Zhaopeng Tu. 2023b. Is chatgpt a good translator? a preliminary study. <i>arXiv preprint arXiv:2301.08745</i> , 1(10).			
742				
743				
744				
745	Julia Kreutzer, Isaac Caswell, Lisa Wang, Ahsan Wahab, Daan van Esch, Nasanbayar Ulzii-Orshikh, Allahsera Tapo, Nishant Subramani, Artem Sokolov, Claytone Sikasote, Monang Setyawan, Supheakmungkol Sarin, Sokhar Samb, Benoît Sagot, Clara Rivera, Annette Rios, Isabel Papadimitriou, Salomey Osei, Pedro Ortiz Suarez, Iroro Orife, Kelechi Ogueji, Andre Niyongabo Rubungo, Toan Q. Nguyen, Mathias Müller, André Müller, Shamsuddeen Hassan Muhammad, Nanda Muhammad, Ayanda Mnyakeni, Jamshidbek Mirzakhlov, Tapiwanashe Matangira, Colin Leong, Nze Lawson, Sneha Kudugunta, Yacine Jernite, Mathias Jenny, Orhan Firat, Bonaventure F. P. Dossou, Sakhile Dlamini, Nisansa de Silva, Sakine Çabuk Ballı, Stella Biderman, Alessia Battisti, Ahmed Baruwa, Ankur Bapna, Pallavi Baljekar, Israel Abebe Azime, Ayodele Awokoya, Duygu Ataman, Orevaoghene Ahia, Oghenefego Ahia, Sweta Agrawal, and Mofetoluwa Adeyemi. 2022. Quality at a glance: An audit of web-crawled multilingual datasets. <i>Transactions of the Association for Computational Linguistics</i> , 10:50–72.			
746				
747				
748				
749				
750				
751				
752				
753				
754				
755				
756				
757				
758				
759				
760				
761				
762				
763				
764				
765				
766				
767	Jiahuan Li, Hao Zhou, Shujian Huang, Shanbo Cheng, and Jiajun Chen. 2024. Eliciting the translation ability of large language models via multilingual finetuning with translation instructions. <i>Transactions of the Association for Computational Linguistics</i> , 12:576–592.			
768				
769				
770				
771				
772				
773	Zehui Lin, Xiao Pan, Mingxuan Wang, Xipeng Qiu, Jiangtao Feng, Hao Zhou, and Lei Li. 2020. Pre-training multilingual neural machine translation by leveraging alignment information. <i>arXiv preprint arXiv:2010.03142</i> .			
774				
775				
776				
777				
778	Jean Maillard, Cynthia Gao, Elahe Kalbassi, Kaushik Ram Sadagopan, Vedanuj Goswami, Philipp Koehn, Angela Fan, and Francisco Guzmán. 2023. Small data, big impact: Leveraging minimal data for effective machine translation. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2740–2756.			
779				
780				
781				
782				
783				
784				
785				
		Bonan Min, Hayley Ross, Elior Sulem, Amir Pouran Ben Veyseh, Thien Huu Nguyen, Oscar Sainz, Eneko Agirre, Ilana Heintz, and Dan Roth. 2023. Recent advances in natural language processing via large pre-trained language models: A survey. <i>ACM Computing Surveys</i> , 56(2):1–40.		786 787 788 789 790 791
		James Cross Onur Çelebi Maha Elbayad Kenneth Heafield Kevin Heffernan Elahe Kalbassi Janice Lam Daniel Licht Jean Maillard Anna Sun Skyler Wang Guillaume Wenzek Al Youngblood Bapi Akula Loic Barrault Gabriel Mejia Gonzalez Prangthip Hansanti John Hoffman Semarley Jarrett Kaushik Ram Sadagopan Dirk Rowe Shannon Spruit Chau Tran Pierre Andrews Necip Fazil Ayan Shruti Bhosale Sergey Edunov Angela Fan Cynthia Gao Vedanuj Goswami Francisco Guzmán Philipp Koehn Alexandre Mourachko Christophe Ropers Safiyyah Saleem Holger Schwenk Jeff Wang NLLB Team, Marta R. Costa-jussà. 2022. No language left behind: Scaling human-centered machine translation.		792 793 794 795 796 797 798 799 800 801 802 803 804 805
		Pedro Javier Ortiz Su'arez, Benoit Sagot, and Laurent Romary. 2019. Asynchronous pipelines for processing huge corpora on medium to low resource infrastructures. Proceedings of the Workshop on Challenges in the Management of Large Corpora (CMLC-7) 2019. Cardiff, 22nd July 2019, pages 9 – 16, Mannheim. Leibniz-Institut f"ur Deutsche Sprache.		806 807 808 809 810 811 812
		Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>Advances in Neural Information Processing Systems</i> , 35:27730–27744.		813 814 815 816 817 818
		Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, July 6-12, 2002, Philadelphia, PA, USA</i> , pages 311–318. ACL.		819 820 821 822 823 824
		Hieu Pham, Xinyi Wang, Yiming Yang, and Graham Neubig. 2021. Meta back-translation. In <i>International Conference on Learning Representations</i> .		825 826 827
		Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. 2020. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In <i>Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining</i> , pages 3505–3506.		828 829 830 831 832 833
		Ricardo Rei, Craig Stewart, Ana C. Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020</i> , pages 2685–2702. Association for Computational Linguistics.		834 835 836 837 838 839 840
		Zhenwei Shao, Zhou Yu, Meng Wang, and Jun Yu. 2023. Prompting large language models with answer heuristics for knowledge-based visual question answering.		841 842 843

