

# Reducing Translationese via Iterative Translation Refinement with Large Language Models

Anonymous ARR submission

## Abstract

Translations created by machines or humans can suffer from *translationese*—an awkward or unnatural output due to the translation process. We argue that the advent of large language models offers a means to mitigate translationese via iterative refinement, which is infeasible for conventional encoder-decoder models. Our experiments show that refinement reduces string-based metric scores, but neural metrics suggest comparable or improved quality. Human evaluations demonstrate that translationese is lessened compared to initial translations and even human references, while maintaining quality. Ablation studies underscore the importance of anchoring the refinement to the source and a reasonable seed translation. We also discuss current challenges in measuring translationese.

## 1 Introduction

Large language models (LLMs), e.g. generative pre-trained Transformers (GPT), have made notable advancements in natural language processing (Radford et al., 2019; Brown et al., 2020; Kaplan et al., 2020; Ouyang et al., 2022). In machine translation (MT), where the convention is to use an encoder-decoder architecture to deal with source and target sentences respectively (Bahdanau et al., 2015; Vaswani et al., 2017), recent papers have examined the feasibility of LLM prompting (Vilar et al., 2023; Zhang et al., 2023; Hendy et al., 2023).

Prior research combining LLMs and MT did not extensively explore the phenomenon of “translationese”, which refers to a translation that does not read as naturally as an original text. This is due to both source language interference and the translation process itself (Gellerstam, 1986; Baker, 1996; Teich, 2003). It appears in various stages of MT: data (Riley et al., 2020), machine outputs (Freitag et al., 2020a), human post-edited translations (“post-edited”, Toral, 2019), and even human references (Freitag et al., 2020b). Moreover, MT

training typically relies on parallel data, which in the first place would come from human translators, or in some cases, other MT systems. Thus the data naturally exhibit translationese patterns to a certain extent, which in turn propagates into MT training. Even LLMs might be translationese-prone as their translation power is associated with implicit bilingual signals (Briakou et al., 2023). These imperfections not only damage translation performance but also undermine the evaluation process (Toral et al., 2018; Graham et al., 2020).

Going beyond the current translation paradigm, we propose a simple way to refine translations iteratively with LLMs, building on automatic post-editing which imitates human corrections (APE, Knight and Chander, 1994; Chatterjee et al., 2018). We prompt an LLM for a translation and feed the source-translation pair back for a refined translation in multiple rounds. Our method offers two insights for combating translationese: (1) Unlike translation or post-editing models, LLMs have been exposed to datasets that are orders of magnitude larger and less translationese. We thus indirectly incorporate genuine texts to pursue natural translations. (2) Our prompting mechanism allows for iterative and arbitrary rephrasing compared to APE which is limited to token-level error correction without style editing (Ive et al., 2020).

Empirical results show that the refinement process introduces significant textual changes reflected by the drop in BLEU and chrF++, but attains similar or higher COMET scores compared to initial translations. Native speakers prefer the refined outputs in terms of reduced translationese, which is more prevalent in GPT translations and even the human references. Referenced-based human evaluation confirms that such gains are made without sacrificing general quality. As corroborated by recent works, these are challenging to capture by automatic metrics like BLEU or COMET alone (Freitag et al., 2019, 2022).

Mode	Prompt
<i>Translate</i>	Source: $\{\text{source}\}$ Please give me a translation in $\{\text{lang}\}$ without any explanation.
<i>Refine</i>	Source: $\{\text{source}\}$ Translation: $\{\text{prev\_translation}\}$ Please give me a better $\{\text{lang}\}$ translation without any explanation.
<i>Refine</i> <sub>Contrast</sub>	Source: $\{\text{source}\}$ Bad translation: $\{\text{prev\_translation}\}$ Please give me a better $\{\text{lang}\}$ translation without any explanation.
<i>Refine</i> <sub>Random</sub>	Source: $\{\text{source}\}$ Bad translation: $\{\text{random\_target}\}$ if first-round, else $\{\text{prev\_translation}\}$ Please give me a better $\{\text{lang}\}$ translation without any explanation.
<i>Paraphrase</i>	Sentence: $\{\text{prev\_translation}\}$ Please give me a paraphrase in $\{\text{lang}\}$ without any explanation.

Table 1: Prompts used in our work, where  $\{\text{variable}\}$  is substituted with its corresponding content.

## 2 Methodology

Having an input source sentence  $x$  and an optimizable model  $\theta_{mt}$ , the process to obtain a translation  $y$  can be modelled as  $y = \text{argmax}_y P(y|x, \theta_{mt})$ . Next, an automatic post-editor  $\theta_{ape}$  creates a refined translation  $y'$  through  $y' = \text{argmax}_{y'} P(y'|x, y, \theta_{ape})$ . Conventional translation or automatic post-editing models are trained on  $(x, y)$  or  $(x, y, y')$  data pairs.

