



# Building accurate translation-tailored LLMs with language aware instruction tuning

Changtong Zan<sup>1</sup>, Liang Ding<sup>†2</sup>, Li Shen<sup>3,4</sup>, Yibing Zhan<sup>3</sup>, Xinghao Yang<sup>1</sup>, Weifeng Liu<sup>†1</sup>

<sup>1</sup>College of Control Science and Engineering, China University of Petroleum (East China), Shandong 266580, China

<sup>2</sup>School of Computer Science, University of Sydney, NSW 2006, Australia

<sup>3</sup>JD Explore Academy, JD.com Inc., Beijing 100101, China

<sup>4</sup>School of Cyber Science and Technology, Shenzhen Campus of Sun Yat-sen University, Guangdong 518107, China

<sup>†</sup>E-mail: liangding.liam@gmail.com; liuwf@upc.edu.cn

Received May 30, 2024; Revision accepted Nov. 27, 2024; Crosschecked

**Abstract:** Large language models (LLMs) exhibit remarkable capabilities in various natural language processing tasks, such as machine translation. However, the large number of LLM parameters incurs significant costs during inference. Previous studies have attempted to train translation-tailored LLMs with moderately sized models by fine-tuning them on translation data. Nevertheless, when applying zero-shot translation directions not included in the fine-tuning data, the issue of ignoring instructions and thus translating into the wrong language, i.e., the off-target translation issue, remains unsolved. In this work, we design a two-stage fine-tuning algorithm to improve the instruction-following ability of translation-tailored LLMs, particularly for maintaining accurate translation directions. We first fine-tune LLMs on the translation dataset to elicit basic translation capabilities. In the second stage, we construct instruction-conflicting samples by randomly replacing the instructions with incorrect ones. Then, we introduce an extra unlikelihood loss to reduce the probability assigned to those samples. Experiments on IWSLT and WMT benchmarks using the LLaMA2 and LLaMA3 models, spanning 16 zero-shot directions, demonstrate that, compared to the competitive baseline – translation-finetuned LLaMA, our method could effectively reduce the off-target translation ratio (up to  $-62.4\%$ ), thus improving translation quality (up to  $+9.7$  BLEU). Analysis shows that our method can preserve the model's performance on other tasks, such as supervised translation and general tasks. Code is released at: [https://github.com/alphadl/LanguageAware\\_Tuning](https://github.com/alphadl/LanguageAware_Tuning).

**Key words:** Zero-shot machine translation; Off-target issue; Large language model; Language-aware instruction tuning; Instruction-conflicting sample

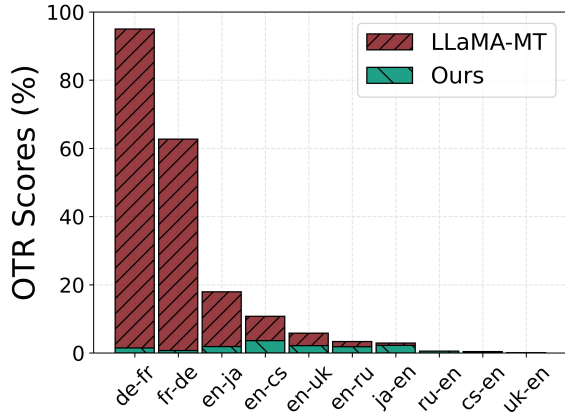
<https://doi.org/10.1631/FITEE.2400458>

**CLC number:** TP

## 1 Introduction

Large language models (LLMs) (Min et al., 2023; Li et al., 2024a) have demonstrated excellent performance on a wide range of natural language processing (NLP) tasks, including reasoning (Wei et al., 2022), summarization (Wang et al., 2023a; Huang et al., 2024), translation (Hendy et al., 2023), understanding (Zhang et al., 2022; Zhong et al., 2023), and evaluation (Lu et al., 2023). LLMs exemplified by GPT-3 (Brown et al., 2020), OPT (Zhang et al., 2022),

LLaMA (Touvron et al., 2023a), and LLaMA2 (Touvron et al., 2023b), leverage large-scale monolingual data through pre-training with the causal language modeling task and exhibit strong zero-shot capabilities with few demonstration examples. Instruction tuning (Wei et al., 2021; Mishra et al., 2022) further elicits the capacity of LLMs to address general tasks directly with proper guidance, such as task definition. However, due to the significant cost to call the state-of-the-art proprietary LLMs, such as GPT-4 (OpenAI, 2024), it is attractive to explore



**Fig. 1** OTR (%) ↓ in ZST of WMT dataset. We present a comparison between LLaMA2-MT, the LLaMA2 model fine-tuned on translation data, and our model. OTR: Off-target translation ratio; ZST: Zero-shot translation.

strategies for effectively fitting suitably sized LLMs into specific tasks, such as machine translation (Fu et al., 2023; ?).

In zero-shot translation (ZST) (Gu et al., 2019; Chen et al., 2023; Zan et al., 2023), the objective is to translate sentences from a source language to a target language, where there is either a lack of direct mapping between source and target languages in the training data, or the target or source languages themselves are absent during training. Addressing the ZST problem is both vital and challenging, especially for low-resource languages with limited data resources. Recent research demonstrate that building translation-tailored LLMs by fine-tuning translation data can achieve superior translation performance (Zeng et al., 2023; Liu et al., 2023; ?). However, as illustrated in Fig. 1, our preliminary study shows that, when tackling zero-shot directions, translation-tailored LLMs often encounter the off-target translation problem, where the generated translations are in the wrong languages. For example, in De→Fr, the off-target ratio reaches up to 95.0%. We attribute this problem to the fact that pre-training LLMs in the fashion of predicting the next token may lead to overlooking the key information contained in instructions.

Previous studies (Peng et al., 2023; ?) indicate that introducing more informative prompts during inference, such as preemptively translating prompts into the target language or incorporating few-shot demonstration samples, can be beneficial. Sennrich

et al. (2024) modified the decoding process by introducing language-contrastive samples to constrain the decoding process, thus alleviating the off-target problem. Different from the above approaches that focus on the inference stage, our motivation is to fundamentally improve the instruction-following ability (especially the awareness of translation direction) of LLMs themselves.

In this paper, we introduce a simple and effective two-stage fine-tuning algorithm to enhance the effect of instructions in translation-tailored LLMs. This is accomplished by introducing unlikelihood loss on instruction-conflicting samples in which the translation sentence pairs deviate from the specified tasks associated with the given instructions. Initially, we fine-tune the LLMs using a multilingual translation dataset, unlocking the inherent translation capabilities of LLMs. Subsequently, we build upon the pre-tuned model by incorporating translation data along with instruction-conflicting samples. We create instruction-conflicting samples by randomly replacing the task instruction with an incorrect one. These data are used to fine-tune the model, leveraging the unlikelihood training paradigm. Our approach can be viewed as emphasizing the effect of instructions, thereby guiding the model to produce translation in the correct language.

We apply our method in the experiments to fine-tune the currently competitive open-source LLaMA2 and LLaMA3 models. Compared with direct fine-tuning on translation data, the results reveal substantial reductions in the off-target translation ratio (OTR), with 62.3% and 30.5% for the IWSLT benchmark and 29.9% and 18.6% for the WMT benchmark, respectively, for each model. This leads to notable enhancements in translation quality, as evidenced by increases of average +9.7/ +6.2 and +6.2/ +3.9 BLEU scores in the IWSLT and WMT datasets. Also, our method maintains the capability to perform other tasks, such as supervised directions and general tasks. The main contributions are as follows:

- We reveal the heavy off-target problem of translation-tailored LLMs in ZST settings, and we attribute this problem to the weak instruction (translation direction) following ability.
- To fundamentally improve the translation direction-following ability, we introduce a two-stage fine-tuning algorithm for LLMs that lever-

ages instruction-conflicting samples.

- Extensive experiments illustrate the effectiveness of our approach in mitigating the off-target translation problem and producing better translations. Analyses show that our method will not affect the ability of other tasks, e.g., general task performance on AlpacaEval and supervised translation performance.

## 2 Preliminary

**Instruction tuning** Instruction tuning aims to refine LLMs by fine-tuning a diverse collection of data characterized by explicit instructions. This refinement process significantly enhances zero-shot performance on previously unseen tasks (Wei et al., 2021). Each instance in the instruction tuning dataset comprises three fundamental components: (1) *Instruction*: A textual representation that describes NLP tasks in natural language. (2) *Input* (optional): Supplementary contextual information that provides additional context for the given task. (3) *Output*: The expected response that LLMs should generate. During the tuning process, the model is trained using a teacher-forcing approach (Cho et al., 2014). It models the distribution of output tokens conditioned on the instruction and, optionally, the input. This training methodology empowers the model to understand and follow instructions effectively. Subsequently, the instruction-tuned model is capable of directly performing unseen tasks by following the appropriate task instructions in a zero-shot manner. In this study, our primary focus is translation-tailored LLMs, where we fine-tune LLMs on multilingual translation data.

