

# Word Alignment as Preference for Machine Translation

Anonymous ACL submission

## Abstract

The problem of hallucination and omission, a long-standing problem in machine translation (MT), is more pronounced when a large language model (LLM) is used in MT because an LLM itself is susceptible to these phenomena. In this work, we mitigate the problem in an LLM-based MT model by guiding it to better word alignment. We first study the correlation between word alignment and the phenomena of hallucination and omission in MT. Then we propose to utilize word alignment as preference to optimize the LLM-based MT model. The preference data are constructed by selecting chosen and rejected translations from multiple MT tools. Subsequently, direct preference optimization is used to optimize the LLM-based model towards the preference signal. Given the absence of evaluators specifically designed for hallucination and omission in MT, we further propose selecting hard instances and utilizing GPT-4 to directly evaluate the performance of the models in mitigating these issues. We verify the rationality of these designed evaluation methods by experiments, followed by extensive results demonstrating the effectiveness of word alignment-based preference optimization to mitigate hallucination and omission.

## 1 Introduction

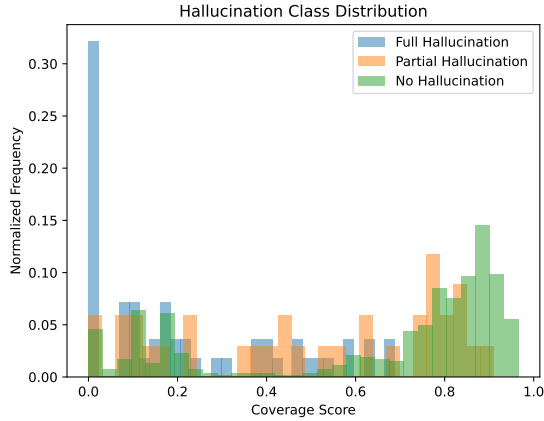
Large language models (LLMs) have been evolving rapidly and showing predominant performance in many natural language processing (NLP) tasks (Brown et al., 2020; Achiam et al., 2023; Touvron et al., 2023). However, in machine translation (MT), the use of a decoder-only LLM is still limited due to issues such as model size (Xu et al., 2024a) and low-resource languages (Hendy et al., 2023). Conventional encoder-decoder MT models trained on parallel corpora still dominate in practice (Costa-jussà et al., 2022). One of the primary concerns of applying an LLM to MT is reliability. Although it does not happen frequently, an LLM is known to hallucinate (Dhuliawala et al., 2023;

Zhang et al., 2023a; Bang et al., 2023) as it is pre-trained to predict the next token in very large-scale raw texts. Specifically in MT, LLM-based translation systems therefore could have the phenomena of hallucination and omission, which is also a long-term challenge in the field of MT (Yang et al., 2019; Vamvas and Sennrich, 2022), known as over- and under-translation. In this work, we attempt to mitigate the hallucination and omission in LLM-based MT to improve its practicality.

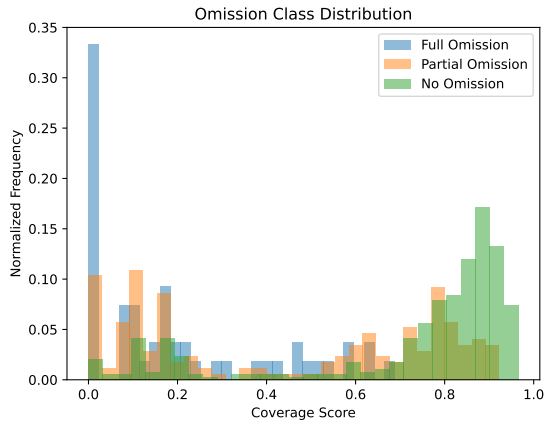
Hallucination in MT occurs when information not present in the source text is generated in the translation, and omission occurs when some of the information in the source text is missed in the translation. As a related tool that explicitly aligns the source text and translation at the word level, word alignment is potentially positive for MT due to the nature of align and translate (Bahdanau et al., 2015). The degree of coverage of the source text in translation could be a direct signal to identify the hallucination and omission in MT (Tu et al., 2016). Figure 1 shows the normalized frequency of the coverage scores predicted by a word aligner. The examples that are annotated as “no hallucination or omission” tend to have a higher coverage score, while those in “full hallucination or omission” are more likely to have an extremely low coverage score. “small hallucination or omission” and “partial hallucination or omission” distribute in the middle. As the annotations are carefully made by humans and highly correlates to the coverage scores from the word aligner, this indicates that word alignment is a simple but promising direction to mitigate these phenomena.

Consequently, we propose Word Alignment Preference (WAP) that utilizes word alignment as a signal to optimize LLM-based MT models. WAP consists of three steps: diverse translation collection, preference data construction, and preference optimization. Specifically, we collect diverse translations with multiple existing translation tools, se-

043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083



(a) Coverage distribution of different hallucination degree.



(b) Coverage distribution of different omission degree.

Figure 1: A preliminary experiment shows that higher coverage scores correlates to less hallucination and omission. The coverage scores are predicted by a word aligner (Wu et al., 2023). The human annotation of hallucination and omission is from HalOmi benchmark (Dale et al., 2023b). Details about the dataset and word alignment model can be found in §3.1.

lect chosen and rejected examples with the word aligner (Wu et al., 2023), and optimize the model on preference data using direct preference optimization (DPO) (Rafailov et al., 2024).

Furthermore, the evaluation of hallucination and omission is challenging, and there is no existing evaluator specifically designed for this. Improving the BLEU and COMET score does not necessarily mean reducing hallucination and omission because there are other factors such as mistranslation and fluency. In addition, hallucination is relatively infrequent, although very severe once happens. Hence, to effectively evaluate it, we design extensive experiments that include testing on instances that potentially have the problem of hallucination and omission, and using GPT-4 as the evaluator with comprehensive analysis. Experi-

tal analysis demonstrates the effectiveness of WAP in mitigating hallucination and omission in MT.

In summary, the contributions of this work include the following:

- We studied the correlation between the coverage score by word alignment and the phenomena of hallucination and omission in MT. From the preliminary experiments in Figure 1 we found that word alignment is a promising signal to mitigate it.
- In §2 we propose a novel approach, namely WAP, to construct a word alignment-based preference dataset, and use DPO to optimize the LLM-based MT model. The validity of the preference dataset is also demonstrated by direct fine-tuning on preferred and rejected translations in §4.4.
- As there is no benchmark particularly for evaluating the performance of MT models on hallucination and omission. We design various experiments, including selecting hard instances and utilizing GPT-4 as the evaluator in §3.2. The effectiveness of the evaluation, as well as the proposed WAP has been validated through experiments and analysis in §4

## 2 Proposed approach

### 2.1 Gathering translation candidates

To steer the MT model to avoid hallucination and omission using preference optimization, we first need comparable but different translations. Starting with a source text  $x$ , we utilize  $K$  methods to produce translations, notated as  $\pi^1, \dots, \pi^K$ . Then we can get a set of translations  $Y$ , in which  $y^k \in Y$  is obtained by  $y^k = \pi^k(x)$  and  $|Y| = K$ .