Since translationese naturally arises during the translation process, we hypothesize that we can alleviate it via refinement using LLMs to bypass the direct translation formality. Our study uses zero-shot prompting by affixing a task description to form a prompt  $p$  and querying an LLM  $\theta_{LLM}$  to elicit a response (Brown et al., 2020). We introduce five prompts in our study:

1. *Translate*: this queries for a translation of a source input, extending the translation process with a prompt  $p$ :  $y = \text{argmax}_y P(y|p, x, \theta_{LLM})$
2. *Refine*: similar to APE, the LLM is given the source sentence and the previous translation to produce a better translation  $y' = \text{argmax}_{y'} P(y'|p, x, y, \theta_{LLM})$ .
3. *Refine*<sub>Contrast</sub>: as a contrasting prompt to the above, we insert the word “bad” to hint that the previously translated text is unwanted, regardless of its actual quality.
4. *Refine*<sub>Random</sub>: same prompt as *Refine*<sub>Contrast</sub>, but in the first iteration, a random sentence is fed instead of a translation to imitate a genuinely “bad translation”.
5. *Paraphrase*: to ablate the translation process, we prompt to rephrase a translation without feeding the source sentence  $x$ :  $y'' = \text{argmax}_{y''} P(y''|p, y, \theta_{LLM})$ .

Our study proposes to iteratively call the refine-

ment prompts, where the source stays the same but the previous translation is updated with the latest, to understand how quality changes. Through ab-  
lative prompts, we can analyse to what degree the  
source input and seed translations are helpful. The  
exact prompt texts are displayed in Table 1.

## 3 Experiments

### 3.1 Data and model details

We experiment with language pairs from the trans-  
lation tasks hosted at WMT 2021 and 2022 (Farhad  
et al., 2021; Kocmi et al., 2022). In total, we tested  
seven translation directions: English $\leftrightarrow$ {German,  
Chinese}, German $\rightarrow$ French, English $\rightarrow$ Japanese,  
and Ukrainian $\rightarrow$ Czech. We directly benchmark  
on the test sets, and in situations where multiple  
references are available, we use human reference  
“A” released by the WMT organizers.

We experiment with GPT-3.5, a powerful API  
from OpenAI that can be accessed by all users.<sup>1</sup>  
As the API is very slow to query, we randomly  
sample 200 instances from the official test set to  
form our own test. Similar to the black-box con-  
dition in APE, we do not keep the query history,  
in order to prevent an LLM from seeing that the  
previous translation is produced by itself. Overall,  
translation refinement is iterated four times.

### 3.2 Evaluation setup

We consider four automatic metrics: string-based  
BLEU (Papineni et al., 2002) and chrF++ (Popović,  
2017), as well as embedding-based COMET<sub>DA</sub> and  
COMET<sub>QE</sub> (Rei et al., 2020). The difference be-  
tween DA and QE versions is that COMET<sub>DA</sub> re-

<sup>1</sup>We accessed gpt-3.5-turbo which has training data up to Sep 2021, so it should not have seen WMT 2021 or 2022 test references. Nevertheless, our findings are mostly drawn from reference-free metrics and human evaluation.

	WMT21 de→en		WMT21 en→de		WMT21 zh→en		WMT21 en→zh		WMT22 de→fr		WMT22 en→ja		WMT22 uk→cs	
	BLEU	COMET <sub>QE</sub>	BLEU	COMET <sub>QE</sub>	BLEU	COMET <sub>QE</sub>	BLEU	COMET <sub>QE</sub>	BLEU	COMET <sub>QE</sub>	BLEU	COMET <sub>QE</sub>	BLEU	COMET <sub>QE</sub>
Reference <sub>A</sub>	-	.0919	-	.1127	-	.0708	-	.0956	-	.0772	-	.1345	-	.1273
Translate	<b>30.90</b>	.1128	<b>25.39</b>	.1083	<b>25.64</b>	.0867	<b>29.28</b>	.0761	<b>36.25</b>	.0807	<b>23.00</b>	.1255	<b>29.91</b>	.1173
Refine	23.14	.1116	22.35	<b>.1153</b>	20.26	.0921	28.26	.0870	32.47	<b>.0851</b>	22.63	<b>.1305</b>	28.60	<b>.1183</b>
Refine <sub>Contrast</sub>	22.88	<b>.1162</b>	22.54	.0929	24.81	<b>.1132</b>	<b>29.28</b>	<b>.0881</b>	33.12	.0805	22.82	.1282	28.90	.1151
Refine <sub>Random</sub>	18.83	.0770	19.36	.0832	24.24	.1022	25.71	.0763	-	-	-	-	-	-
Paraphrase	11.01	.0919	13.60	.1006	12.76	.0885	21.95	.0716	16.06	.0682	17.69	.1086	13.59	.0969

Table 2: Automatic scores of different strategies on translation directions from WMT 2021 and 2022 news translation.

quires a source, a translation, and a human reference, whereas COMET<sub>QE</sub> is reference-free.<sup>2</sup>

Although these metrics are widely used to measure translation quality, there is no effective measure for translationese thus far. Freitag et al. (2020a) hint that too high a single-reference BLEU cannot imply high quality; we see it as an indicator of text variations from the reference. Further, we argue that since human references could be translationese-prone, evaluation should not anchor to them. We hence rely on the reference-free COMET<sub>QE</sub>, which correlates well with human judgements (Freitag et al., 2022). We report BLEU and COMET<sub>QE</sub> scores in the main content but also attach chrF++ and COMET<sub>DA</sub> in Appendix A.

### 3.3 Refinement results

**WMT21** We first experiment with en↔de and en↔zh from WMT21, and display results in Table 2. For iterative experiments, the best iteration is picked according to COMET<sub>QE</sub>. We observe that the refined translations record a drastic drop in string-based metrics compared to initial translations, indicating lexical and structural variations. In terms of COMET<sub>QE</sub>, refined outputs surpass all initial GPT translations, with substantial improvement for into-English directions. The ablative *Paraphrase* method sees a decline in all metrics, suggesting the importance of feeding the input as an anchor during iterations to prevent semantic drift.