**Unlikelihood training** Welleck et al. (2020) explored a novel approach that encourages the model to assign lower probabilities to improbable generations, in contrast to the traditional likelihood training, which focuses on the overall probability distribution of correct sequences. The general training framework comprises two types of updates:

1. Likelihood update: Maximizing the probability of the ground-truth sequence with likelihood loss:

$$\mathcal{L}_{MLE} = - \sum_{t_n \in y} \log P(t_n | t_0, \dots, t_{n-1})$$

where  $t_n$  is the ground-truth token at position  $n$  and  $y$  is the ground-truth sequence.

2. Unlikelihood update: Minimizing the probability of negative samples (tokens that should not occur) with unlikelihood loss:

$$\mathcal{L}_{UL} = - \sum_{t_n \in y} \sum_{k \in \mathcal{C}_n} \log(1 - P(t_n | t_0, \dots, t_{n-1}))$$

where  $\mathcal{C}_n$  is the set of negative candidate tokens at position  $n$ . We extend this approach to the domain of ZST based on translation-tailored LLMs. We introduce instruction-conflicting samples for unlikelihood updates, emphasizing the impact of translation instructions (especially the translation direction and language) and addressing off-target problems.

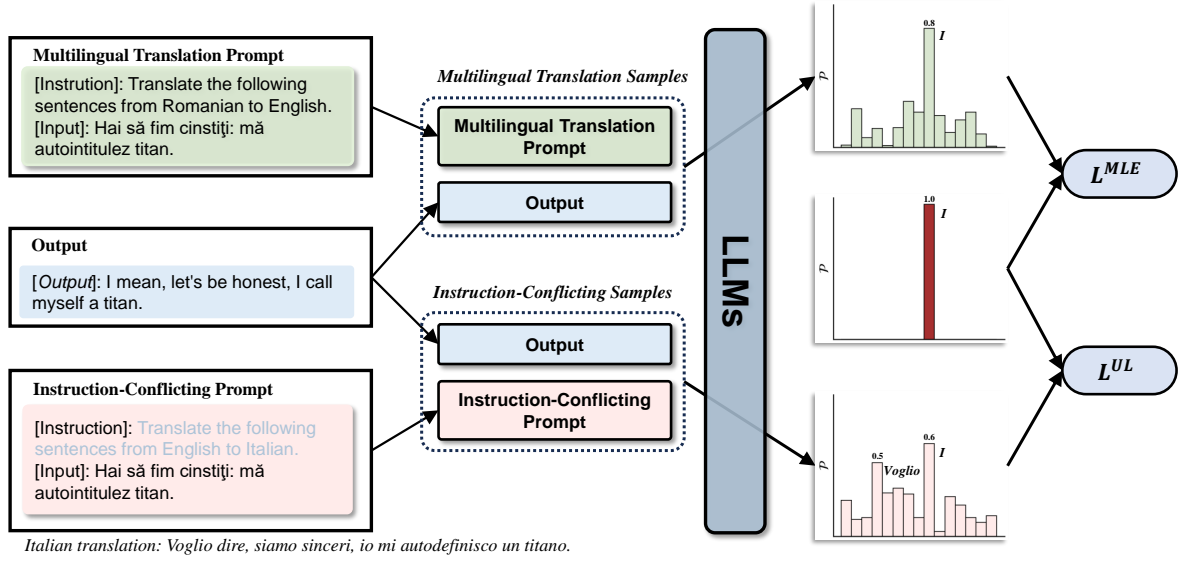
## 3 Methodology

### 3.1 Pre-tuning on multilingual translation samples

To unlock the translation capabilities of LLMs, we employ a pre-tuning stage using multilingual translation examples. Given a collection of instruction samples  $\mathcal{D} = \{\mathcal{D}_1, \dots, \mathcal{D}_i, \dots, \mathcal{D}_N\}$  covering  $N$  language pairs, where  $\mathcal{D}_i$  denotes a translation parallel corpus of the  $i$ -th language pair. For each sample  $\mathcal{S}_j^i = \{\mathit{ins}_j^i, \mathbf{x}_j^i, \mathbf{y}_j^i\}$  in  $\mathcal{D}_i$ , the model is trained with the common approach, i.e. maximum likelihood estimation (MLE). As depicted in Fig. 2, the instruction *ins* for the multilingual translation sample is “Translate the following sentences from Romanian to English,” and the corresponding input source sentence  $\mathbf{x}$  is “Hai să fim cinstiți: mă autointitulez titan.” We train the model to model the mapping from both *ins* and  $\mathbf{x}$  to the expected translation output  $\mathbf{y}$ , i.e., “I mean, let’s be honest, I call myself a titan.” Formally, the training loss for each sample  $\mathcal{S}_j^i$  is as follows:

$$\mathcal{L}_{\mathcal{S}_j^i}^{\text{MLE}}(\theta) = - \sum_{t_n \in \mathbf{y}_j^i} \log P(t_n | \mathit{ins}_j^i, \mathbf{x}_j^i, t_0, \dots, t_{n-1}; \theta),$$

where  $t_n$  is the  $n$ -th token of the translation sequence  $y$ , such as the first token of  $\mathbf{y}$  in the above example is “I.”  $\theta$  represents the trainable model parameters. Minimizing the loss of each token in the translation output unlocks the ability of LLMs to perform translation tasks by following the provided instructions. Then, the final objective of our first fine-tuning stage, pre-tuning on the multilingual translation dataset,



**Fig. 2** Overview of our fine-tuning framework for ZST. (a) In the first stage, we pre-tune LLMs on multilingual translation samples, focusing on unlocking the translation ability of LLMs. (b) Subsequently, we introduce instruction-conflicting samples by randomly substituting the instruction component with a different one. We then train the model with  $\mathcal{L}^{MLE}$  on translation data and incorporate an unlikelihood loss  $\mathcal{L}^{UL}$  on the instruction-conflicting samples to assign lower probabilities to wrong language tokens. LLMs: Large language models; ZST: Zero-shot translation.

can be formulated as follows:

$$\mathcal{L}_{\mathcal{D}}^{MLE}(\theta) = - \sum_{\mathcal{D}_i \in \mathcal{D}} \sum_{\mathcal{S}_j^i \in \mathcal{D}_i} \mathcal{L}_{\mathcal{S}_j^i}^{MLE}(\theta).$$

It is worth noting that the instruction  $ins$  is not limited to translation task definition. Samples from diverse NLP tasks can be formulated into this unified format with proper instruction.

### 3.2 Unlikelihood training with instruction-conflicting samples

In the second stage, our objective is to improve the LLM's instruction-following ability and enhance the ZST capability. We employ a dual optimization approach, which involves training the model with likelihood loss on multilingual translation samples and unlikelihood loss on samples with conflicting instructions.

**Instruction-conflicting samples** Although LLMs can achieve impressive performance on tasks within the multilingual translation tuning dataset, they may encounter issues such as ignoring instructions and generating translations in the wrong language during ZST, commonly referred to as the off-target problem. To address the off-target issue

with unlikelihood training, we define the negative candidate samples by substituting the instruction with a different one while keeping the input and output unchanged. We call this type of samples *instruction-conflicting samples* as the translation pairs deviate from the task associated with the given instructions. As shown in Fig. 2, given a multilingual translation sample  $\mathcal{S}_j^i = \{ins_j^i, x_j^i, y_j^i\}$  from the instruction tuning dataset  $\mathcal{D}$ , we randomly select a sample with different instruction  $\widetilde{ins}_j^i$ , e.g., “Translate the following sentences from English to Italian.” Then, we replace the original correct  $ins_j^i$  to get the instruction-conflicting sample  $\widetilde{\mathcal{S}}_j^i = \{\widetilde{ins}_j^i, x_j^i, y_j^i\}$ , e.g. “[**Instruction**]: Translate the following sentences from English to Italian. [**Input**]: Hai să fim cinstiți: mă autointitulez titan. [**Output**]: I mean, let’s be honest, I call myself a titan.” in the example. The instruction-conflicting samples contain incorrect outputs and do not align with the task defined by the instruction.

**Unlikelihood training with instruction-conflicting samples** Unlikelihood training aims to reduce the probability assigned by the model to negative candidate tokens. Based on the previously defined instruction-conflicting samples, we can extend the

unlikelihood training to translation-tailored LLMs and enhance their ability to follow instructions for translation tasks.