**Details of gathered translations** We start with the parallel training data in ALMA (Xu et al., 2024a). This parallel data encompasses five language pairs with human translations in both directions:  $cs \leftrightarrow en$ ,  $de \leftrightarrow en$ ,  $is \leftrightarrow en$ ,  $zh \leftrightarrow en$  and  $ru \leftrightarrow en$ . We employ ISO 639 language codes<sup>1</sup> to denote languages. Specifically, “*cs*” corresponds to Czech, “*de*” to German, “*is*” to Icelandic, “*zh*” to Chinese and “*ru*” and “*en*” to Russian and English, respectively. To generate the translations we require, this dataset is translated in both directions

<sup>1</sup>[https://en.wikipedia.org/wiki/List\\_of\\_ISO\\_639\\_language\\_codes](https://en.wikipedia.org/wiki/List_of_ISO_639_language_codes)

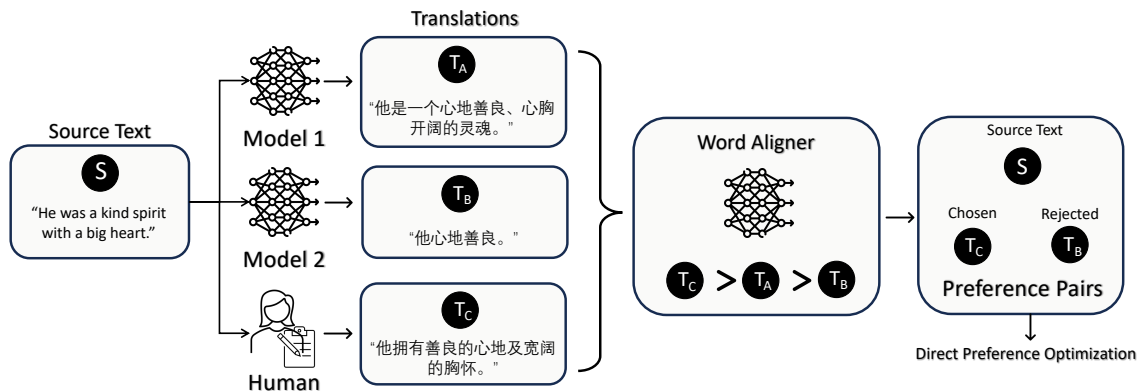


Figure 2: An illustration of WAP framework. The source is first translated by multiple MT tools, including human translation. An external word aligner is then utilized to predict the coverage score for each translation. Finally, translation with the highest and lowest coverage score are selected as preference pairs for preference optimization.

using two well-known MT tools, including DeepL<sup>2</sup> and ChatGPT (gpt-3.5-turbo-0613)<sup>3</sup>. The prompt for ChatGPT that we utilize to translate sentences is shown in Figure 6. The original human-written translation in the training set is also utilized. In particular, Icelandic (*is*) is not supported by the DeepL API, therefore, we use the Google Translate API<sup>4</sup> as an alternative.

## 2.2 Selecting chosen and rejected translation

After obtaining the translation candidates ( $y^1, \dots, y^K$ ), we use a state-of-the-art public word aligner, namely WSPAlign<sup>5</sup>, to automatically annotate the degree of coverage for each translation. We follow the usage setting in the original paper of WSPAlign (Wu et al., 2023). In particular, WSPAlign performs a bidirectional alignment and uses a threshold to filter out low-confident alignment of word pairs. Then, the ratio of the source words, *that are aligned with at least one word*, in the translation is taken as the coverage score, which will be used for the following preference annotation. The whole process predicting the coverage score is notated as  $C(\cdot, \cdot)$ . Formally, the coverage score for a translation  $y^k$  can be calculated by  $C(x, y^k) \in [0.0, 100.0]$ . Subsequently, the preferred translation and the rejected translation are selected as follows:

<sup>2</sup><https://www.deepl.com/en/translator>

<sup>3</sup><https://platform.openai.com/docs/models/gpt-3-5-turbo>

<sup>4</sup><https://cloud.google.com/translate/docs/basic/translate-text-basic>

<sup>5</sup><https://github.com/qiyuw/WSPAlign>

$$\begin{aligned}
 y^w &= \arg \max_{y^k \in Y} C(x, y^k) \\
 y^l &= \arg \min_{y^k \in Y} C(x, y^k)
 \end{aligned}
 \tag{1}$$

where  $y^w$  is the chosen translation and  $y^l$  is the rejected one. Then a triplet  $(x, y^w, y^l)$  is constructed for the following preference optimization.

## 2.3 Filtering

Note that the whole pipeline of constructing the preference data is automatic, and existing MT and word alignment models are not perfect. Even for human-annotated translation, the quality of it is also an issue that cannot be ignored (Xu et al., 2024b), and may affect the performance of the model trained on it. Hence, noises are inevitable in both the translated texts and the preference choices. On the other hand, the MT tools we choose generally have good performance, it could happen that the generated translations are not diverse enough, leading to the preference signal being disrupted. To improve the quality of the constructed preference datasets as much as possible, multiple strategies are applied to filter out potential bad training instances:

- Remove the instance when the chosen and rejected translations only have a marginal difference in coverage score. The difference threshold is empirically set as 5.0, that is,  $(x, y^w, y^l)$  is excluded from the dataset if  $C(x, y^w) - C(x, y^l) < 5.0$ .
- Remove the instance when the chosen and rejected translations are too semantically similar.

Sentence embedding is a widely used technique for sentence similarity with low computation cost (Gao et al., 2021; Wu et al., 2022; Zhao et al., 2024). LaBSE (Feng et al., 2022)<sup>6</sup> is used in our experiments. We notate it as  $LB(\cdot)$ . The similarity threshold is empirically set as 0.9, i.e.  $(x, y^w, y^l)$  is excluded from the dataset if  $\text{sim}(LB(y^w), LB(y^l)) > 0.9$ .  $\text{sim}(\cdot, \cdot) \in [0.0, 1.0]$  is cosine similarity.

- One possible failure case for word alignment is when the MT models directly copy the original texts, which is bad translation, but gets a high alignment score because the wrong translation is partially the same with the original texts. To remove this part of the noise, we calculate the BLEU score (Papineni et al., 2002)<sup>7</sup> for the chosen translation and exclude it if the BLEU score  $> 20.0$ .

The details and analysis of the final preference dataset after filtering is introduced in §3.1.

## 2.4 Optimization LLM-based MT model

The final step is to optimize the LLM-based MT model on our preference data. Direct preference optimization (DPO) (Rafailov et al., 2024) is a simple but effective approach that directly optimizes the preference model on a pre-constructed static dataset. DPO has been applied to optimize LLM in preference data (Tunstall et al., 2023; Xu et al., 2024b) recently. We also utilize DPO as an optimization approach. Formally, the training objective is as follows,

$$l = -\log \sigma\left(\beta \log \frac{\pi(y^w|x)}{\pi_{ref}(y^w|x)} - \beta \log \frac{\pi(y^l|x)}{\pi_{ref}(y^l|x)}\right) \quad (2)$$

where  $\sigma$  is the sigmoid function,  $\pi$  is the model to optimize and  $\pi_{ref}$  is the reference model. We use ALMA-13B<sup>8</sup> as our base model, i.e., the starting point of  $\pi$ , in the experiments. ALMA-13B is also used as a reference model  $\pi_{ref}$ , but note that  $\pi_{ref}$  will not be updated during training.

## 3 Evaluation

The experimental setup is introduced in §A.

<sup>6</sup><https://huggingface.co/sentence-transformers/LaBSE>

<sup>7</sup><https://github.com/mjpost/sacrebleu>

<sup>8</sup><https://github.com/felixxu/ALMA>

## 3.1 Baselines and evaluation datasets

We choose ALMA-13B<sup>9</sup> as the baseline for all experiments in this paper, as well as the starting point of optimization. ALMA (Xu et al., 2024a) was trained from Llama (Touvron et al., 2023) in two steps: initial fine-tuning on monolingual data and subsequent fine-tuning on a small set of high-quality parallel data.

For fairly studying the effect of word alignment preference, we use the data used in the supervised fine-tuning in ALMA as the source dataset to construct our preference data in §2. Specifically, the source data was collected from WMT’17 (Bojar et al., 2017) to WMT’20 (Barrault et al., 2020), in addition to the development and text dataset from Flores-200 (Costa-jussà et al., 2022). After filtering, we finally make 20,074 and 2,226 preference triplets for training and development, respectively. For evaluation, the test set is from WMT22, except that *is*  $\leftrightarrow$  *en* is from WMT21. The remaining data from WMT21 (except *is*  $\leftrightarrow$  *en*) is used as the development set. Specifically, 3485, 4021, 2000, 3912, 4053 examples are included in the test set for *cs*  $\leftrightarrow$  *en*, *de*  $\leftrightarrow$  *en*, *is*  $\leftrightarrow$  *en*, *zh*  $\leftrightarrow$  *en*, and *ru*  $\leftrightarrow$  *en*, respectively.