To investigate the behaviour of different refinement strategies, we plot BLEU, COMET<sub>DA</sub>, and COMET<sub>QE</sub> at different iterations in Appendix C Figure 2. We see that *Refine* and *Refine*<sub>Contrast</sub> usually attain their best after the first iteration, but in almost all *Paraphrase* experiments, scores decrease monotonically, indicating that semantics drift away as paraphrasing iterates. Moreover, *Refine*<sub>Random</sub> results start low, gradually catch up, but never reach

<sup>2</sup>BLEU and chrF++ are as in the sacrebleu toolkit (Post, 2018). For COMET, we use wmt-2022-da and wmt-2021-qe-da respectively. We document details in Appendix E.

as high as *Refine* or *Refine*<sub>Contrast</sub>. This means that iterative refinement is indeed useful in fixing translations, but starting with a reasonable translation is also crucial for obtaining a strong result.

**WMT22** For non-English translation, we pick three directions from WMT22. Since *Refine*<sub>Random</sub> results are not desirable for WMT21, we omit experiments with this. We find that *Refine* works best, obtaining higher COMET<sub>QE</sub> than vanilla translations and *Refine*<sub>Contrast</sub>. Also, the reduction in string-based scores becomes less obvious, which might be attributed to seed GPT translations in lesser-resourced languages being lower in quality.

**WMT system refinement** Finally, in addition to translation refinement from GPT-3.5 itself, we also apply our refinement calls to outputs from conventional MT systems and human translators. These translations can represent genuine errors, if any, introduced during the translation process. We experiment with seven different submissions in the WMT 2021 German-to-English news translation track as a starting point. Due to the space constraint, we introduce the systems and report automatic metric scores in Appendix B.

A pattern similar to previous GPT refinement is noticed. For five out of seven WMT entries, the refinement strategy reaches a higher COMET<sub>QE</sub> score, surprisingly, with up to one-third drop in BLEU. *Refine*<sub>Contrast</sub> in all but one system surpass *Refine*, and without the initial translation, *Paraphrase* iterations record the lowest scores compared to the original submissions and refinements.

## 4 Human Evaluation

String-based and neural scores are observed to vary in opposite directions, which may suggest changes in texts without affecting meaning (Freitag et al., 2020b). As there is no automatic metric for translationese, we set up human evaluations to measure two characteristics in the refined translations: the translationese degree and overall quality.

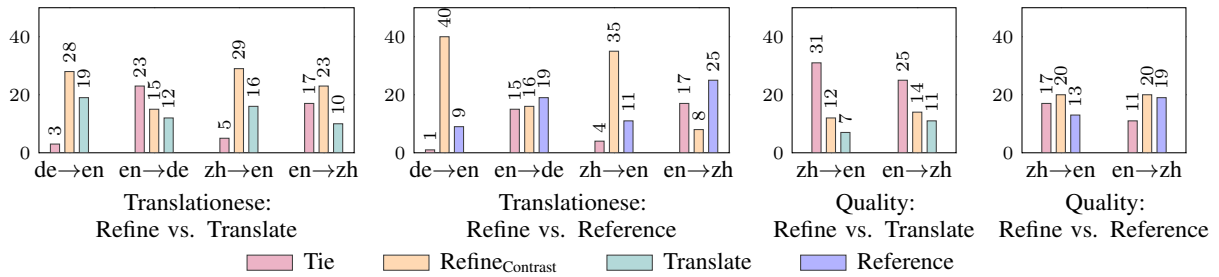


Figure 1: Human preferences on reduced translationese (source-free, left) and overall quality (source-based, right).

## 4.1 Translationese

Since the term “translationese” is not commonly known, we mimic an established work on translationese detection (Lembersky et al., 2012). We present native speakers with two translations but without the source sentence; then we ask “Please choose the translation that is more fluent, natural, and reflecting better use of  $\{\text{language}\}$ ”. The evaluators can select one of the two translations, or a “tie” if they consider both equally (un)natural. We conduct such pairwise evaluation to compare the first-round output from *Refine*<sub>Contrast</sub> against human references, as well as against *Translate* separately.

We evaluate 50 samples from  $\text{en} \leftrightarrow \text{de}$  and  $\text{en} \leftrightarrow \text{zh}$  experiments in Section 3.3, and report results in Figure 1 (left). Native speakers prefer *Refine*<sub>Contrast</sub> to vanilla *Translate* in all four directions, and even favour *Refine*<sub>Contrast</sub> over human references when translating into English. The results demonstrate that our simple strategy enhances the naturalness of GPT translations, and that human references could be more translationese than GPT outputs for into-English directions, thus making reference-based metrics like BLEU or COMET<sub>DA</sub> less reliable.

## 4.2 Overall quality

We then evaluate for general translation quality. In this setup, a source sentence and two translations are given to an evaluator who is fluent in both languages. They are asked to pick the translation with better quality or indicate a tie. We only evaluated two translation directions, English to and from Chinese, due to the limited availability of bilingual speakers. Similar to the previous evaluation, we compare *Refine*<sub>Contrast</sub> against human references, as well as *Refine*<sub>Contrast</sub> against *Translate* separately.