In each update, we optimize the unlikelihood loss for instruction-conflicting samples  $\widetilde{\mathcal{S}}_j^i$ :

$$\mathcal{L}_{\widetilde{\mathcal{S}}_j^i}^{\text{UL}}(\theta) = - \sum_{t_n \in \widetilde{\mathbf{y}}_j^i} \log(1 - P(t_n | \widetilde{\mathbf{ins}}_j^i, \mathbf{x}_j^i, \mathbf{t}_0, \dots, \mathbf{t}_{n-1}; \theta)),$$

where  $t_n$  represents the  $n$ -th token of output. Due to the relatively small difference between  $\widetilde{\mathbf{ins}}$  and  $\mathbf{ins}$ , such as only two words being different in the example, the model pre-tuned on multilingual translation samples usually tends to assign high probability on the output. As shown in Fig. 2, the first token ‘‘P’’ has a high probability, and the model is more confident to generate ‘‘P’’ rather than the correct word ‘‘Voglio.’’ The objective of unlikelihood loss on the whole dataset is as follows:

$$\mathcal{L}_{\mathcal{D}}^{\text{UL}}(\theta) = - \sum_{\mathcal{D}_i \in \mathcal{D}} \sum_{\widetilde{\mathcal{S}}_j^i \in \mathcal{D}_i} \mathcal{L}_{\widetilde{\mathcal{S}}_j^i}^{\text{UL}}(\theta).$$

To prevent the potential overfitting on the unlikelihood objective while maintaining the supervised translation ability, we incorporate multilingual translation samples to simultaneously train the model with likelihood loss. The overall objective function in unlikelihood training involves a mixture of likelihood and unlikelihood loss, defined as:

$$\mathcal{L}_{\mathcal{D}}(\theta) = \mathcal{L}_{\mathcal{D}}^{\text{MLE}}(\theta) + \alpha \mathcal{L}_{\mathcal{D}}^{\text{UL}}(\theta),$$

where  $\alpha$  is the mixing hyper-parameter.

## 4 Experiments

In this section, we conduct a series of experiments spanning 16 ZST directions to assess the effectiveness of our algorithm.

### 4.1 Experimental setup

**Datasets** We consider the following two widely used datasets:

- **WMT**: Following Jiao et al. (2023) and Liu et al. (2023), we use the development sets from WMT2017 to WMT2020 for instruction tuning. This includes four language directions: En $\leftrightarrow$ Zh

and En $\leftrightarrow$ De, encompassing a total of 51k translation sentence pairs. Then, we assess translation performance on WMT22 test sets, including En $\leftrightarrow$ Cs, En $\leftrightarrow$ Ja, En $\leftrightarrow$ Ru, En $\leftrightarrow$ Uk, and Fr $\leftrightarrow$ De. All these translation language pairs do not exist in the training set, thus allowing for the evaluation of ZST performance.

- **IWSLT**: Following Qu and Watanabe (2022), we use the IWSLT-17 dataset to evaluate the performance of the models. We consider the four languages (‘‘En, Ro, It, Nl’’) from MMCR4NLP (Dabre and Kurohashi, 2019). For training, we randomly select 12k sentence pairs from the training set, spanning six directions: En $\leftrightarrow$ Nl, En $\leftrightarrow$ It, En $\leftrightarrow$ Ro. The evaluation is conducted on the test set of IWSLT-17, including Ro $\leftrightarrow$ Nl, Ro $\leftrightarrow$ It, Nl $\leftrightarrow$ It. All these translation language pairs do not exist in the training set.

**Base models** We employ both the 7B size LLaMA2 (Touvron et al., 2023b) and the 8B size LLaMA3<sup>1</sup> as the base models. LLaMA2 is a robust language model that has undergone training on 2 trillion tokens, with less than 2% non-English data. In contrast, LLaMA3 represents a significant advancement, having been pretrained on more than 7 $\times$  pretrain data used for LLaMA2, specifically 15 trillion tokens. This consists of more multilingual data (over 5% non-English data). Moreover, LLaMA3 has a larger vocabulary than LLaMA2. These setups ensure that LLaMA3 models have better capability to process multilingual tasks.

**Baselines** We consider the following baselines:

- **LLaMA**: We employ the pre-trained LLMs directly for inference. The prompt is the same as those used by the below MT.
- **MT**: Following previous works (Jiao et al., 2023; ?; Zeng et al., 2023), we fine-tune the LLMs on multilingual translation samples, establishing this configuration as our main baseline. We format translation sentence pairs into a unified translation template.
- **Post-ins**: Following Liu et al. (2023), we switch the positions of instruction and input, where the

<sup>1</sup><https://llama.meta.com/llama3/>

Base model	Methods	Cs-En		Ja-En		Ru-En		Uk-En		Fr-De		AVG
		←	→	←	→	←	→	←	→	←	→	
<i>BLEU score</i> ↑												
LLaMA2	LLaMA	0.2	1.3	0.1	0.4	0.3	1.3	0.2	1.9	0.8	0.6	<u>0.7</u>
	MT	18.2	<b>39.0</b>	9.6	16.1	21.2	36.4	9.6	<b>34.3</b>	4.3	3.2	<u>19.2</u>
	Post-ins	18.8	38.0	12.4	15.8	<b>22.1</b>	<b>36.5</b>	14.6	34.1	<b>30.7</b>	5.3	<u>22.8</u>
	PTL	17.6	<b>39.0</b>	10.6	16.1	20.0	36.4	17.8	<b>34.3</b>	24.7	<b>24.6</b>	<u>24.1</u>
	1-shot	18.8	37.2	11.4	15.5	20.9	34.9	17.7	32.9	3.9	3.2	<u>19.6</u>
	5-shot	18.3	37.0	12.2	15.1	20.9	34.2	<b>18.4</b>	31.8	3.7	3.2	<u>19.5</u>
	$\mathcal{C}_{src+lang}$	3.6	35.5	2.6	13.1	3.2	33.9	2.1	31.8	1.4	0.7	<u>12.8</u>
	Ours	<b>18.8</b>	38.9	<b>12.8</b>	<b>16.3</b>	20.9	35.8	18.0	32.6	29.8	19.4	<u>24.3</u>
<i>OTR score (%)</i> ↓												
LLaMA2	LLaMA	90.0	24.2	98.7	16.7	86.6	20.0	90.8	14.9	84.0	85.6	<u>61.1</u>
	MT	10.7	0.3	27.9	2.1	6.7	<b>0.3</b>	60.4	0.1	90.8	99.3	<u>29.9</u>
	Post-ins	6.3	1.9	9.7	2.7	2.5	0.4	25.1	<b>0.0</b>	<b>4.0</b>	87.4	<u>14.0</u>
	PTL	18.5	<b>0.3</b>	16.0	2.1	13.1	<b>0.3</b>	6.2	0.1	19.7	<b>11.4</b>	<u>8.8</u>
	1-shot	10.1	0.4	12.6	4.0	6.5	0.5	7.6	0.1	97.7	99.4	<u>23.9</u>
	5-shot	13.8	0.6	11.8	4.9	8.3	0.4	7.1	0.1	98.9	99.6	<u>24.6</u>
	$\mathcal{C}_{src+lang}$	<b>3.1</b>	<b>0.3</b>	19.2	<b>1.4</b>	<b>2.4</b>	0.5	12.5	<b>0.0</b>	87.5	97.9	<u>22.5</u>
	Ours	6.6	<b>0.3</b>	<b>2.3</b>	1.9	2.6	<b>0.3</b>	<b>3.2</b>	0.1	6.7	31.7	<u>5.6</u>
<i>BLEU score</i> ↑												
LLaMA3	LLaMA	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	<u>0.0</u>
	MT	20.6	<b>39.7</b>	11.5	17.3	<b>25.2</b>	<b>35.9</b>	19.7	<b>34.2</b>	15.3	4.8	<u>22.4</u>
	Post-ins	21.7	<b>39.7</b>	12.0	16.4	24.9	<b>35.9</b>	18.6	34.1	<b>33.4</b>	11.8	<u>24.8</u>
	PTL	20.6	<b>39.7</b>	5.8	17.3	24.6	<b>35.9</b>	19.6	<b>34.2</b>	15.4	10.1	<u>22.3</u>
	1-shot	21.5	38.4	11.7	15.4	24.6	34.0	19.2	32.7	19.1	5.6	<u>22.2</u>
	5-shot	22.0	38.0	13.1	14.9	24.1	33.7	<b>19.8</b>	32.4	24.6	8.4	<u>23.1</u>
	$\mathcal{C}_{src+lang}$	3.0	35.1	2.2	14.5	2.6	31.6	3.0	29.4	1.4	1.7	<u>12.4</u>
	Ours	<b>23.5</b>	39.6	<b>15.6</b>	<b>17.7</b>	24.9	34.8	19.6	33.3	31.0	<b>27.1</b>	<u>26.7</u>
<i>OTR score (%)</i> ↓												
LLaMA3	LLaMA	100.0	0.1	100.0	0.2	100.0	0.0	100.0	0.0	99.2	99.3	<u>59.9</u>
	MT	10.8	0.4	17.9	2.9	3.3	0.5	5.8	<b>0.0</b>	62.7	95.0	<u>19.9</u>
	Post-ins	8.8	0.3	20.6	3.4	2.8	0.5	8.8	<b>0.0</b>	1.3	62.7	<u>10.9</u>
	PTL	12.8	0.4	58.7	2.9	4.0	0.5	5.8	<b>0.0</b>	60.2	79.5	<u>22.5</u>
	1-shot	10.8	0.4	10.6	3.1	3.5	0.5	4.8	0.1	48.6	92.6	<u>17.5</u>
	5-shot	11.1	0.8	10.2	3.3	4.2	0.5	6.3	0.2	27.6	79.9	<u>14.4</u>
	$\mathcal{C}_{src+lang}$	<b>2.2</b>	<b>0.1</b>	8.1	<b>1.8</b>	<b>0.5</b>	<b>0.3</b>	<b>0.8</b>	<b>0.0</b>	24.4	71.3	<u>10.9</u>
	Ours	3.6	0.3	<b>1.9</b>	2.2	1.8	0.5	2.2	0.1	<b>0.7</b>	<b>1.5</b>	<u>1.5</u>

Table 1 ZST performance achieved on WMT benchmark. **Bold:** The better results, except OTR scores of LLaMA. Underline: Average scores obtained for all directions. OTR: Off-target translation ratio; PTL: Prompt in the target language; ZST: Zero-shot translation.

model pays more attention to the instruction. Other settings remain consistent with MT.