**HalOmi** In particular, we want to validate whether our proposed method is capable of mitigating hallucination and omission in MT. Hence, we also utilize HalOmi (Dale et al., 2023b) in the experiments. HalOmi is an evaluation benchmark for the detection of hallucination and omission in MT. It contains fine-grained sentence-level and token-level annotations of full and partial hallucinations and omissions that cover 18 language directions. Each instance in the data set was annotated in “No hallucination and omission”, “Small hallucination and omission”, “Partial hallucination and omission” or “Full hallucination and omission” by humans. In this paper, we use it to test the performance of GPT-4 as an evaluator. Details are in §3.2.

## 3.2 The design of evaluation

We focus on optimizing LLM-based MT models to avoid hallucination and omission. However, to our best knowledge, there is no benchmark measuring MT models specifically for this issue, making the evaluation very challenging. Improving the BLEU or COMET score does not necessarily mean reducing hallucination and omission because there are

<sup>9</sup><https://huggingface.co/haoranxu/ALMA-13B>



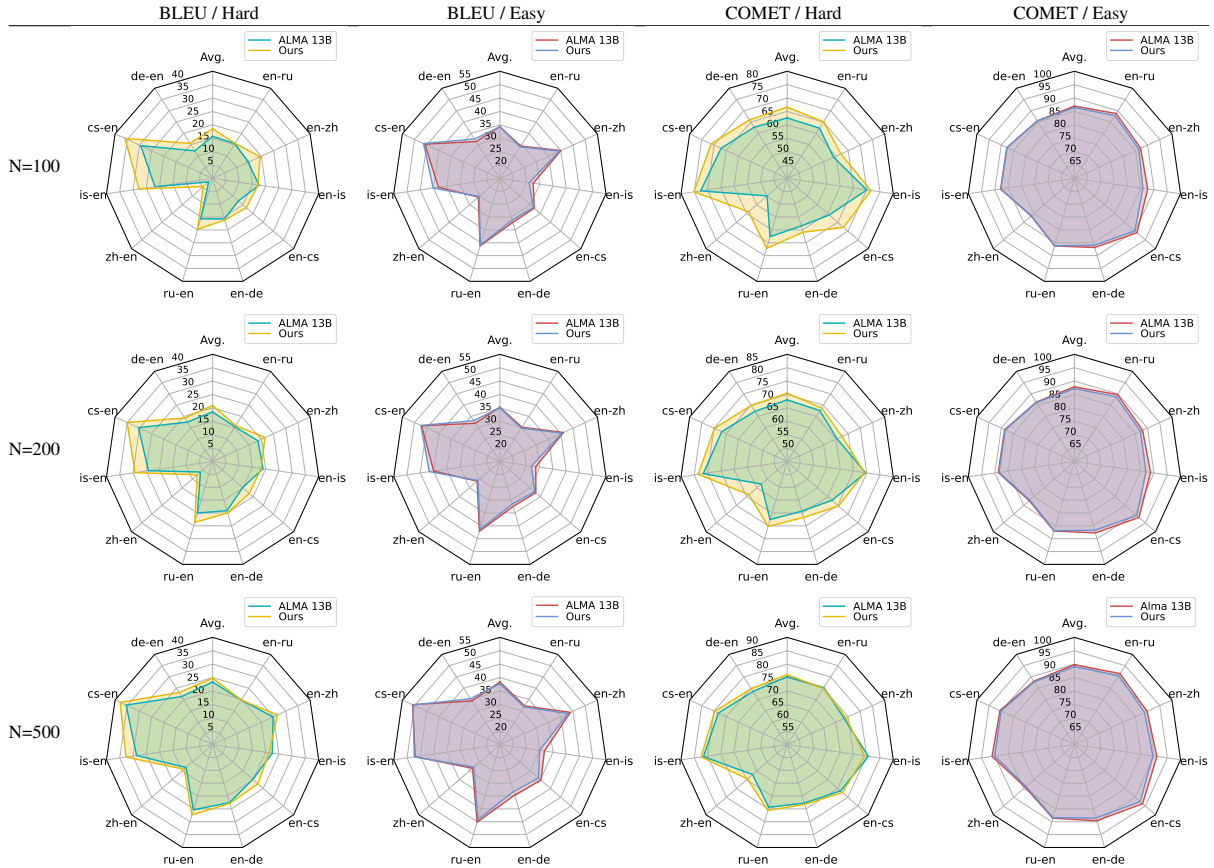


Figure 3: Comparison of WAP and baseline in hard and easy instances.  $N$  instances with the lowest COMET score by the baseline are selected from the test set as hard instances, and the remaining are easy instances. Results when  $N = 100, 200$  and  $500$  are presented. Refer to §D for the full numeric results of the entire test.

other factors such as mistranslation and fluency. In addition, hallucination is relatively infrequent, although very severe once happens. To intuitively validate whether our approach is capable of mitigating hallucination and omission in MT, we design several evaluation strategies in this section.

**Select hard instances.** We first select instances that the baseline model does not perform well on. This subset of instances is labeled as *hard instances* in this paper. The subset of the remaining examples is labeled as *easy instances*. Specifically,  $N$  instances with the lowest COMET score are selected from the test set for each translation direction. As hard examples tend to include more hallucination and omission, we report the comparison of models on hard examples and remaining examples, respectively. In the experiment, we sample three subsets where  $N = 100, N = 200$  and  $N = 500$ . The experimental analysis can be found in §4.1. Note that the hard instances are only selected for evaluation. We do not differentiate hard or easy instances in the training set. Only word alignment signal is used to

select preferred dataset for a fair comparison.

	Hallucination			Omission		
	No	Partial	Full	No	Partial	Full
# of examples	817	42	65	627	237	60
Avg. score	84.19	45.95	3.84	87.97	66.28	1.66
Pearson Corr.	0.5969			0.5686		

Table 1: Average coverage score calculated by GPT-4 for different level of hallucination or omission. The Pearson Correlation between the annotated labels and GPT-4 coverage scores is also reported. Ideally, higher score should correlate to less hallucination and omission.

**Utilize LLM as the evaluator for hallucination and omission.** Besides the BLEU and COMET in hard instances, a direct estimate of the degree of hallucination and omission in translation is still needed. As we mentioned earlier that improving the BLEU and COMET score does not necessarily mean reducing hallucination and omission because there are other factors such as mistranslation and fluency, we utilize the generalization and reasoning ability of LLM (Kojima et al., 2022; Mitchell

et al., 2023; Wei et al., 2023) to achieve this direct evaluation. We use one of the most powerful LLM, gpt-4-0613<sup>10</sup>, as the evaluator. LLM is prompted to check whether the given translation has hallucination or omission referring to the given source texts. A coverage score between 0 and 100 is output as the degree metric. The prompt used is shown in Figure 7.

**Is LLM really capable of evaluating hallucination and omission in MT?** Despite the fact that LLMs have shown impressive zero-shot performance in various tasks (Kojima et al., 2022; Mitchell et al., 2023; Wei et al., 2023), the assessment of LLM in the evaluation of hallucination and omission is still important because it has not been widely used on this task. We use HalOmi datasets introduced in §3.1 to assess this ability of GPT-4. The examples in *de* ↔ *en*, *zh* ↔ *en*, and *ru* ↔ *en* are selected, then GPT-4 is used to predict the coverage score for these examples.

Table 1 shows the average score of the degree of coverage predicted by GPT-4. The examples from HalOmi are split into three subsets based on the labels. We merged the “Partial hallucination and omission” and “Small hallucination and omission” in the original because the number of examples in these two categories is small. It clearly demonstrates that examples annotated as “No hallucination and omission” have a higher coverage score predicted by GPT-4 and those in “Full hallucination and omission” have an extremely low coverage score. As a result, using GPT-4 is an effective way to assess whether a translation has the problem of hallucination or omission.