We plot the human preference results in Figure 1 (right). It reveals that GPT *Refine* attains slightly better performance in  $\text{zh} \rightarrow \text{en}$  and similar performance in  $\text{en} \rightarrow \text{zh}$  when compared with human references. On the other hand, it is more favourable

than GPT *Translate* in terms of human judgements. Combining the findings with translationese evaluation, we conclude that the refinement strategy could improve the naturalness of target translations without undermining the general quality.

## 5 Discussions

In Appendix D Table 5 we show outputs from different strategies for a single source input, where a native speaker marked preference for *Refine*<sub>Contrast</sub>, in both German→English and Chinese→English. We use different colours for phrase-level alignments to highlight the lexical variations. It illustrates that the word choice is diverse for both directions, and specifically for Chinese→English, there are substantial structural changes. The huge variety in expressions across translations can result in low BLEU against human references, but without much change in meaning as we observed, for instance, in Table 2 where BLEU can decline up to one-third, but neural metric scores change little.

Integrating LLMs into MT could benefit advances in both translationese reduction and translationese detection, yet we show the inability to measure translationese using automatic metrics at the moment. Finally, although the concepts of iterative refinement, post-editing, or translationese are not new, we use a combination of these to explore translationese reduction, instead of focusing on achieving state-of-the-art metric scores. Apart from the key related works in the introduction, we detail other works in Appendix F.

## 6 Conclusion and Future Work

We presented a simple way of including a powerful LLM in the process of translation refinement, which significantly reduces translationese in the outputs. It is shown that our method maintains translation quality and introduces lexical and structural changes, especially for high-resource into-English translation. Future work can explore sentence-level refinement decisions to reduce cost.

## 7 Limitations

Translationese is an interesting phenomenon in the field of translation studies, but it is difficult to quantify. Our work uses automatic scores to show changes in wording but not meaning. Then we rely on assessing the translations' naturalness as well as quality to show that translationese is reduced without hurting overall quality. We did not use any direct measure for translationese, but this is due to the lack of such at the moment.

We only experimented with GPT-3.5 without replicating with open-source LLMs. However, we argue that our intention is not to achieve state-of-the-art translation results, but to pose a new perspective on translationese reduction. Therefore, using a powerful LLM is necessary, and open-sourced models might not be as effective. Finally, involving GPT in an iterated process is costly. We think that GPT is useful in showcasing our proposed approach, but smarter refinement strategies need to be investigated for practical use cases.

## 8 Ethical Statement

The contents we analyse are machine-generated. We are not able to manually examine all model outputs, but we are fairly confident that the generated texts do not include harmful or inappropriate elements that will make readers uncomfortable. Our human evaluators are university students recruited by the authors. They are paired with an hourly rate higher than their local legal minimum wage.

## References

- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. [Neural machine translation by jointly learning to align and translate](#). In *3rd International Conference on Learning Representations*.
- Mona Baker. 1996. *Corpus-based Translation Studies: The Challenges that Lie Ahead*, Benjamins Translation Library. John Benjamins Publishing Company.
- Eleftheria Briakou, Colin Cherry, and George Foster. 2023. [Searching for needles in a haystack: On the role of incidental bilingualism in PaLM's translation capability](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*.

- Rajen Chatterjee, Matteo Negri, Raphael Rubino, and Marco Turchi. 2018. [Findings of the WMT 2018 shared task on automatic post-editing](#). In *Proceedings of the Third Conference on Machine Translation*.
- Kehai Chen, Masao Utiyama, Eiichiro Sumita, Rui Wang, and Min Zhang. 2022. [Synchronous refinement for neural machine translation](#). In *Findings of the Association for Computational Linguistics: ACL 2022*.
- Pinzhen Chen, Jindřich Helcl, Ulrich Germann, Laurie Burchell, Nikolay Bogoychev, Antonio Valerio Miceli Barone, Jonas Waldendorf, Alexandra Birch, and Kenneth Heafield. 2021. [The University of Edinburgh's English-German and English-Hausa submissions to the WMT21 news translation task](#). In *Proceedings of the Sixth Conference on Machine Translation*.
- Shamil Chollampatt, Raymond Hendy Susanto, Liling Tan, and Ewa Szymanska. 2020. [Can automatic post-editing improve NMT?](#) In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. [PaLM: Scaling language modeling with pathways](#). *arXiv preprint*.
- Koel Dutta Chowdhury, Rricha Jalota, Cristina España-Bonet, and Josef Genabith. 2022. [Towards debiasing translation artifacts](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- Akhbardeh Farhad, Arkhangorodsky Arkady, Biesialska Magdalena, Bojar Ondřej, Chatterjee Rajen, Chaudhary Vishrav, Marta R Costa-jussa, España-Bonet Cristina, Fan Angela, Federmann Christian, et al. 2021. [Findings of the 2021 conference on machine translation \(WMT21\)](#). In *Proceedings of the Sixth Conference on Machine Translation*.
- Markus Freitag, Isaac Caswell, and Scott Roy. 2019. [APE at scale and its implications on MT evaluation biases](#). In *Proceedings of the Fourth Conference on Machine Translation*.
- Markus Freitag, George Foster, David Grangier, and Colin Cherry. 2020a. [Human-paraphrased references improve neural machine translation](#). In *Proceedings of the Fifth Conference on Machine Translation*.
- Markus Freitag, David Grangier, and Isaac Caswell. 2020b. [BLEU might be guilty but references are not innocent](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*.
- Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André F. T. Martins. 2022. [Results of WMT22 metrics shared task: Stop](#)