- **Prompt in the target language (PTL):** Instead of using the English prompt during inference, we translate the prompt into the target language during inference, which could potentially provide more guidance information for generating target language words. This setup uses the MT models for inference

- **$K$ -shot:** In-context learning (Brown et al., 2020) has proven to be an effective way to improve the performance of LLMs. We report the few-shot performance for comprehensive comparison, including 1-shot and 5-shot. MT models are used for inference.

- $\mathcal{C}_{lang}$ : Following Sennrich et al. (2024), we employ the decoding method by contrasting the translation sentence with language-contrastive

input and  $\lambda_{lang}$  0.5. The inference relies on MT models and uses a greedy decoding strategy.

**Model training** We conduct experiments on the **Huggingface Transformers** (Wolf et al., 2020) toolkit. During the pre-tuning phase, we set the learning rate (lr) to  $2e-5$ , the warmup ratio to 0.03, and the batch size to 128. For the IWSLT dataset, we performed training over 3 epochs, while for the WMT dataset, training was conducted for 1 epoch. During the second stage of training, we set the mixing parameter denoted as  $\alpha$  to 0.05, the lr to  $2e-6$ , the batch size to 8, and the training step to 100. We use the final model for evaluation.

**Evaluation** We adopt **SacreBLEU** (Post, 2018) to evaluate the translation accuracy, where translations are generated with a beam size of 4. Besides, we compute the ratio of wrong language translation in the generated outputs, i.e., **OTR**, with a publicly available language detector<sup>2</sup> (Joulin et al., 2016a,b). We utilize vLLM (Kwon et al., 2023) to accelerate during inference. Unless otherwise stated, we report results in the ZST directions.

## 4.2 Main results

We compare the ZST performance of our models and other baseline methods on the WMT and IWSLT benchmarks, as depicted in Tables 1 and 2. Obviously, our models generate more translations in the expected languages with the lowest average OTR scores across 16 directions for both base models. Especially notable is the performance of our method with the LLaMA2 base model, which demonstrates a significant average OTR gain of  $-62.4\%$  in the IWSLT benchmark compared with MT. We attribute this improvement to the better translation direction instruction following the ability of our models. The more the model tends to omit the instruction of language, the greater the benefit from our method. It is worth noting that LLaMA3 achieves 0 BLEU scores in our experiments because it generates a terminator as the first token, which may be attributed to its specific pre-training settings. Furthermore, our models outperform baseline approaches that mitigate off-target problems during inference, such as

<sup>2</sup><https://fasttext.cc/docs/en/language-identification.html>

PTL,  $K$ -shot, and  $\mathcal{C}_{lang}$ . Our method achieves improvements of up to  $+14.3/ +14.3$  average BLEU scores and  $-21.0\%/ -63.3\%$  average OTR scores in the WMT/ IWSLT datasets compared with these methods. Regarding baseline adjustments during the tuning stage, our model achieves improvements over Post-ins, up to  $+1.6/ +6.9$  average BLEU and  $-9.5\%/ -55.6\%$  average OTR in the WMT/ IWSLT datasets. In IWSLT, we also observe that our model shows a slightly lower BLEU score than Post-ins (15.9 vs. 16.1 average BLEU) when LLaMA3 is used as the base model. Additionally, our model surpasses other robust baseline models in these evaluations.

## 5 Analysis

To provide deeper insight into the proposed algorithm, we analyze to investigate the following issues: (1) Can increasing the number of training steps lead to greater benefits? (2) What is the impact of the mixing hyperparameter? (3) How does the size of the translation data in pre-tuning affect the outcome? (4) What role does the model size play? (5) Can continuing training with instruction-conflicting samples maintain supervised translation performance? (6) What is the source of the observed improvements? (7) If we introduce additional general task tuning data, can the model still retain its general task ability?

### 5.1 Effect of unlikelihood training steps

To provide insight into the impact of unlikelihood training steps in our second fine-tuning stage, Fig. 3a presents the ZST performance of our LLaMA2-based models on the IWSLT dataset. As observed, the model produces fewer wrong language translations and higher quality translations with more unlikelihood training steps. From the figure, it can be seen that the model achieves its best performance, denoted by near-zero OTR scores, after about 40 updates, and this performance is consistently maintained even with further training extending up to 100 steps. This shows that **our method can improve the ZST performance without a significant number of updates.**

Base model	Methods	It→Nl	Nl→It	It→Ro	Ro→It	Nl→Ro	Ro→Nl	AVG
<i>BLEU score</i> ↑								
LLaMA2	LLaMA	0.9	0.5	0.5	1.0	0.3	0.7	<u>0.7</u>
	MT	7.9	8.3	2.3	4.8	2.8	4.2	<u>5.0</u>
	Post-ins	9.0	11.2	5.6	7.7	8.1	5.1	<u>7.8</u>
	PTL	10.2	9.4	6.8	7.8	5.3	6.6	<u>7.7</u>
	1-shot	11.5	10.8	3.0	8.5	3.0	6.9	<u>7.3</u>
	5-shot	8.5	9.4	1.7	6.8	1.4	4.5	<u>5.4</u>
	$\mathcal{C}_{src+lang}$	2.3	2.0	0.7	1.7	0.7	1.8	<u>1.5</u>
	Ours	<b>17.5</b>	<b>16.2</b>	<b>15.4</b>	<b>12.8</b>	<b>12.0</b>	<b>14.5</b>	<b><u>14.7</u></b>
<i>OTR score (%)</i> ↓								
LLaMA2	LLaMA	86.2	85.5	91.3	84.3	94.7	88.7	<u>88.5</u>
	MT	49.8	39.8	85.8	65.7	80.1	68.8	<u>65.0</u>
	Post-ins	56.8	32.4	71.3	64.7	46.2	77.7	<u>58.2</u>
	PTL	34.4	34.2	49.7	50.2	60.0	52.3	<u>46.8</u>
	1-shot	32.0	26.7	82.2	47.1	81.8	50.1	<u>53.3</u>
	5-shot	46.2	34.4	95.1	57.5	92.3	69.8	<u>65.9</u>
	$\mathcal{C}_{src+lang}$	31.8	46.0	85.4	65.4	82.8	56.0	<u>61.2</u>
	Ours	<b>2.7</b>	<b>1.5</b>	<b>3.8</b>	<b>1.5</b>	<b>3.8</b>	<b>2.5</b>	<b><u>2.6</u></b>
<i>BLEU score</i> ↑								
LLaMA3	LLaMA	0.0	0.0	0.0	0.0	0.1	0.0	<u>0.0</u>
	MT	12.0	15.1	3.5	15.9	3.1	8.2	<u>9.6</u>
	Post-ins	17.7	<b>17.9</b>	13.5	<b>19.9</b>	<b>12.1</b>	<b>15.6</b>	<b><u>16.1</u></b>
	PTL	14.5	17.3	8.7	18.4	8.9	10.4	<u>13.0</u>
	1-shot	15.4	17.3	4.0	19.3	8.3	12.6	<u>12.8</u>
	5-shot	15.0	17.8	1.3	19.3	9.6	14.3	<u>12.9</u>
	$\mathcal{C}_{src+lang}$	1.3	1.5	1.2	2.3	1.1	2.5	<u>1.6</u>
	Ours	<b>18.2</b>	<b>17.9</b>	<b>15.6</b>	18.4	10.5	14.7	<u>15.9</u>
<i>OTR score (%)</i> ↓								
LLaMA3	LLaMA	99.0	99.2	98.9	98.1	99.9	99.7	<u>99.2</u>
	MT	27.4	4.5	80.7	5.7	53.3	30.4	<u>33.6</u>
	Post-ins	5.1	3.6	13.2	4.0	11.7	5.4	<u>7.2</u>
	PTL	19.4	7.4	45.3	5.9	32.1	25.6	<u>22.6</u>
	1-shot	16.3	4.9	79.4	5.7	34.1	19.9	<u>26.7</u>
	5-shot	16.8	3.5	97.5	3.5	32.5	13.9	<u>27.9</u>
	$\mathcal{C}_{src+lang}$	21.7	2.4	49.9	4.7	34.2	25.6	<u>23.1</u>
	Ours	<b>2.5</b>	<b>1.5</b>	<b>4.8</b>	<b>1.7</b>	<b>6.4</b>	<b>2.1</b>	<b><u>3.2</u></b>

Table 2 ZST performance on the IWSLT dataset. **Bold**: The best results. Underline: Average scores obtained for all directions. OTR: Off-target translation ratio; PTL: Prompt in the target language; ZST: Zero-shot translation.