## 4 Experimental results

### 4.1 Evaluation on hard instances

In §3.2 we introduce how to select hard instances from the test set and explain why hard instances are suitable to assess hallucination and omission. In this section, we evaluate our model on these hard instances and the remaining examples, respectively. Figure 3 demonstrates the results when the number of hard instances  $N = 100, 200,$  and  $500,$  respectively. The following findings can be concluded:

- WAP consistently outperforms the baseline in hard instances in most translation directions, for both BLEU and COMET metrics.

<sup>10</sup><https://platform.openai.com/docs/models/gpt-4-turbo-and-gpt-4>

- WAP generally reaches competitive performance compared to the baseline for both BLEU and COMET.
- With increasing the number of hard instances, the improvement gained by WAP gets smaller.

These results indicate that WAP mitigates hallucination and omission to a certain extent, because these issues are more likely to occur in hard instances. In addition, with the improvement in the hard instances, our model remains competitive to the baseline in the remaining easy instances. It is reasonable that there is no significant difference in the easy instances because the compared models are generally good. The challenging part should be in the hard ones. Moreover, it is also observed that with increasing  $N$ , the improvement gets narrower. The reason is that more relatively easy instances are included in the subset. This is another evidence that WAP provides gains particularly for hallucination and omission in MT. The specific numeric results and the overall results for the entire test set are shown in §D.

### 4.2 Direct evaluation of hallucination and omission by GPT-4

In addition to improving hard examples, which is more likely to have hallucination and omission, direct evaluations of them are also needed to confirm the effectiveness of the proposed WAP. In §3.2 we have verified the usefulness of GPT-4 as an evaluator with experiments. In this section, we prompt GPT-4 to directly predict a coverage score as the metric of hallucination and omission. The results are demonstrated in Table 2. The reported number is the average of the coverage scores in hard examples. The results show that our model outperforms the baseline in all translation directions except *en* ↔ *is*. Specifically in the average score of all translation directions, WAP outperforms the baseline model by **4.96, 1.63** and **1.24** when  $N=100, 200$  and  $500,$  respectively. The trend is similar to that of §4.1, which directly indicates that the LLM-based MT model is steered to avoid generating hallucination and omission in MT with the preference dataset we constructed.

### 4.3 Human evaluation

Although the validity of GPT-4 as evaluator for hallucination and omission has been demonstrated in §3.2 and Table 1, we conduct a human evaluation to further verify our findings, as LLM could

	de-en	cs-en	is-en	zh-en	ru-en	en-de	en-cs	en-is	en-zh	en-ru	Avg.
N=100											
Baseline	94.30	92.95	94.90	63.08	89.85	92.85	82.75	<b>97.05</b>	84.65	90.53	88.29
+WAP	<b>95.85</b>	<b>94.65</b>	<b>96.05</b>	<b>80.23</b>	<b>91.75</b>	<b>96.25</b>	<b>91.85</b>	96.10	<b>92.90</b>	<b>96.87</b>	<b>93.25(+4.96)</b>
N=200											
Baseline	95.71	95.05	95.45	74.83	92.83	94.20	89.95	<b>97.70</b>	89.19	94.25	91.92
+WAP	<b>97.10</b>	<b>96.55</b>	<b>97.48</b>	<b>85.63</b>	<b>95.53</b>	<b>95.18</b>	<b>91.84</b>	96.73	<b>92.81</b>	<b>96.66</b>	<b>94.55(+2.63)</b>
N=500											
Baseline	97.18	96.74	97.29	87.85	96.16	97.35	94.46	98.21	91.64	96.10	95.30
+WAP	<b>98.10</b>	<b>97.79</b>	<b>98.12</b>	<b>90.76</b>	<b>97.82</b>	<b>97.36</b>	<b>96.05</b>	<b>98.22</b>	<b>94.07</b>	<b>97.13</b>	<b>96.54(+1.24)</b>

Table 2: Coverage score output by GPT-4. The range of the score is [0.0, 100.0]. The average score is reported for each translation direction. Higher scores are highlighted in bold.

	Translation Quality	Hallucination				Omission			
		No	Small	Partial	Full	No	Small	Partial	Full
Baseline	11.33%	64.00%	21.00%	11.33%	3.66%	56.00%	25.33%	13.66%	4.33%
+WAP	<b>39.66%</b>	<b>75.66%</b>	<b>17.33%</b>	<b>7.00%</b>	<b>0.00%</b>	<b>80.00%</b>	<b>16.66%</b>	<b>5.33%</b>	<b>0.00%</b>

Table 3: Human evaluation on “zh-en” when N=100. Translation quality is the measured by ratio of examples where WAP beats the baseline. The remaining columns present the ratio of examples in which the corresponding degree of hallucination or omission occurs. Better model is highlighted with bold fonts.

418 still be unreliable. The subset of “N=100” on “zh-  
419 en” is selected. Three volunteers who speak Chi-  
420 nese and English are asked to assess the quality  
421 of the translation and the degree of hallucination  
422 and omission for the baseline and our model, with-  
423 out knowing which model generates the transla-  
424 tions. Table 3 demonstrates the results. In general,  
425 our model generates better translation in 39.66%  
426 of the examples, while the percentage for ALMA  
427 is 11.33%. Furthermore, it is observed that with  
428 DPO on word-alignment preferred data fine-tuning,  
429 the degree of both hallucination and omission de-  
430 creases. Specifically, the percentage of “no hallu-  
431 cination” increases from 64% to 75.66%, and that  
432 of “small, partial, and full hallucination” decreases  
433 accordingly. The decrease in omission is more  
434 distinct, in which the percentage of “no omission”  
435 increase by 24%. Notably, for both hallucination  
436 and omission, the percentage of “full hallucination  
437 and omission” has decreased to 0 for our model.  
438 These results indicate that omission is more fre-  
439 quent than hallucination, and WAP can mitigate  
440 them in LLM-based MT model.

#### 441 4.4 Ablation study

442 In this section, we conduct in-depth investigation  
443 for our word alignment preference, as we use the  
444 same training data as our baseline ALMA, i.e., hu-  
445 man translation, but extra translations from DeepL  
446 and ChatGPT are included to conduct our prefer-

ence data. To investigate where the improvement  
447 comes from, we introduce two variants without  
448 preference tuning to compare with WAP.  
449

- *FT\_reject*: directly fine-tuning ALMA with  
450 the rejected translations in the dataset.  
451
- *FT\_prefer*: directly fine-tuning ALMA with  
452 the preferred translations in the dataset.  
453

The comparison is demonstrated in Figure 4.  
454

**Does the preferred data really better contribute  
455 to the training?** It is observed that *FT\_prefer* sig-  
456 nificantly outperforms *FT\_reject* in both hard and  
457 easy instances. This indicates that our proposed  
458 pipeline ensures that the samples are selected, lead-  
459 ing to better translation quality.  
460

**Is the DPO preference tuning necessary?** Par-  
461 ticularly, the filled area demonstrates the necessity  
462 of preference tuning using DPO. In hard instances  
463 *FT\_prefer* can reach a competitive performance  
464 with a small gap. However, in easy instances,  
465 *FT\_prefer* largely underperforms WAP and ALMA,  
466 which limits the practicality of it. The possible rea-  
467 son for the different performance in the hard and  
468 easy instances is the direct fine-tuning. Directly  
469 fine-tuning on the preferred data without the com-  
470 parison with rejected examples could cause a hard  
471 fitting to the word-aligned preference but ignore  
472 the general translation quality.  
473

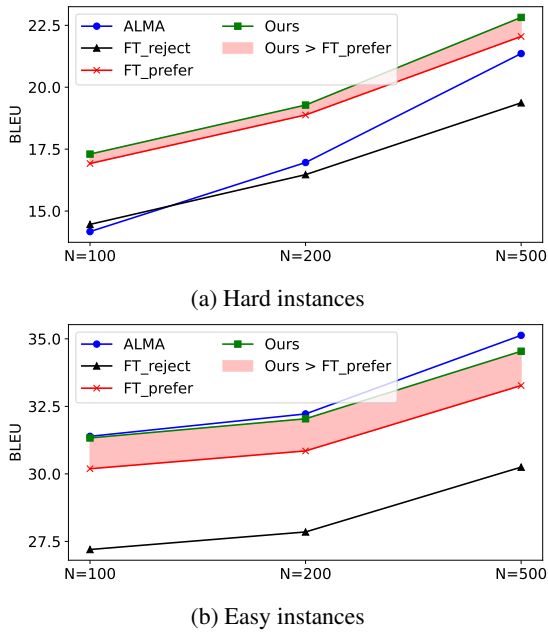


Figure 4: Ablation study. Results in BLEU is demonstrated. Higher BLEU is better. For fair comparison the range of y-axis are the same for hard instances and easy instances. The result in COMET is in the same trend, which can be found in Figure 8 in the Appendix.