416	using BLEU – neural metrics are better and more robust. In <i>Proceedings of the Seventh Conference on Machine Translation</i> .	Tom Kocmi and Christian Federmann. 2023. Large language models are state-of-the-art evaluators of translation quality. <i>arXiv preprint</i> .	470
417			471
418			472
419	Martin Gellerstam. 1986. Translationese in Swedish novels translated from English. In L. Wollin and H. Lindquist, editors, <i>Translation studies in Scandinavia: Proceedings from the Scandinavian Symposium on Translation Theory II</i> . CWK Gleerup.	Jason Lee, Elman Mansimov, and Kyunghyun Cho. 2018. Deterministic non-autoregressive neural sequence modeling by iterative refinement. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> .	473
420			474
421			475
422			476
423			477
424	Yvette Graham, Barry Haddow, and Philipp Koehn. 2020. Statistical power and translationese in machine translation evaluation. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing</i> .	Gennadi Lembersky, Noam Ordan, and Shuly Wintner. 2012. Language Models for Machine Translation: Original vs. Translated Texts. <i>Computational Linguistics</i> .	478
425			479
426			480
427			481
428			
429	Jiatao Gu, Changhan Wang, and Junbo Zhao. 2019. Levenshtein transformer. In <i>Advances in Neural Information Processing Systems</i> .	Qingyu Lu, Baopu Qiu, Liang Ding, Liping Xie, and Dacheng Tao. 2023. Error analysis prompting enables human-like translation evaluation in large language models: A case study on ChatGPT. <i>arXiv preprint</i> .	482
430			483
431			484
432	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</i> .	Jan Niehues, Eunah Cho, Thanh-Le Ha, and Alex Waibel. 2016. Pre-translation for neural machine translation. In <i>Proceedings of the 26th International Conference on Computational Linguistics</i> .	485
433			486
434			487
435			488
436			489
437			490
438	Julia Ive, Lucia Specia, Sara Szoc, Tom Vanallemeersch, Joachim Van den Bogaert, Eduardo Farah, Christine Maroti, Artur Ventura, and Maxim Khalilov. 2020. A post-editing dataset in the legal domain: Do we underestimate neural machine translation quality? In <i>Proceedings of the Twelfth Language Resources and Evaluation Conference</i> .	Roman Novak, Michael Auli, and David Grangier. 2016. Iterative refinement for machine translation. <i>arXiv preprint</i> .	491
439			492
440			493
441			
442			494
443			495
444			496
445	Wenxiang Jiao, Wenxuan Wang, Jen tse Huang, Xing Wang, and Zhaopeng Tu. 2023. Is ChatGPT a good translator? Yes with GPT-4 as the engine. <i>arXiv preprint</i> .	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. In <i>Advances in Neural Information Processing Systems</i> .	497
446			498
447			499
448			
449	Marcin Junczys-Dowmunt and Roman Grundkiewicz. 2018. MS-UEdin submission to the WMT2018 APE shared task: Dual-source transformer for automatic post-editing. In <i>Proceedings of the Third Conference on Machine Translation</i> .	Santanu Pal, Hongfei Xu, Nico Herbig, Sudip Kumar Naskar, Antonio Krüger, and Josef van Genabith. 2020. The transference architecture for automatic post-editing. In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> .	500
450			501
451			502
452			503
453			504
454	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. <i>arXiv preprint</i> .	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics</i> .	505
455			506
456			507
457			508
458			509
459	Kevin Knight and Ishwar Chander. 1994. Automated postediting of documents. In <i>Proceedings of the Twelfth AAAI National Conference on Artificial Intelligence</i> .	Maja Popović. 2017. chrF++: words helping character n-grams. In <i>Proceedings of the Second Conference on Machine Translation</i> .	510
460			511
461			512
462			
463	Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, et al. 2022. Findings of the 2022 conference on machine translation (WMT22). In <i>Proceedings of the Seventh Conference on Machine Translation</i> .	Matt Post. 2018. A call for clarity in reporting BLEU scores. In <i>Proceedings of the Third Conference on Machine Translation</i> .	513
464			514
465			515
466			
467			516
468			517
469			518
			519
			520
			521
			522
			523