## 5.2 Effect of $\alpha$

As mentioned in Section 3.2, our algorithm incorporates a mixing hyper-parameter  $\alpha$  to balance MLE loss and UL loss. This is an ablation to evaluate the effect of different  $\alpha$ . Fig. 3b shows the ablation results of our LLaMA2-based models on the IWSLT dataset. As expected, the higher  $\alpha$  highlights the UL loss, resulting in fewer wrong language translations.

Models fine-tuned with  $\alpha$  exceeding 0.02 are unlikely to produce translations in the wrong languages. However, when  $\alpha$  increases, there is a slight decrease in translation quality as reflected by the BLEU score. This decline may be attributed to potential overfitting on the unlikelihood loss. Future research efforts could aim to alleviate the effects of this potential overfitting issue. In summary, our experimental results indicate



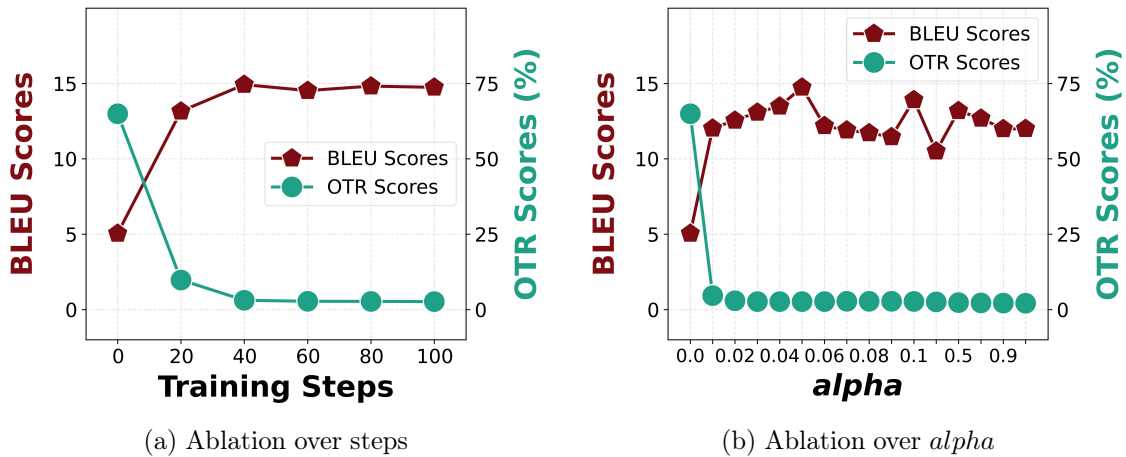


Fig. 3 Ablation studies. (a) Ablation study on unlikelihood training steps. (b) Ablation study on the mixing hyper-parameter  $\alpha$ .

Models	Size	It→Nl	Nl→It	It→Ro	Ro→It	Nl→Ro	Ro→Nl	AVG
<i>BLEU Score</i> ↑								
LLaMA2-MT	7B	7.9	8.3	2.3	4.8	2.8	4.2	<u>5.0</u>
	13B	11.2	9.2	7.0	6.7	4.8	7.6	<u>7.7</u>
Ours	7B	<b>17.5</b>	<b>16.2</b>	<b>15.4</b>	<b>12.8</b>	<b>12.0</b>	<b>14.5</b>	<u><b>14.7</b></u>
	13B	<b>15.5</b>	<b>18.8</b>	<b>13.8</b>	<b>20.4</b>	<b>12.9</b>	<b>17.4</b>	<u><b>16.5</b></u>
<i>OTR Score (%)</i> ↓								
LLaMA2-MT	7B	49.8	39.8	85.8	65.7	80.1	68.8	<u>65.0</u>
	13B	38.9	55.1	57.3	69.0	67.1	59.7	<u>57.9</u>
Ours	7B	<b>2.7</b>	<b>1.5</b>	<b>3.8</b>	<b>1.5</b>	<b>3.8</b>	<b>2.5</b>	<u><b>2.6</b></u>
	13B	<b>3.4</b>	<b>1.1</b>	<b>4.6</b>	<b>1.2</b>	<b>4.3</b>	<b>2.7</b>	<u><b>2.9</b></u>

Table 3 The impact of model size. We report the BLEU and OTR scores on the IWSLT dataset. Bold: The better results. Underline: Average scores obtained for all directions. OTR: Off-target translation ratio.

that our method exhibits robustness to varying values of the mixing parameter,  $\alpha$ .

### 5.3 Results with different size of LLMs

To examine the influence of model size, we fine-tune the 13B size LLaMA2 on the IWSLT dataset, employing the same experimental setup as the 7B size model. We report the results of MT and our models, which are summarized in Table 3. The 13B model consistently outperforms the 7B model in terms of both reducing the wrong language translation ratio ( $-7.1\%$  average OTR score) and improving translation quality ( $+2.7$  average BLEU score). This finding is consistent with prior research (Kaplan et al., 2020), which suggests that increasing the number of training parameters yields benefits. However, the off-target problem still exists in the 13B size LLaMA2-MT model. And, our model achieves a significantly lower

OTR (2.9% vs. 57.9% average OTR score), leading to a higher quality translation (16.5 vs. 7.7 average BLEU score). This result demonstrates that our algorithm remains effective with larger LLMs.

### 5.4 Results with different amounts of translation data

Fig. 4 illustrates the ZST performance of LLaMA2-based models, as measured by BLEU and OTR scores, with multilingual translation datasets of varying sizes, denoted as  $n$  samples. Four settings are considered, with dataset sizes of 12k, 24k, 48k, and 96k samples. Although the focus is on ZST performance, increasing the translation data size also yields improvements. For LLaMA-MT models, tuning with 96k samples achieves the best performance (14.2 BLEU score and 24.7% OTR score). And, our model attains its peak performance with 96k training

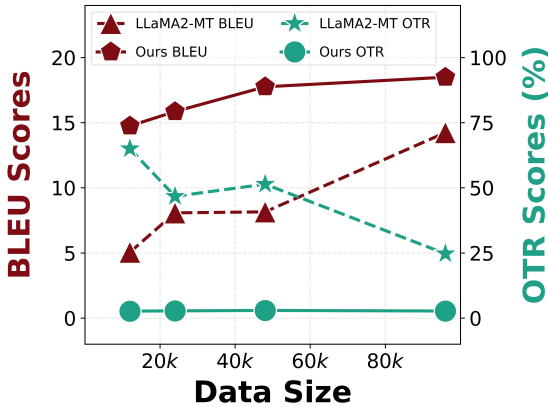


Fig. 4 The impact of fine-tuning translation data size. We report the BLEU and OTR scores on the IWSLT dataset. The x-axis represents the fine-tuning data size, denoted as  $n$ . OTR: Off-target translation ratio.

samples (18.5 BLEU score and 2.7% OTR score). Furthermore, our method demonstrates robustness across different translation data sizes and consistently achieves OTR scores close to zero across all four settings, resulting in significantly higher BLEU scores. **This confirms that our method remains effective with various sizes of fine-tuning data.**

### 5.5 Performance on supervised translation

Base	Methods	IWSLT	WMT	AVG
LlaMA2	MT	29.4	27.6	<u>28.5</u>
	Ours	28.8	<b>27.1</b>	<u>28.0</u>
LlaMA3	MT	29.5	<b>29.1</b>	<u>29.3</u>
	Ours	29.4	28.3	<u>28.9</u>

Table 4 Supervised translation performance. **Bold:** The best results. Underline: Average scores. *Take-away:* Our model successfully achieved the goal of improving ZST performance without compromising the effectiveness of supervised translation. ZST: Zero-shot translation.