## 5 Related work

**Hallucination and omission in MT.** Hallucinations are cases in which the model generates output that is partially or completely unrelated to the source sentence, while omissions are translations that do not include some of the input information (Dale et al., 2023b). Dale et al. (2023a) explore methods that leverage the internal workings of models and external tools, such as cross-lingual sentence similarity and natural language inference models, to detect and mitigate hallucinations in MT. HalOmi (Dale et al., 2023b) introduces an annotated dataset specifically designed to detect hallucinations and omissions. In Figure 1 and §3.2 we use HalOmi as a reference to assess how these two phenomena correlate to the coverage output of the GPT-4 evaluator and the word aligner, respectively. In particular, Yang et al. (2019) introduce using word alignment to reduce omission in MT, which partially inspires our idea.

**Preference tuning for LLMs.** LLMs are capable of completing tasks in the zero-shot or few-shot manner (Radford et al., 2019; Brown et al., 2020). In addition, performance in downstream tasks can also be enhanced by fine-tuning them with instruction datasets (Wei et al., 2022; Chung

et al., 2024; Ouyang et al., 2022). However, acquiring instruction datasets is costly, while obtaining preferences for LLM responses is relatively easier (Rafailov et al., 2024). DPO (Rafailov et al., 2024) directly optimize LLM with preference data by removing an extra reward model. We utilize DPO in this work due to the ease of use and effectiveness. A contemporaneous preference-based MT model ALMA-R (Xu et al., 2024b), introduces contrastive preference optimization to fine-tune LLMs specifically using reference-free MT metrics and human annotation as preference. ALMA-R focuses on improving general LLM-based MT but we attempt to mitigate the hallucination and omission in MT. In addition, our preference data are entirely made automatically, which also draws the difference between ALMA-R and our work.

**Word alignment.** Word-level information is useful in many NLP tasks such as language pre-training (Chi et al., 2021; Wu et al., 2021), cross-lingual sentence embedding (Zhang et al., 2023b; Li et al., 2023; Miao et al., 2024), and particularly for MT (Bahdanau et al., 2015; Tu et al., 2016). Word aligners based on pre-trained language models (Jalili Sabet et al., 2020; Dou and Neubig, 2021; Nagata et al., 2020; Chousa et al., 2020) have outperformed previous ones based on statistical MT (Och and Ney, 2003; Dyer et al., 2013). WSPAlign (Wu et al., 2023) is a pre-trained word aligner outperforming most previous ones, hence we use it in the experiments.

## 6 Conclusion

The problem of hallucination and omission, a long-standing problem in MT, could become more severe when an LLM is used because an LLM itself could hallucinate or omit in nature. In this paper, our aim is to mitigate this problem in LLM-based MT by optimizing the model toward a preference for better word alignment. We construct preference datasets by collecting translations using multiple MT tools and selecting the preference pair with a higher coverage score output by a word aligner. DPO is then utilized to optimize the model towards the word-aligned preference. As evaluation of hallucination and omission is challenging, we design experiments that include selecting hard instances and using GPT-4 to directly predict coverage score, ensuring an effective evaluation, which indicates that the proposed WAP mitigates hallucination and omission, especially in hard instances.



## 550 Limitation

551 The primary limitation of our method stems from  
552 the imperfections of the word alignment model.  
553 Within our approach, it is inevitable to encounter  
554 some alignment errors, which we address through a  
555 filtering method. However, this solution adds com-  
556 plexity and clutter to the method. Additionally, the  
557 effectiveness of our method is diminished for low-  
558 resource language translations due to the limited  
559 number of parallel sentences available. Lastly, our  
560 reliance on the GPT-4 API to evaluate the results  
561 introduces a significant cost factor. We aim to find  
562 a cost-free alternative for this evaluation process in  
563 future work.

## 564 Ethical Statement

565 All datasets and checkpoints used in this paper  
566 are copyright-free for research purposes. Previous  
567 studies are properly cited and discussed. This re-  
568 search aims to improve LLM-based machine trans-  
569 lation models with word alignment preference data,  
570 and the preference is made by an automatic word  
571 aligner. We do not introduce additional bias to par-  
572 ticular communities. We have obtained the consent  
573 of the annotation volunteers for this study.

## 574 References

575 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama  
576 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
577 Diogo Almeida, Janko Altenschmidt, Sam Altman,  
578 Shyamal Anadkat, et al. 2023. Gpt-4 technical report.  
579 *arXiv preprint arXiv:2303.08774*.

580 Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Ben-  
581 gio. 2015. [Neural machine translation by jointly  
582 learning to align and translate](#). In *3rd International  
583 Conference on Learning Representations, ICLR 2015,  
584 San Diego, CA, USA, May 7-9, 2015, Conference  
585 Track Proceedings*.

586 Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wen-  
587 liang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei  
588 Ji, Tiezheng Yu, Willy Chung, et al. 2023. A multi-  
589 task, multilingual, multimodal evaluation of chatgpt  
590 on reasoning, hallucination, and interactivity. In *Pro-  
591 ceedings of the 13th International Joint Conference  
592 on Natural Language Processing and the 3rd Confer-  
593 ence of the Asia-Pacific Chapter of the Association  
594 for Computational Linguistics (Volume 1: Long Pa-  
595 pers)*, pages 675–718.

596 Loïc Barrault, Magdalena Biesialska, Ondřej Bo-  
597 jar, Marta R. Costa-jussà, Christian Federmann,  
598 Yvette Graham, Roman Grundkiewicz, Barry Had-  
599 dow, Matthias Huck, Eric Joanis, Tom Kocmi,  
600 Philipp Koehn, Chi-kiu Lo, Nikola Ljubešić, Christof

Monz, Makoto Morishita, Masaaki Nagata, Toshi-  
aki Nakazawa, Santanu Pal, Matt Post, and Marcos  
Zampieri. 2020. [Findings of the 2020 conference on  
machine translation \(WMT20\)](#). In *Proceedings of  
the Fifth Conference on Machine Translation*, pages  
1–55, Online. Association for Computational Linguis-  
tics.

Ondřej Bojar, Rajen Chatterjee, Christian Federmann,  
Yvette Graham, Barry Haddow, Shujian Huang,  
Matthias Huck, Philipp Koehn, Qun Liu, Varvara  
Logacheva, Christof Monz, Matteo Negri, Matt Post,  
Raphael Rubino, Lucia Specia, and Marco Turchi.  
2017. [Findings of the 2017 conference on machine  
translation \(WMT17\)](#). In *Proceedings of the Second  
Conference on Machine Translation*, pages 169–214,  
Copenhagen, Denmark. Association for Computa-  
tional Linguistics.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
Askell, Sandhini Agarwal, Ariel Herbert-Voss,  
Gretchen Krueger, Tom Henighan, Rewon Child,  
Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens  
Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma-  
teusz Litwin, Scott Gray, Benjamin Chess, Jack  
Clark, Christopher Berner, Sam McCandlish, Alec  
Radford, Ilya Sutskever, and Dario Amodei. 2020.  
[Language models are few-shot learners](#). In *Ad-  
vances in Neural Information Processing Systems*,  
volume 33, pages 1877–1901. Curran Associates,  
Inc.