524	Vikas Raunak, Amr Sharaf, Hany Hassan Awadallah, and Arul Menezes. 2023b. <a href="#">Leveraging GPT-4 for automatic translation post-editing</a> . <i>arXiv preprint</i> .		
525		<a href="#">tion shared task</a> . In <i>Proceedings of the Sixth Conference on Machine Translation</i> .	579
526			580
527	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. <a href="#">COMET: A neural framework for MT evaluation</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing</i> .		581
528			582
529			583
530			584
531			585
532	Parker Riley, Isaac Caswell, Markus Freitag, and David Grangier. 2020. <a href="#">Translationese as a language in “multilingual” NMT</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> .		586
533			587
534			588
535			589
536			590
537	Michel Simard, Cyril Goutte, and Pierre Isabelle. 2007. <a href="#">Statistical phrase-based post-editing</a> . In <i>Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics</i> .		591
538			592
539			593
540			594
541			595
542	Elke Teich. 2003. <i>Cross-Linguistic Variation in System and Text: A Methodology for the Investigation of Translations and Comparable Texts</i> . De Gruyter Mouton.		596
543			597
544			598
545			599
546	Antonio Toral. 2019. <a href="#">Post-editeese: An exacerbated translationese</a> . In <i>Proceedings of Machine Translation Summit XVII</i> .		600
547			
548			
549	Antonio Toral, Sheila Castilho, Ke Hu, and Andy Way. 2018. <a href="#">Attaining the unattainable? reassessing claims of human parity in neural machine translation</a> . In <i>Proceedings of the Third Conference on Machine Translation</i> .		601
550			602
551			603
552			604
553			605
554	Chau Tran, Shruti Bhosale, James Cross, Philipp Koehn, Sergey Edunov, and Angela Fan. 2021. <a href="#">Facebook AI’s WMT21 news translation task submission</a> . In <i>Proceedings of the Sixth Conference on Machine Translation</i> .		606
555			607
556			608
557			609
558			610
559	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. <a href="#">Attention is all you need</a> . In <i>Advances in Neural Information Processing Systems</i> .		611
560			
561			
562			
563			
564	David Vilar, Markus Freitag, Colin Cherry, Jiaming Luo, Viresh Ratnakar, and George Foster. 2023. <a href="#">Prompting PaLM for translation: Assessing strategies and performance</a> . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics</i> .		612
565			613
566			614
567			615
568			616
569			617
570	Longyue Wang, Mu Li, Fangxu Liu, Shuming Shi, Zhaopeng Tu, Xing Wang, Shuangzhi Wu, Jiali Zeng, and Wen Zhang. 2021. <a href="#">Tencent translation system for the WMT21 news translation task</a> . In <i>Proceedings of the Sixth Conference on Machine Translation</i> .		618
571			619
572			620
573			621
574			622
575	Daimeng Wei, Zongyao Li, Zhanglin Wu, Zhengzhe Yu, Xiaoyu Chen, Hengchao Shang, Jiixin Guo, Minghan Wang, Lizhi Lei, Min Zhang, et al. 2021. <a href="#">HW-TSC’s participation in the WMT 2021 news translation shared task</a> . In <i>Proceedings of the Sixth Conference on Machine Translation</i> .		623
576			624
577			
578			

## A Additional scores for GPT refinement

Due to the space constraint, we are not able to display all metric scores in the main content, so we attach chrF++ and COMET<sub>DA</sub> scores here for reference. We observe the same patterns in BLEU and chrF++ across all language pairs. Regarding COMET<sub>DA</sub>, as we have discussed, it is conditioned on the human reference, which (1) can be translationese-prone itself, and (2) is a subject in our comparison. Hence it might be not indicative. The Additional scores for GPT refinement experiments are listed in Table 3.

## B WMT system refinement

Out of the seven WMT21 submissions, we select outputs from four models built by research labs that, based on human evaluation, have been ranked at significantly different positions on the German-to-English leaderboard: Tencent (Wang et al., 2021), Facebook AI (Tran et al., 2021), Edinburgh (Chen et al., 2021), and Huawei TSC (Wei et al., 2021). These are competitive systems built with data augmentation, multilingualism, ensembling, re-ranking, etc. We then include two online commercial systems tested in WMT 2021: Online-A and Online-Y.<sup>3</sup> Finally, human reference “B”

<sup>3</sup>The online systems were anonymized by WMT21 organizers, so we do not have knowledge about them. The time of access is believed to be in 2021.

	WMT21 de→en		WMT21 en→de		WMT21 zh→en		WMT21 en→zh		WMT22 de→fr		WMT22 en→ja		WMT22 uk→cs	
	chrF++	COMET <sub>DA</sub>	chrF++	COMET <sub>DA</sub>	chrF++	COMET <sub>DA</sub>	chrF++	COMET <sub>DA</sub>	chrF++	COMET <sub>DA</sub>	chrF++	COMET <sub>DA</sub>	chrF++	COMET <sub>DA</sub>
Reference <sub>A</sub>	-	-	-	-	-	-	-	-	-	-	-	-	-	-
Translate	<b>57.55</b>	<b>.8606</b>	<b>53.54</b>	.8427	<b>53.74</b>	.8199	<b>20.61</b>	.8300	<b>59.50</b>	<b>.8395</b>	25.89	.8863	<b>54.64</b>	<b>.9074</b>
Refine	51.91	.8525	50.57	<b>.8478</b>	49.06	.8156	19.28	<b>.8417</b>	55.83	.8353	<b>27.30</b>	<b>.8941</b>	53.06	.9040
Refine <sub>Contrast</sub>	52.47	.8452	51.21	.8211	51.77	<b>.8538</b>	19.69	.8395	56.37	.8308	26.71	.8928	54.29	.9036
Refine <sub>Random</sub>	51.79	.7777	46.56	.7906	47.11	.8323	17.49	.8126	-	-	-	-	-	-
Paraphrase	40.05	.8044	43.54	.8197	40.92	.7931	17.14	.8144	44.28	.7937	23.18	.8592	40.04	.8625

Table 3: Additional automatic scores of different strategies on translation directions from WMT 2021 and 2022 news translation.