Our algorithm primarily enhances ZST performance through unlikelihood training on instruction-conflicting samples. It raises a question: Does the supervised translation ability persist even after unlikelihood training? As shown in Table 4, we report the performance of MT and ours on the IWSLT and WMT benchmarks. Remarkably, **our models successfully retain the supervised translation ability after unlikelihood training with instruction-conflicting samples.** Specifically, our final model achieves comparable results compared with LLaMA-MT for LLaMA2-based (28.0 vs 28.5 in BLEU score) and LLaMA3-based (28.9 vs. 29.3 BLEU score) mod-

els.

Base	Methods	IWSLT	WMT	AVG
<i>XComet</i> ↑				
LlaMA2	MT	<b>88.6</b>	88.2	<u>88.4</u>
	Ours	88.4	<b>88.3</b>	<u>88.3</u>
LlaMA3	MT	85.6	89.1	<u>87.3</u>
	Ours	<b>85.7</b>	<b>89.3</b>	<u>87.5</u>
<i>CometKiwi</i> ↑				
LlaMA2	MT	<b>63.1</b>	66.9	<u>65.0</u>
	Ours	<b>63.1</b>	<b>67.8</b>	<u>65.5</u>
LlaMA3	MT	64.0	68.8	<u>66.4</u>
	Ours	<b>64.3</b>	<b>69.4</b>	<u>66.9</u>

Table 5 Evaluation with XComet and CometKiwi on samples translated into correct language by MT. **Bold:** The best results. Underline: Average score.

### 5.6 Comparison with MT models

To provide more insights into the improvements brought by our approach, we introduce model-based translation evaluation metrics to compare MT with our models. The model-based metrics evaluate the semantic similarity of model generations and human translations, including reference-based XComet<sup>3</sup> and reference-free CometKiwi<sup>4</sup>. Both of these metrics achieve better correlations with human evaluation than SacreBLEU. As they may assign high scores to language-mismatch translations, such as directly copying the source sentence or incorrectly translating into another language, we split the test set into two parts: those where MT models translate into the correct language and those where it translates into an incorrect language. We report the scores for the former and provide a case study for the latter.

As shown in Table 5, we observe that our models' generated translations achieve comparable semantic similarity to those of MT models when evaluated using model-based metrics. Moreover, our models achieve higher CometKiwi scores (+0.5) in both LLaMA2- and LLaMA3-based experiments. In Table 6, we present a case study analyzing incorrect language translations by LLaMA3-based models on the WMT En→Ja translation task. Compared to the MT model, our model generates translations that are more consistent with human translations and are

<sup>3</sup><https://huggingface.co/Unbabel/XCOMET-XL>

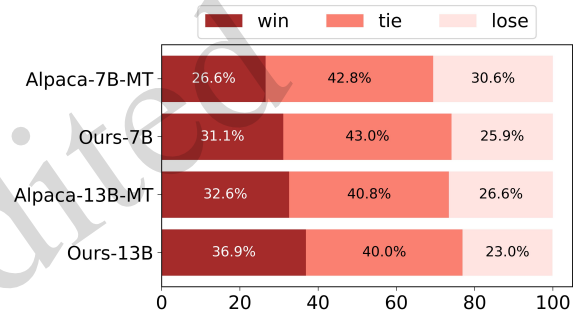
<sup>4</sup><https://huggingface.co/Unbabel/wmt23-cometkiwi-da-xl>

**Table 6** Case study of LLaMA3-based models on the WMT En→Ja translation task. We present translation examples for MT generating incorrect language translations, which are highlighted in yellow. “Input” and “Human” indicate the source English sentences and the ground-truth target Japanese sentences, respectively.

	<i>Sentence</i>		<i>Sentence</i>
<b>Input</b>	Critics say he remains inseparable from apartheid-era crimes and could have been held accountable for them had he lived longer.	<b>Input</b>	Sources have told The Age and The Sydney Morning Herald that the Gabba is the only major stadium in Australian cricket where ...
<b>Human</b>	家によると、同氏はアパルトヘイト代の犯罪とは切っても切れないにあり、もっとく生きていればその任をわけていたかもしれない。	<b>Human</b>	情筋は、世界模の映像配信に必要な膨大な数の中やデバイスを十分に稼させる会の生源の力供が不十分なのは、オーストラリアの大模スタジアムでもガバだけだと...
<b>MT</b>	批判家表示, 他始终与种族隔离时期的罪行不可分割, 倘若他活得更长一点, 就可能被追究责任。 (Translation in Chinese)	<b>MT</b>	Sources have told The Age and The Sydney Morning Herald that the Gabba is the only major stadium ... (Copy the Input)
<b>Ours</b>	批家によると、彼はアパルトヘイト代の犯罪かられられず、もし生存していたならば任を担する必要があった。	<b>Ours</b>	《The Age》と《Sydney Morning Herald》(SMH) の取材によると、ガバはオーストラリアの主要なクリケットスタジ...

	BLEU $\uparrow$	OTR (%) $\downarrow$
LLaMA2-7B-MT	5.0	65.0
LLaMA2-13B-MT	7.7	57.9
Alpaca-7B-MT	10.9	26.3
<b>Ours-7B</b>	<b>14.8</b>	<b>3.1</b>
Alpaca-13B-MT	12.8	24.7
<b>Ours-13B</b>	<b>16.2</b>	<b>2.8</b>

(a) ZST performance



(b) Comparative winning rates

**Fig. 5** Performance after combining with general tasks data. We combine the Alpaca and IWSLT datasets for fine-tuning. (a) We report the ZST performance on the IWSLT test set. **Bold:** The best results. (b) We also present the winning rates (%) on the AlpacaEval dataset. The higher win rate is better. OTR: Off-target translation ratio; ZST: Zero-shot translation.

presented in the correct target language.

These observations demonstrate that our method enhances the model’s ability to follow language-specific translation instructions while preserving its existing translation capabilities.

### 5.7 Effect on general task performance

Inspired by Jiao et al. (2023), we examine LLMs fine-tuned on a mixed-task dataset, encompassing both the general task and the translation task. Our objective is to investigate whether our approach can enhance ZST capabilities without compromising the general task performance. We construct the instruction tuning dataset by combining Alpaca<sup>5</sup> with the IWSLT translation dataset and denote models

<sup>5</sup>[https://github.com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca)

tuned on this dataset as Alpaca-(\*)-MT. The Alpaca dataset consists of 52k samples generated by Self-Instruct (Wang et al., 2023b) across various task types. For the second fine-tuning stage of our method, only the translation data part is utilized. The same hyperparameters as the main IWSLT experiments are used for fine-tuning. We use the BLEU and OTR scores on IWSLT test sets for translation evaluation. Following AlpacaEval<sup>6</sup>, we assess the general task performance with an LLM-based metric. Specifically, we use OpenAI ChatGPT (gpt-3.5-turbo) to perform the judgment automatically while taking LLaMA2-7B tuned on Alpaca data as the reference model to compute the win rate % on the AlpacaEval dataset. The higher win rate is better. The responses are

<sup>6</sup>[https://github.com/tatsu-lab/alpaca\\_eval](https://github.com/tatsu-lab/alpaca_eval)

generated with a temperature of 0.3 and a repetition penalty of 1.2.

As shown in Fig. 5a, both Alpaca-7B-MT and Alpaca-13B-MT achieve better ZST performance than models solely fine-tuning with translation samples (10.9/ 12.8 vs 5.0/ 7.7 BLEU score and 26.3/ 24.7 vs. 65.0/ 57.9 OTR score). This aligns with prior findings (Chung et al., 2024), which suggest that more training tasks bring gains for tasks not present in the training data. By tuning with UL loss, our models achieve better ZST performance than Alpaca-7B-MT/ Alpaca-13B-MT (14.8/ 16.2 vs. 10.9/ 12.8 BLEU and 3.1/ 2.8 vs. 26.3/ 24.7 OTR). For general tasks, as depicted in Fig. 5b, our models achieve comparable general task performance with Alpaca-7B-MT/ Alpaca-13B-MT. **These results confirm the effectiveness of our algorithm in addressing off-target translation issues without compromising general task performance.**

## 6 Related work

**Translation-tailored LLMs** Due to the huge cost to call the state-of-the-art LLMs, such as GPT-4 (OpenAI, 2024), there is a need to investigate how to effectively fit a smaller LLM into specific tasks, e.g., machine translation. Note that although there are some powerful sequence-to-sequence-style large-scale pre-trained machine translation models (Liu et al., 2020; Zan et al., 2022), this paper mainly focuses on the decoder-only LLMs due to their flexible interaction modes and rich world knowledge. In the field of LLM-based translation, various approaches have been proposed to optimize translation performance (Jiao et al., 2023; Zeng et al., 2023; ?; Xu et al., 2024a; Liu et al., 2023; Feng et al., 2024; Stap et al., 2024; Li et al., 2024b; Zhang et al., 2024). Parrot (Jiao et al., 2023) proposed fine-tuning the model on machine translation data with a hint that incorporated extra requirements to regulate the translation process. TIM (Zeng et al., 2023) introduced translation samples in comparisons to compute additional preference loss for regularization, exhibiting superior translation ability in both supervised and zero-shot directions. ALMA (?) proposed a two-stage approach that first fine-tuned on monolingual data of downstream languages followed by fine-tuning on high-quality translation data, which achieved significant improvement of translation quality. ALMA-R (Xu et al., 2024a)

aligned the ALMA models with the preferences of translation evaluators using contrastive preference optimization, achieving better performance than human translation on the test set. Liu et al. (2023) presented the position of instruction matters, that just moving the location of the instruction closer to the output can alleviate the instruction forgetting issue. Zhang et al. (2024) showed the overlooking of source sentences in LLMs and proposed both unsupervised and supervised approaches to improve this issue.