Zewen Chi, Li Dong, Bo Zheng, Shaohan Huang, Xian-  
Ling Mao, Heyan Huang, and Furu Wei. 2021. [Im-  
proving pretrained cross-lingual language models via  
self-labeled word alignment](#). In *Proceedings of the  
59th Annual Meeting of the Association for Compu-  
tational Linguistics and the 11th International Joint  
Conference on Natural Language Processing (Vol-  
ume 1: Long Papers)*, pages 3418–3430, Online. As-  
sociation for Computational Linguistics.

Katsuki Chousa, Masaaki Nagata, and Masaaki Nishino.  
2020. [SpanAlign: Sentence alignment method based  
on cross-language span prediction and ILP](#). In *Pro-  
ceedings of the 28th International Conference  
on Computational Linguistics*, pages 4750–4761,  
Barcelona, Spain (Online). International Committee  
on Computational Linguistics.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret  
Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi  
Wang, Mostafa Dehghani, Siddhartha Brahma, et al.  
2024. Scaling instruction-finetuned language models.  
*Journal of Machine Learning Research*, 25(70):1–53.

Marta R Costa-jussà, James Cross, Onur Çelebi, Maha  
Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe  
Kalbassi, Janice Lam, Daniel Licht, Jean Maillard,  
et al. 2022. No language left behind: Scaling  
human-centered machine translation. *arXiv preprint  
arXiv:2207.04672*.

659	David Dale, Elena Voita, Loic Barrault, and Marta R. Costa-jussà. 2023a. <a href="#">Detecting and mitigating hallucinations in machine translation: Model internal workings alone do well, sentence similarity Even better</a> . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 36–50, Toronto, Canada. Association for Computational Linguistics.	
660		
661		
662		
663		
664		
665		
666		
667	David Dale, Elena Voita, Janice Lam, Prangthip Hansanti, Christophe Ropers, Elahe Kalbassi, Cynthia Gao, Loic Barrault, and Marta Costa-jussà. 2023b. <a href="#">HalOmi: A manually annotated benchmark for multilingual hallucination and omission detection in machine translation</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 638–653, Singapore. Association for Computational Linguistics.	
668		
669		
670		
671		
672		
673		
674		
675		
676	Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Jason Weston. 2023. <a href="#">Chain-of-verification reduces hallucination in large language models</a> . <i>ArXiv</i> , abs/2309.11495.	
677		
678		
679		
680		
681	Zi-Yi Dou and Graham Neubig. 2021. Word alignment by fine-tuning embeddings on parallel corpora. In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 2112–2128.	
682		
683		
684		
685		
686	Chris Dyer, Victor Chahuneau, and Noah A Smith. 2013. A simple, fast, and effective reparameterization of ibm model 2. In <i>Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 644–648.	
687		
688		
689		
690		
691		
692	Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. <a href="#">Language-agnostic BERT sentence embedding</a> . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 878–891, Dublin, Ireland. Association for Computational Linguistics.	
693		
694		
695		
696		
697		
698		
699	Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. <a href="#">SimCSE: Simple contrastive learning of sentence embeddings</a> . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 6894–6910, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	
700		
701		
702		
703		
704		
705		
706	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>arXiv preprint arXiv:2302.09210</i> .	
707		
708		
709		
710		
711		
712	Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2021. Lora: Low-rank adaptation of large language models. In <i>International Conference on Learning Representations</i> .	
713		
714		
715		
716		
	Masoud Jalili Sabet, Philipp Dufter, François Yvon, and Hinrich Schütze. 2020. SimAlign: High quality word alignments without parallel training data using static and contextualized embeddings. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 1627–1643, Online. Association for Computational Linguistics.	717
		718
		719
		720
		721
		722
		723
	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. <i>Advances in neural information processing systems</i> , 35:22199–22213.	724
		725
		726
		727
		728
	Ziheng Li, Shaohan Huang, Zihan Zhang, Zhi-Hong Deng, Qiang Lou, Haizhen Huang, Jian Jiao, Furu Wei, Weiwei Deng, and Qi Zhang. 2023. <a href="#">Dual-alignment pre-training for cross-lingual sentence embedding</a> . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3466–3478, Toronto, Canada. Association for Computational Linguistics.	729
		730
		731
		732
		733
		734
		735
		736
	Zhongtao Miao, Qiyu Wu, Kaiyan Zhao, Zilong Wu, and Yoshimasa Tsuruoka. 2024. Enhancing cross-lingual sentence embedding for low-resource languages with word alignment. <i>arXiv preprint arXiv:2404.02490</i> .	737
		738
		739
		740
		741
	Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D. Manning, and Chelsea Finn. 2023. <a href="#">Detectgpt: Zero-shot machine-generated text detection using probability curvature</a> . In <i>International Conference on Machine Learning</i> .	742
		743
		744
		745
		746
	Masaaki Nagata, Katsuki Chousa, and Masaaki Nishino. 2020. A supervised word alignment method based on cross-language span prediction using multilingual bert. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 555–565.	747
		748
		749
		750
		751
		752
	Franz Josef Och and Hermann Ney. 2003. A systematic comparison of various statistical alignment models. <i>Computational linguistics</i> , 29(1):19–51.	753
		754
		755
	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. 2022. <a href="#">Training language models to follow instructions with human feedback</a> . In <i>Advances in Neural Information Processing Systems</i> , volume 35, pages 27730–27744. Curran Associates, Inc.	756
		757
		758
		759
		760
		761
		762
		763
		764
		765
	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. <a href="#">Bleu: a method for automatic evaluation of machine translation</a> . In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics</i> , pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.	766
		767
		768
		769
		770
		771
		772