	BLEU	chrF++	COMET <sub>DA</sub>	COMET <sub>QE</sub>
Reference <sub>A</sub>	-	-	-	.0919
Submission	<b>30.05</b>	<b>56.00</b>	.8497	.1050
Refine	23.39	51.80	.8527	<b>.1123</b>
Refine <sub>Contrast</sub>	25.10	53.82	<b>.8566</b>	.1116
Paraphrase	12.52	41.03	.8031	.0894
Submission	<b>34.45</b>	<b>60.78</b>	<b>.8582</b>	.1061
Refine	23.37	51.67	.8494	.1098
Refine <sub>Contrast</sub>	25.14	52.84	.8534	<b>.1137</b>
Paraphrase	12.22	41.34	.8097	.0942
Submission	<b>32.70</b>	<b>59.32</b>	.8500	.0981
Refine	22.92	50.85	<b>.8522</b>	.1080
Refine <sub>Contrast</sub>	24.40	53.32	.8517	<b>.1134</b>
Paraphrase	11.97	40.29	.8054	.0892
Submission	<b>35.35</b>	<b>61.28</b>	<b>.8584</b>	.1055
Refine	23.75	52.16	.8488	.1095
Refine <sub>Contrast</sub>	26.89	54.75	.8553	<b>.1116</b>
Paraphrase	12.43	41.35	.8116	.0947
Submission	<b>34.67</b>	<b>60.78</b>	<b>.8677</b>	<b>.1146</b>
Refine	22.97	51.05	.8505	.1113
Refine <sub>Contrast</sub>	25.74	53.88	.8548	.1130
Paraphrase	11.80	40.99	.8099	.0922
Submission	<b>34.20</b>	<b>60.03</b>	<b>.8588</b>	.1087
Refine	22.04	50.29	.8496	.1097
Refine <sub>Contrast</sub>	25.24	52.87	.8546	<b>.1147</b>
Paraphrase	12.79	40.18	.8067	.0921
Submission	<b>35.13</b>	<b>61.17</b>	<b>.8643</b>	<b>.1126</b>
Refine	22.24	50.82	.8519	.1097
Refine <sub>Contrast</sub>	24.95	52.47	.8560	<b>.1124</b>
Paraphrase	12.20	40.74	.8078	.0909

Table 4: Automatic scores of refining WMT 2021 news shared task German-to-English submissions.

is added so that we can experiment with our refinement strategy with human translations.<sup>4</sup> References “A” and “B” are sourced from different translation agencies (Farhad et al., 2021).

We report automatic scores from the refinement process in Table 4. We explain the results in the main content Section 3.3. Overall, we ob-

<sup>4</sup>The overview paper of WMT 2021 states that “for German↔English, the ‘B’ reference was found to be a post-edited version of one of the participating online systems”. We discover that it refers to English→German only, and German→English is not affected.

serve patterns similar to refining GPT translations. The string-based metrics see significant drops, but COMET<sub>QE</sub> improves for five out of seven original entries.

### C Score changes through iterations

We plot the changes in BLEU, COMET<sub>DA</sub>, and COMET<sub>QE</sub> in Figure 2. Apart from scores from our translate and refinement queries, we also include the human reference performance in the COMET<sub>QE</sub> plot.

### D Example outputs

We place two examples in Table 5 as a case study. The cases illustrate significant string changes, but the meaning of sentences does not vary too much. This signifies the inability to use automatic string-based metrics in distinguishing translation quality or the degree of translationese when the outputs are relatively high-quality.

### E Evaluation metric details

BLEU and chrF++ are as implemented in the sacrebleu toolkit.<sup>5</sup> We also use this toolkit to obtain test sets with references as well as past WMT systems’ outputs. Specifically for tokenization in BLEU calculation, we use “zh” for Chinese, “jamecab” for Japanese, and “13a” for the rest. The BLEU signature is nrefs:1 | case:mixed | eff:no | smooth:exp | version:2.3.1, and the chrF++ signature is nrefs:1 | case:mixed | eff:yes | nc:6 | nw:2 | space:no | version:2.3.1. For COMET metrics, we used the official implementation released by the authors.<sup>6</sup>