In contrast, we focus on the off-target problem of ZST, where the model fails to follow translation instructions, generating sequences not in the target language. Additionally, we show how instruction-conflicting samples can enhance the influence of instruction, thus mitigating the off-target problem.

**Unlikelihood training** Unlikelihood training (Welleck et al., 2020) aims to force the model to assign a lower probability to unlikely tokens. This method has been further explored in dialog tasks by Li et al. (2020), who demonstrated its effectiveness in generating more consistent and coherent human-like dialog. Santos et al. (2020) used the unlikelihood loss for ranking and proposed a generative information retrieval approach. Hosseini et al. (2021) proposed the combination of an unlikelihood objective with a reference-based setup for input sentences to model negation with pre-trained BERT (Devlin et al., 2019). Hu et al. (2023) took the semantically-similar or ambiguous tokens as negative information and acquired them via inherent uncertainty for the ASQP task. Zan et al. (2023) proposed a method to alleviate the off-target problem in translation models trained on multilingual translation corpora, which consist of millions of sentence pairs. They analyzed the impact of exposure bias on off-target translation issue and utilized incorrect language ID sentences as negative samples.

In this work, we focus on translation-tailored LLMs that reduce the reliance of translation models on expensive, large-scale parallel corpora. We take the instances where translation pairs conflict with the given instructions as negative samples for ZST. Furthermore, we consider a new case that enhances the ability of LLMs to better follow translation instructions and generate translations in the correct language.

## 7 Conclusion

We propose a simple two-stage fine-tuning strategy to enhance the instruction-following ability of LLM for translation. The core procedure consists of two main steps: (1) creating instruction-conflicting samples by replacing the translation directions with incorrect ones, and (2) training on these samples using an additional unlikelihood loss. Experimental results on IWSLT and WMT, spanning 16 ZST directions, demonstrate the effectiveness of the proposed method, which reduces the OTR and produces translations of higher quality. Furthermore, our approach exerts a negligible influence on other aspects of LLMs, such as supervised translation performance and general task performance.

## Limitations

The proposed method contains a mixing hyperparameter  $\alpha$  to balance MLE loss and UL loss in unlikelihood training on instruction-conflicting samples, and the high  $\alpha$  may overfit the model on UL loss. In future work, we may focus on how to balance them adaptively.

This work only focuses on the off-target problem in the ZST of LLMs, which could be seen as a specific type of input-conflicting hallucination. In future work, we will continue to explore the application of the unlikelihood training on general tasks, such as programming, math, dialog, etc., and more types of hallucinations (Zhang et al., 2023), such as fact-conflicting hallucinations, and context-conflicting hallucinations. Also, we will apply our method to enhance the instruction-following abilities of different interesting LLM-based scenarios, e.g., safety (Miao et al., 2024; Zhang et al., 2024; Zhong et al., 2024), debiasing (Xu et al., 2024b), multimodal analysis (Wang et al., 2024b), healthcare (Ren et al., 2024), and difficult code generation (Wang et al., 2024a).

## Acknowledge

This work was supported by National Natural Science Foundation of China (Grant No. 62372468), Shandong Natural Science Foundation (Grant No. ZR2023MF008), Major Basic Research Projects in Shandong Province (Grant No. ZR2023ZD32) and Qingdao Natural Science Foundation (Grant No. 23-

2-1-161-zyyd-jch).

## Contributors

Changtong Zan and Liang Ding designed the research. Changtong Zan processed the data. Changtong Zan drafted the paper. Li Shen, Yibing Zhan, and Xinghao Yang helped organize the paper. Changtong Zan, Liang Ding, and Weifeng Liu revised and finalized the paper.

## Compliance with ethics guidelines

Changtong Zan, Liang Ding, Li Shen, Yibing Zhan, Xinghao Yang, and Weifeng Liu declare that they have no conflict of interest.

## References

- Brown TB, Mann B, Ryder N, et al., 2020. Language models are few-shot learners. Proc Conf of Advances in Neural Information Processing Systems, p.1877-1901.
- Chen L, Ma SM, Zhang DD, et al., 2023. On the off-target problem of zero-shot multilingual neural machine translation. Findings of the Association for Computational Linguistics, p.9542-9558. <https://doi.org/10.18653/v1/2023.findings-acl.608>
- Cho K, van Merriënboer B, Gulcehre C, et al., 2014. Learning phrase representations using rnn encoder-decoder for statistical machine translation. Proc Conf on Empirical Methods in Natural Language Processing, p.1724-1734. <https://doi.org/10.3115/v1/D14-1179>
- Chung HW, Hou L, Longpre S, et al., 2024. Scaling instruction-finetuned language models. *J Mach Learn Res*, 25(70):1-53.
- Dabre R, Kurohashi S, 2019. Mmc4nlp: multilingual multi-way corpora repository for natural language processing. <https://doi.org/10.48550/arXiv.1710.01025>
- Devlin J, Chang MW, Lee K, et al., 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. Proc Conf of the North American Chapter of the Association for Computational Linguistics, p.4171-4186. <https://doi.org/10.18653/v1/N19-1423>
- Feng ZP, Chen RZ, Zhang Y, et al., 2024. Ladder: A model-agnostic framework boosting llm-based machine translation to the next level. Findings of the Association for Computational Linguistics, p.15377-15393. <https://doi.org/10.18653/v1/2024.emnlp-main.860>
- Fu Y, Peng H, Ou LT, et al., 2023. Specializing smaller language models towards multi-step reasoning. Proc 40<sup>th</sup> Int Conf on Machine Learning, p.10421-10430.
- Gu JT, Wang Y, Cho K, et al., 2019. Improved zero-shot neural machine translation via ignoring spurious correlations. Proc 57<sup>th</sup> Annual Meeting of the Association for Computational Linguistics, p.1258-1268. <https://doi.org/10.18653/v1/P19-1121>
- Hendy A, Abdelrehim M, Sharaf A, et al., 2023. How good are gpt models at machine translation? a comprehensive