773	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	United Arab Emirates. Association for Computa-	830
774	Dario Amodei, Ilya Sutskever, et al. 2019. Language	tional Linguistics.	831
775	models are unsupervised multitask learners. <i>OpenAI</i>		
776	<i>blog</i> , 1(8):9.		
777	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christo-	Qiyu Wu, Chen Xing, Yatao Li, Guolin Ke, Di He, and	832
778	pher D Manning, Stefano Ermon, and Chelsea Finn.	Tie-Yan Liu. 2021. <a href="#">Taking notes on the fly helps</a>	833
779	2024. Direct preference optimization: Your language	<a href="#">language pre-training</a> . In <i>International Conference</i>	834
780	model is secretly a reward model. <i>Advances in Neu-</i>	<i>on Learning Representations</i> .	835
781	<i>ral Information Processing Systems</i> , 36.		
782	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Has-	836
783	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	san Awadalla. 2024a. <a href="#">A paradigm shift in machine</a>	837
784	Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti	<a href="#">translation: Boosting translation performance of</a>	838
785	Bhosale, et al. 2023. Llama 2: Open founda-	<a href="#">large language models</a> . In <i>The Twelfth International</i>	839
786	tion and fine-tuned chat models. <i>arXiv preprint</i>	<i>Conference on Learning Representations</i> .	840
787	<i>arXiv:2307.09288</i> .		
788	Zhaopeng Tu, Zhengdong Lu, Yang Liu, Xiaohua Liu,	Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan,	841
789	and Hang Li. 2016. Modeling coverage for neural	Lingfeng Shen, Benjamin Van Durme, Kenton Mur-	842
790	machine translation. In <i>Proceedings of the 54th An-</i>	ray, and Young Jin Kim. 2024b. Contrastive pref-	843
791	<i>annual Meeting of the Association for Computational</i>	erence optimization: Pushing the boundaries of llm	844
792	<i>Linguistics (Volume 1: Long Papers)</i> , pages 76–85.	performance in machine translation. <i>arXiv preprint</i>	845
793	Lewis Tunstall, Edward Beeching, Nathan Lambert,	<i>arXiv:2401.08417</i> .	846
794	Nazneen Rajani, Kashif Rasul, Younes Belkada,	Zonghan Yang, Yong Cheng, Yang Liu, and Maosong	847
795	Shengyi Huang, Leandro von Werra, Cl��mentine	Sun. 2019. <a href="#">Reducing word omission errors in neural</a>	848
796	Fourrier, Nathan Habib, Nathan Sarrazin, Omar San-	<a href="#">machine translation: A contrastive learning approach</a> .	849
797	seviero, Alexander M. Rush, and Thomas Wolf. 2023.	In <i>Proceedings of the 57th Annual Meeting of the As-</i>	850
798	<a href="#">Zephyr: Direct distillation of lm alignment</a> . <i>ArXiv</i> ,	<i>sociation for Computational Linguistics</i> , pages 6191–	851
799	abs/2310.16944.	6196, Florence, Italy. Association for Computational	852
800	Jannis Vamvas and Rico Sennrich. 2022. <a href="#">As little as</a>	Linguistics.	853
801	<a href="#">possible, as much as necessary: Detecting over- and</a>	Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu,	854
802	<a href="#">undertranslations with contrastive conditioning</a> . In	Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang,	855
803	<i>Proceedings of the 60th Annual Meeting of the As-</i>	Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei	856
804	<i>sociation for Computational Linguistics (Volume 2:</i>	Bi, Freda Shi, and Shuming Shi. 2023a. <a href="#">Siren’s song</a>	857
805	<i>Short Papers)</i> , pages 490–500, Dublin, Ireland. As-	<a href="#">in the ai ocean: A survey on hallucination in large</a>	858
806	sociation for Computational Linguistics.	<a href="#">language models</a> . <i>ArXiv</i> , abs/2309.01219.	859
807	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu,	Zhen-Ru Zhang, Chuanqi Tan, Songfang Huang, and	860
808	Adams Wei Yu, Brian Lester, Nan Du, Andrew M.	Fei Huang. 2023b. Veco 2.0: Cross-lingual language	861
809	Dai, and Quoc V Le. 2022. <a href="#">Finetuned language mod-</a>	model pre-training with multi-granularity contrastive	862
810	<a href="#">els are zero-shot learners</a> . In <i>International Confer-</i>	learning. <i>arXiv preprint arXiv:2304.08205</i> .	863
811	<i>ence on Learning Representations</i> .		
812	Xiang Wei, Xingyu Cui, Ning Cheng, Xiaobin Wang,	Kaiyan Zhao, Qiyu Wu, Xin-Qiang Cai, and Yoshimasa	864
813	Xin Zhang, Shen Huang, Pengjun Xie, Jinan Xu,	Tsuruoka. 2024. <a href="#">Leveraging multi-lingual positive</a>	865
814	Yufeng Chen, Meishan Zhang, Yong Jiang, and Wen-	<a href="#">instances in contrastive learning to improve sentence</a>	866
815	juan Han. 2023. <a href="#">Zero-shot information extraction via</a>	<a href="#">embedding</a> . In <i>Proceedings of the 18th Conference of</i>	867
816	<a href="#">chatting with chatgpt</a> . <i>ArXiv</i> , abs/2302.10205.	<i>the European Chapter of the Association for Compu-</i>	868
817	Qiyu Wu, Masaaki Nagata, and Yoshimasa Tsuruoka.	<i>tational Linguistics (Volume 1: Long Papers)</i> , pages	869
818	2023. <a href="#">WSPAlign: Word alignment pre-training via</a>	976–991, St. Julian’s, Malta. Association for Com-	870
819	<a href="#">large-scale weakly supervised span prediction</a> . In	putational Linguistics.	871
820	<i>Proceedings of the 61st Annual Meeting of the As-</i>		
821	<i>sociation for Computational Linguistics (Volume 1:</i>		
822	<i>Long Papers)</i> , pages 11084–11099, Toronto, Canada.		
823	Association for Computational Linguistics.		
824	Qiyu Wu, Chongyang Tao, Tao Shen, Can Xu, Xiubo		
825	Geng, and Daxin Jiang. 2022. <a href="#">PCL: Peer-contrastive</a>		
826	<a href="#">learning with diverse augmentations for unsupervised</a>		
827	<a href="#">sentence embeddings</a> . In <i>Proceedings of the 2022</i>		
828	<i>Conference on Empirical Methods in Natural Lan-</i>		
829	<i>guage Processing</i> , pages 12052–12066, Abu Dhabi,		



## A Experimental setup

The implementation from alignment-handbook<sup>11</sup> is used for the training of DPO. The learning rate is searched based on performance on development set and set to  $5e-6$ . LoRA (Hu et al., 2021) is used.  $r$  is set as 16 and  $\beta$  is set as 0.1. We train the model for 1 epoch and fix the random seed to 42. The model is trained on  $4 \times$  Nvidia A100 80G and the total batch size is 64. For evaluation, we use the implementation of ALMA<sup>12</sup> to calculate the BLEU and COMET scores.

## B Details of dataset

Figure 5 presents the varying proportions of “chosen” and “rejected” preference pairs from three sources: ChatGPT, DeepL, and Human. The figure indicates that the majority of the “chosen” translations originate from ChatGPT, while a significant portion of human-written translations are “rejected”. This observation supports the conclusion that human-written translations can also exhibit quality issues, as discussed in ALMA-R (Xu et al., 2024b). Examples in our constructed preference dataset are presented in §C.1.

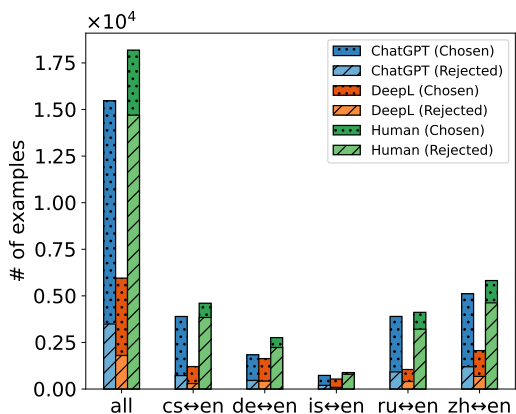


Figure 5: This figure illustrates the proportions of “chosen” and “rejected” preference pairs derived from three sources: ChatGPT, DeepL and Human. “all” represents the overall proportion for the aggregated dataset.  $xx \leftrightarrow en$  is the subset pair of English and another language. Particularly, Google Translate is used for  $is \leftrightarrow en$  as an alternative to DeepL.

**Prompt for Translation**

You are a helpful assistant that translates {SOURCE\_LANG} sentences to {TARGET\_LANG} sentences.

{TEXT}

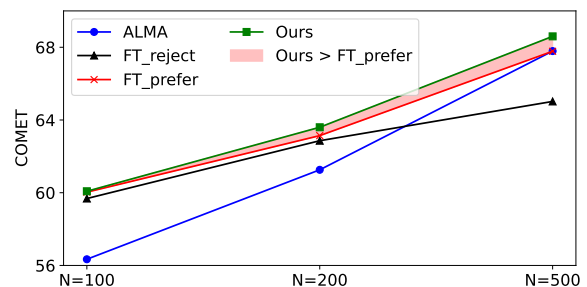
Figure 6: The prompt of ChatGPT that we use to translate sentences.

**Prompt of Coverage Calculation**

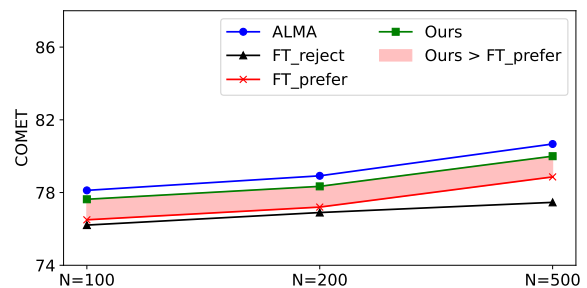
Given a source sentence {SOURCE\_TEXT} in {SOURCE\_LANGUAGE} and a translation {TRANSLATION} in {TARGET\_LANGUAGE}, does the translation has hallucination or omission to the source contents?

\*\*\*You MUST answer with only a coverage percentage score\*\*\*.