<Chinese> (千站云繁殖池是蜘蛛池的升级版, 轮链繁殖池(每一个池可以放10万链接) 2. 查一下网站不是优化过度了, 如果网站优化, 过度蜘蛛是会到网站进行屏蔽的。 <Chinese> (目录繁殖池(可日租周租月租), 联系官方qq: 3. 网站设计的结构不合理, 如果网站的页面独立, 互相没有关联也会引致收录不理想的。 <English> (Address breeding pool (Baidu, 360, Sogou, are independent pool) 4), remember to often check the content on your website is not there some sensitive words, involving some sensitive issues search engine is not included. <Chinese> <(参数繁殖池(我们引的都是有效蜘蛛) 5. 要时常分析网站内容。 <Translation> <English> (Parameter breeding pool (we are citing effective spiders)) 5. Analyze website content from time to time.> <Chinese> 如果网站上的内容都是采集来的, 或者是在其他网站上摘抄的类似内容, 蜘蛛往往是不收录的。 <English> Remember to check the stability of the server or space, and check whether there is improper code in your website program.

Figure 4: An example of Code-Switching data from Chinese to English direction.

Feng. 2023b. [Bayling: Bridging cross-lingual alignment and instruction following through interactive translation for large language models](#). *Preprint*, arXiv:2306.10968.

Tianyi Zhang, Faisal Ladhak, Esin Durmus, Percy Liang, Kathleen McKeown, and Tatsunori B Hashimoto. 2024. Benchmarking large language models for news summarization. *Transactions of the Association for Computational Linguistics*, 12:39–57.

Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. 2023. Lima: Less is more for alignment. *arXiv preprint arXiv:2305.11206*.

Kaijie Zhu, Jindong Wang, Jiaheng Zhou, Zichen Wang, Hao Chen, Yidong Wang, Linyi Yang, Wei Ye, Yue Zhang, Neil Zhenqiang Gong, and Xing Xie. 2023. [Promptbench: Towards evaluating the robustness of large language models on adversarial prompts](#). *Preprint*, arXiv:2306.04528.

A Code-Switching Data Details

The special tokens used in constructing Code-Switching data include LID and TRANS. Among them, LID consists of the language name enclosed in angled brackets, with "<Chinese>", "<English>", and "<German>" representing Chinese (LID_{zh}), English (LID_{en}), and German (LID_{de}) respectively. TRANS is "<Translation>".

An example of Code-Switching data from Chinese to English direction is shown in Figure 4.