- evaluation.  
<https://doi.org/10.48550/arXiv.2302.09210>
- Hosseini A, Reddy S, Bahdanau D, et al., 2021. Understanding by understanding not: Modeling negation in language models. Proc Conf of the North American Chapter of the Association for Computational Linguistics, p.1301-1312.  
<https://doi.org/10.18653/v1/2021.naacl-main.102>
- Hu MT, Bai YH, Wu YK, et al., 2023. Uncertainty-aware unlikelihood learning improves generative aspect sentiment quad prediction. Findings of the Association for Computational Linguistics, p.13481-13494.  
<https://doi.org/10.18653/v1/2023.findings-acl.851>
- Huang YX, Gu HL, Yu ZT, et al., 2024. Enhancing low-resource cross-lingual summarization from noisy data with fine-grained reinforcement learning. *Front of Inform Technol & Electron Eng*, 25(1):121-134.  
<https://doi.org/10.1631/FITEE.2300296>
- Jiao WX, Huang JT, Wang WX, et al., 2023. Parrot: Translating during chat using large language models. Findings of Empirical Methods in Natural Language Processing, p.15009-15020.  
<https://doi.org/10.18653/v1/2023.findings-emnlp.1001>
- Joulin A, Grave E, Bojanowski P, et al., 2016a. Bag of tricks for efficient text classification. Proc 15<sup>th</sup> Conf of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers, p.427-431.
- Joulin A, Grave E, Bojanowski P, et al., 2016b. Fasttext.zip: Compressing text classification models.  
<https://doi.org/10.48550/arXiv.1612.03651>
- Kaplan J, McCandlish S, Henighan T, et al., 2020. Scaling laws for neural language models.  
<https://doi.org/10.48550/arXiv.2001.08361>
- Kwon W, Li ZH, Zhuang SY, et al., 2023. Efficient memory management for large language model serving with pagedattention. Proc 29<sup>th</sup> Symp on Operating Systems Principles, p.611-626.  
<https://doi.org/10.1145/3600006.3613165>
- Li B, Yang P, Sun YK, et al., 2024a. Advances and challenges in artificial intelligence text generation. *Front of Inform Technol & Electron Eng*, 25(1):64-83.  
<https://doi.org/10.1631/FITEE.2300410>
- Li JH, Zhou H, Huang SJ, et al., 2024b. Eliciting the translation ability of large language models via multilingual finetuning with translation instructions. *Trans Assoc Comput Linguist*, 12:576-592.  
[https://doi.org/10.1162/tacl\\_a\\_00655](https://doi.org/10.1162/tacl_a_00655)
- Li M, Roller S, Kulikov I, et al., 2020. Don't say that! making inconsistent dialogue unlikely with unlikelihood training. Proc 58<sup>th</sup> Annual Meeting of the Association for Computational Linguistics, p.4715-4728.  
<https://doi.org/10.18653/v1/2020.acl-main.428>
- Liu YJ, Zeng XF, Meng FD, et al., 2023. Instruction position matters in sequence generation with large language models. Findings of the Association for Computational Linguistics, p.11652-11663.  
<https://doi.org/10.18653/v1/2024.findings-acl.693>
- Liu YH, Gu JT, Goyal N, et al., 2020. Multilingual denoising pre-training for neural machine translation. *Trans Assoc Comput Linguist*, 8:726-742.  
[https://doi.org/10.1162/tacl\\_a\\_00343](https://doi.org/10.1162/tacl_a_00343)
- Lu QY, Qiu BP, Ding L, et al., 2023. Error analysis prompting enables human-like translation evaluation in large language models: A case study on chatgpt. Findings of the Association for Computational Linguistics, p.8801-8816.  
<https://doi.org/10.18653/v1/2024.findings-acl.520>
- Miao YC, Zhang S, Ding L, et al., 2024. Inform: Mitigating reward hacking in rlhf via information-theoretic reward modeling. Proc Conf of Advances in Neural Information Processing Systems.
- Min BN, Ross H, Sulem E, et al., 2023. Recent advances in natural language processing via large pre-trained language models: A survey. *ACM Comput Surv*, 56(2):1-40.  
<https://doi.org/10.1145/3605943>
- Mishra S, Khashabi D, Baral C, et al., 2022. Cross-task generalization via natural language crowdsourcing instructions. Proc 60<sup>th</sup> Annual Meeting of the Association for Computational Linguistics, p.3470-3487.  
<https://doi.org/10.18653/v1/2022.acl-long.244>
- OpenAI, 2024. Gpt-4 technical report.  
<https://doi.org/10.48550/arXiv.2303.08774>
- Peng KQ, Ding L, Zhong QH, et al., 2023. Towards making the most of chatgpt for machine translation. Findings of Empirical Methods in Natural Language Processing, p.5622-5633.  
<https://doi.org/10.18653/v1/2023.findings-emnlp.373>
- Post M, 2018. A call for clarity in reporting BLEU scores. Proc Third Conf on Machine Translation, p.186-191.  
<https://doi.org/10.18653/v1/W18-6319>
- Qu Z, Watanabe T, 2022. Adapting to non-centered languages for zero-shot multilingual translation. Proc 29<sup>th</sup> Int Conf on Computational Linguistics, p.5251-5265.
- Ren ZY, Zhan YB, Yu BS, et al., 2024. Healthcare copilot: Eliciting the power of general llms for medical consultation.  
<https://doi.org/10.48550/arXiv.2402.13408>
- Nogueira dos Santos C, Ma XF, Nallapati R, et al., 2020. Beyond [CLS] through ranking by generation. Proc Conf on Empirical Methods in Natural Language Processing, p.1722-1727.  
<https://doi.org/10.18653/v1/2020.emnlp-main.134>
- Sennrich R, Vamvas J, Mohammadshahi A, 2024. Mitigating hallucinations and off-target machine translation with source-contrastive and language-contrastive decoding. Proc Conf of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers, p.21-33.
- Stap D, Hasler E, Byrne B, et al., 2024. The fine-tuning paradox: Boosting translation quality without sacrificing LLM abilities. Proc 62<sup>nd</sup> Annual Meeting of the Association for Computational Linguistics, p.6189-6206.  
<https://doi.org/10.18653/v1/2024.acl-long.336>
- Touvron H, Lavril T, Izacard G, et al., 2023a. Llama: Open and efficient foundation language models.  
<https://doi.org/10.48550/arXiv.2302.13971>
- Touvron H, Martin L, Stone K, et al., 2023b. Llama 2: Open foundation and fine-tuned chat models.  
<https://doi.org/10.48550/arXiv.2307.09288>
- Wang S, Ding L, Shen L, et al., 2024a. Oop: Object-oriented programming evaluation benchmark for large language models. Findings of the Association for Computational

- Linguistics, p.13619-13639.  
<https://doi.org/10.18653/v1/2024.findings-acl.808>
- Wang WB, Ding L, Shen L, et al., 2024b. Wisdom: Improving multimodal sentiment analysis by fusing contextual world knowledge. Proc Conf on ACM Multimedia, p.2282-2291.
- Wang YM, Zhang ZS, Wang R, 2023a. Element-aware summarization with large language models: Expert-aligned evaluation and chain-of-thought method. Proc 61<sup>st</sup> Annual Meeting of the Association for Computational Linguistics, p.8640-8665.  
<https://doi.org/10.18653/v1/2023.acl-long.482>
- Wang YZ, Kordi Y, Mishra S, et al., 2023b. Self-instruct: Aligning language models with self-generated instructions. Proc 61<sup>st</sup> Annual Meeting of the Association for Computational Linguistics, p.13484-13508.  
<https://doi.org/10.18653/v1/2023.acl-long.754>
- Wei J, Bosma M, Zhao V, et al., 2021. Finetuned language models are zero-shot learners. Proc Conf of the International Conference on Learning Representations.
- Wei J, Wang XZ, Schuurmans D, et al., 2022. Chain-of-thought prompting elicits reasoning in large language models. Proc Conf of Advances in Neural Information Processing Systems.
- Welleck S, Kulikov I, Roller S, et al., 2020. Neural text generation with unlikelihood training. Proc Conf of the International Conference on Learning Representations.
- Wolf T, Debut L, Sanh V, et al., 2020. Transformers: State-of-the-art natural language processing. Proc Conf on Empirical Methods in Natural Language Processing: System Demonstrations, p.38-45.  
<https://doi.org/10.18653/v1/2020.emnlp-demos.6>
- Xu HR, Sharaf A, Chen YM, et al., 2024a. Contrastive preference optimization: Pushing the boundaries of LLM performance in machine translation. Proc 41<sup>th</sup> Int Conf on Machine Learning.
- Xu ZY, Peng KQ, Ding L, et al., 2024b. Take care of your prompt bias: Investigating and mitigating the prompt bias in factual knowledge extraction. Proc Joint Int Conf on Computational Linguistics, Language Resources and Evaluation, p.15552-15565.
- Zan CT, Peng KQ, Ding L, et al., 2022. Vega-mt: The jd explore academy machine translation system for wmt22. Proc Seventh Conf on Machine Translation, p.411-422.
- Zan CT, Ding L, Shen L, et al., 2023. Unlikelihood tuning on negative samples amazingly improves zero-shot translation.  
<https://doi.org/10.48550/arXiv.2309.16599>
- Zeng JL, Meng FD, Yin YJ, et al., 2023. Tim: Teaching large language models to translate with comparison. Proc of the AAAI Conference on Artificial Intelligence, p.19488-19496.  
<https://doi.org/10.1609/aaai.v38i17.29920>
- Zhang HB, Chen Q, Zhang WW, 2022. Improving entity linking with two adaptive features. *Front of Inform Technol & Electron Eng*, 23(11):1620-1630.  
<https://doi.org/10.1631/FITEE.2100495>
- Zhang HB, Chen KH, Bai XF, et al., 2024. Paying more attention to source context: Mitigating unfaithful translations from large language model. Findings of the Association for Computational Linguistics, p.13816-13836.  
<https://doi.org/10.18653/v1/2024.findings-acl.821>
- Zhang SS, Roller S, Goyal N, et al., 2022. Opt: Open pre-trained transformer language models.  
<https://doi.org/10.48550/arXiv.2205.01068>
- Zhang YQ, Ding L, Zhang LF, et al., 2024. Intention analysis makes llms a good jailbreak defender.  
<https://doi.org/10.48550/arXiv.2401.06561>
- Zhang Y, Li YF, Cui LY, et al., 2023. Siren's song in the ai ocean: A survey on hallucination in large language models.  
<https://doi.org/10.48550/arXiv.2309.01219>
- Zhong QH, Ding L, Liu JH, et al., 2023. Can chatgpt understand too? a comparative study on chatgpt and fine-tuned bert.  
<https://doi.org/10.48550/arXiv.2302.10198>
- Zhong QH, Ding L, Liu JH, et al., 2024. Rose doesn't do that: Boosting the safety of instruction-tuned large language models with reverse prompt contrastive decoding. Findings of the Association for Computational Linguistics, p.13721-13736.  
<https://doi.org/10.18653/v1/2024.findings-acl.814>