Figure 7: Prompt to calculate the coverage score.



(a) Hard instances



(b) Easy instances

Figure 8: Ablation study. Results in COMET is demonstrated. Higher COMET is better. For fair comparison the range of y-axis are the same for hard instances and easy instances. Refer to §4.4 for discussion.



Example 1 (Chinese-English)		Coverage Score
source	“我想，在考虑重播时，可以解决这个问题”，Coker 说道。	–
chosen (gpt-3.5)	"I think, when considering replay, this issue can be resolved," Coker said.	94.03
rejected (human)	"<<<I think that when I think about>>> the replay, <<<I think that>>> we can probably work it out," Coker said.	79.87
Example 2 (Chinese-English)		Coverage Score
source	<<<富勒>>>在政变图谋失败后	–
chosen (deepl)	<<<Fuller>>> after the failed coup attempt	83.76
rejected (human)	After the failure of the attempted coup,	59.59
Example 3 (English-Chinese)		Coverage Score
source	<<<Originally a one-bedroom property with a convoluted layout - you had to walk through the kitchen to get to the bedroom>>> - Joanne wanted to add storage space and a mezzanine to make the most of the generous ceiling height.'	–
chosen (gpt-3.5)	<<<最初是一个一居室的房产，布局错综复杂- 你必须穿过厨房才能到达卧室>>> - 然而乔安妮想要增加存储空间和一个夹层，以充分利用宽敞的天花板高度。	83.76
rejected (deepl)	乔安妮希望增加储藏空间和一个夹层，充分利用宽敞的天花板高度。	69.97

Table 4: Examples in the preference dataset. The hallucination in rejected examples and omission in the source sentence are highlighted with <<< >>>. The corresponding contents that are omitted in the rejected example are highlighted with <<< >>> in the chosen example. The coverage is calculated by word aligner, refer to §2 for details.

## C Example analysis

### C.1 Examples of the preference dataset

Table 4 includes three examples in our dataset, in which the source sentence, the chosen and rejected translations are shown. Refer to §B for a detailed construction of the dataset. **Example 1:** the rejected translation is from human annotation, in which it repeats the term of “I think” unnaturally. The possible reason could be the resource of the parallel data, e.g., direct collection from transcriptions. **Example 2:** “Fuller” is omitted by human annotation while translated by DeepL. **Example 3:** the chosen translation is from gpt-3.5-turbo that completely translates the source sentence. In contrast, the translation by DeepL omits the first half.

### C.2 Translation examples

Table 5 shows illustrative comparison between translations from the baseline and our model. **Example 1:** “in HBO’s ‘The Gilded Age’” in the source sentence is omitted by the baseline. In

<sup>11</sup><https://github.com/huggingface/alignment-handbook>

<sup>12</sup><https://github.com/felixxu/ALMA>

contrast, our model successfully translate the corresponding part into Chinese. **Example 2:** the baseline generates “卡扣 (fastening)” infinitely in translation. This type of hallucination also occurs in other LLM applications, which emphasizes the need to address the hallucination issue in LLM-based MT models. **Example 3:** “等到什么时候 (when to wait)” is omitted by the baseline model while our model translate that into “how long I have to wait” properly.

## D Specific results

Table 6 shows the numeric results in Figure 3, in which boxes on a blue background highlight the cases where our model outperforms the baseline by a margin  $> 1.0$ , and the boxes in red are the opposite. Boxes without background indicate the cases when our model and the baseline have competitive performance where the margin  $< 1.0$ .

In addition to the main findings in §4.1 that our model generally performs better in harder instances, from the results it can also be observed that our model particularly performs worse on “en-is” than in other translation directions. The reason could be



Model-Metric	de-en	cs-en	is-en	zh-en	ru-en	en-de	en-cs	en-is	en-zh	en-ru	Avg.
<i>N=100</i>											
<i>Easy instances</i>											
ALMA-BLEU	31.38	45.79	38.14	25.64	41.25	32.09	31.95	27.57	40.05	29.37	31.39
Ours-BLEU	32.50	46.32	40.13	25.23	40.80	31.22	31.55	26.00	39.55	29.01	31.33
ALMA-COMET	85.57	87.71	87.82	81.38	86.26	86.84	90.90	87.61	87.14	88.80	78.12
Ours-COMET	85.50	87.67	87.71	81.24	86.17	86.02	89.84	85.80	86.39	87.89	77.63
<i>Hard instances</i>											
ALMA-BLEU	12.25	29.49	21.72	1.95	15.73	15.71	12.79	17.51	14.59	15.45	14.17
Ours-BLEU	15.56	35.93	27.72	4.62	19.77	16.15	16.67	17.13	19.49	15.54	17.30
ALMA-COMET	62.73	67.08	72.62	49.94	62.64	58.50	60.80	70.02	59.07	62.31	56.34
Ours-COMET	65.98	71.16	75.12	58.99	67.19	60.90	67.90	71.57	62.03	65.16	60.08
<i>N=200</i>											
<i>Easy instances</i>											
ALMA-BLEU	31.96	47.11	39.94	26.22	42.13	32.50	32.75	28.54	41.08	30.22	32.22
Ours-BLEU	33.10	47.41	41.60	25.79	41.43	31.52	32.20	26.91	40.48	29.79	32.04
ALMA-COMET	86.34	88.61	88.72	82.31	87.02	87.76	91.85	88.67	87.97	89.67	78.92
Ours-COMET	86.16	88.40	88.43	81.98	86.89	86.75	90.77	86.94	87.12	88.73	78.34
<i>Hard instances</i>											
ALMA-BLEU	17.46	30.39	24.17	6.00	20.03	19.11	14.83	19.02	18.61	15.43	16.96
Ours-BLEU	19.31	35.04	29.25	7.55	23.70	19.96	18.16	18.29	21.52	15.95	19.28
ALMA-COMET	67.24	71.82	76.62	57.84	67.59	64.30	67.13	74.56	65.46	67.59	61.26
Ours-COMET	69.85	74.82	78.52	63.87	70.22	66.77	70.37	74.13	67.50	68.78	63.60
<i>N=500</i>											
<i>Easy instances</i>											
ALMA-BLEU	34.36	50.81	46.92	28.50	45.16	34.61	35.28	31.79	43.91	32.13	35.13
Ours-BLEU	35.33	50.59	47.25	27.82	44.16	33.25	34.07	30.00	42.92	31.67	34.54
ALMA-COMET	88.08	90.54	91.04	84.29	88.62	89.59	93.66	91.08	89.79	91.47	80.67
Ours-COMET	87.80	90.10	90.50	83.86	88.40	88.55	92.48	89.57	88.79	90.61	80.00
<i>Hard instances</i>											
ALMA-BLEU	21.31	35.46	28.66	13.08	25.4	22.53	19.82	22.52	24.81	19.78	21.36
Ours-BLEU	23.09	37.91	32.66	14.04	27.32	22.89	22.38	21.32	26.58	19.78	22.82
ALMA-COMET	73.56	78.24	81.55	67.07	74.39	72.74	76.38	80.61	73.38	75.29	67.79
Ours-COMET	74.77	79.75	82.41	69.56	75.63	73.24	77.34	79.19	74.12	74.97	68.60
<i>Overall performance, i.e., N=infinite when all instances are included.</i>											
ALMA-BLEU	30.73	44.68	36.46	24.15	40.37	31.37	31.12	26.67	39.05	28.76	30.46
Ours-BLEU	31.93	45.60	38.85	23.94	40.09	30.64	30.91	25.22	38.76	28.43	30.59
ALMA-COMET	84.42	86.29	86.30	79.70	85.09	85.45	89.42	85.85	85.76	87.50	76.83
Ours-COMET	84.50	86.53	86.45	80.05	85.22	84.78	88.75	84.38	85.19	86.77	76.59

Table 6: Specific results on 10 translation directions. The size of models are 13B. BLEU and COMET are reported. Cells where the difference is larger than 1.0 are highlighted with colored background. Blue indicates our model outperforms ALMA and red indicates the opposite.