<sup>5</sup><https://github.com/mjpost/sacrebleu>

<sup>6</sup><https://github.com/Unbabel/COMET>



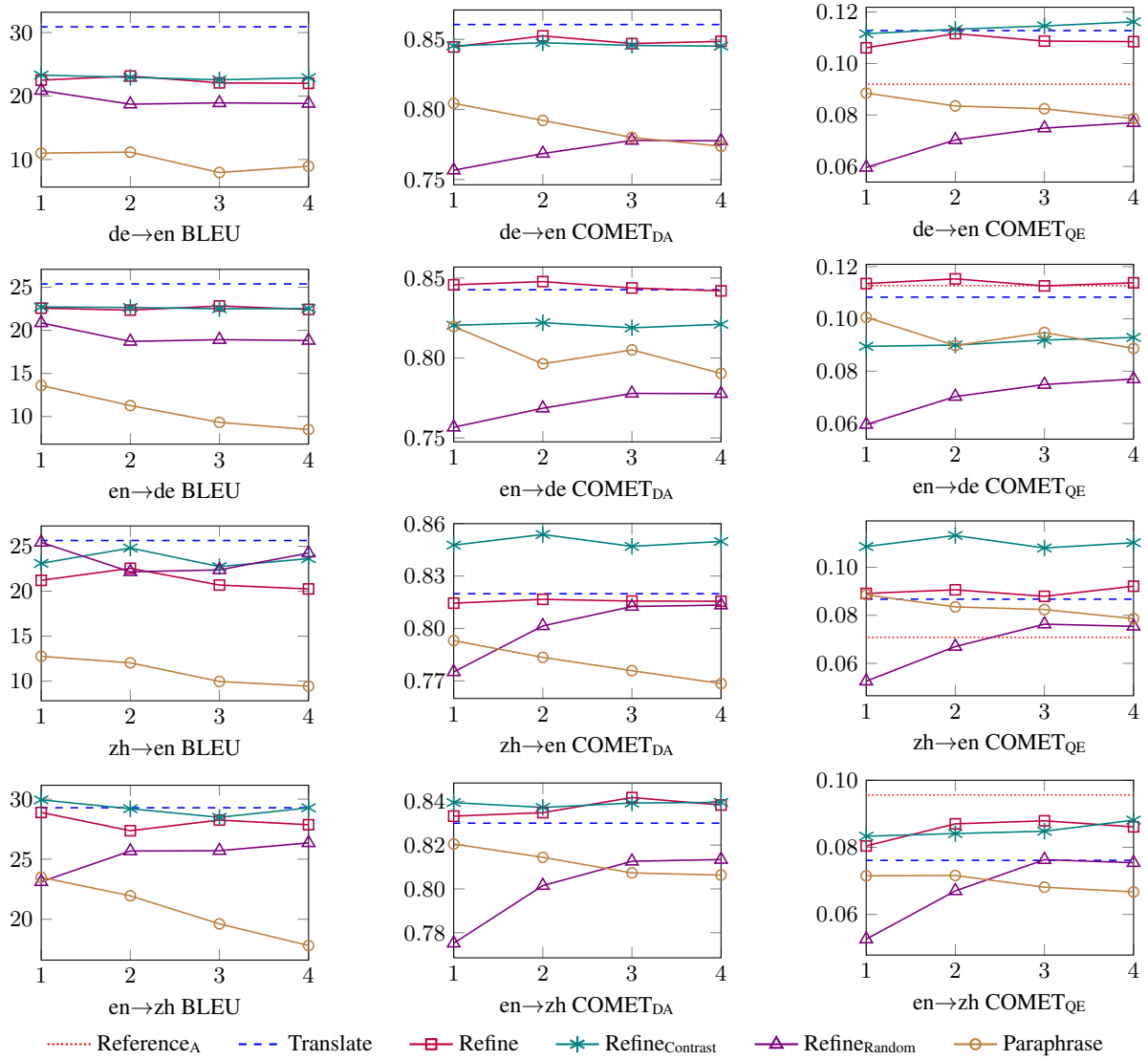


Figure 2: BLEU, COMET<sub>DA</sub>, and COMET<sub>QE</sub> at different refinement and paraphrase iterations for high-resource translation.

Source	Der 17-Jährige floh zunächst vom Tatort, seine Personalien konnten aber im Nachhinein ermittelt werden.
Reference	The 17 year-old proceeded to flee the crime scene, however, his personal details could be retrieved later.
Translate	The 17-year-old initially fled from the crime scene, but his personal information was later determined.
Refine <sub>Contrast</sub>	The 17-year-old initially fled from the scene of the crime, but his personal details could later be identified.
Paraphrase	At first, the 17-year-old ran away from where the crime occurred, but eventually, the authorities were able to identify him by his personal details.
Source	新法令规定，坎帕尼亚大区自即日起室内公共场所必须戴口罩，违者最高可处以1000欧元罚金。
Reference	According to a new decree, people must wear masks in indoor public places in Campania from now on, and offenders can be fined up to 1,000 euros.
Translate	A new regulation stipulates that in Campania, indoor public places must wear masks. Violators can be fined up to 1000 euros.
Refine <sub>Contrast</sub>	A new regulation states that in the Campania region, masks must be worn in indoor public places, with a maximum fine of 1000 euros for those who violate the rule.
Paraphrase	A new rule in Campania requires people to wear masks in indoor public places, and those who don't follow this rule may be charged up to 1000 euros.

Table 5: German→English and Chinese→English examples showing rich lexical variations across translation strategies.

## F Other related works

### F.1 Translation post-editing

Closely related to translation refinement is automatic post-editing (APE), which trains a neural network to fix translation errors by learning from human correction data (Knight and Chander, 1994). While it has shown notable developments in statistical machine translation, it could become less effective in the deep learning era due to original translations being high-quality and lack of post-editing data (Junczys-Dowmunt and Grundkiewicz, 2018; Chatterjee et al., 2018). Whilst one way to facilitate this is more data provision (Chollamatt et al., 2020; Ive et al., 2020), our workaround utilizes a large language model, which possesses the post-editing capability without being specifically tuned. Furthermore, post-editing models have limited power to alleviate translationese, because human editing data is collected from annotators who are usually instructed to not make style improvements (Ive et al., 2020). Compared to APE, our method allows LLMs to re-generate an entirely different translation, which could escape the “post-editese” phenomenon, where Toral (2019) demonstrated that human-edited machine translations still exhibit translationese features.

Some post-editing works do not rely on the source translation or human editing data (Simard et al., 2007). For instance, Freitag et al. (2019) trained a post-editor solely on monolingual data by reconstructing the original text given its round-trip translation. In our work, we incorporate stronger natural language modelling into post-editing by employing LLMs. Other translation refinement research includes combining statistical and neural systems (Novak et al., 2016; Niehues et al., 2016), merging APE into the NMT framework (Pal et al., 2020; Chen et al., 2022), and debiasing translationese in the latent embedding space (Dutta Chowdhury et al., 2022). The iterative editing mechanism is not commonly employed in autoregressive translation or translation editing. Its use cases mostly lie in non-autoregressive translation, where each output token is independent of other target positions and iterative decoding enhances output quality (Lee et al., 2018; Gu et al., 2019; Xu and Carpuat, 2021).

### F.2 Large language models

Large language models have recently become highly effective tools for various NLP tasks (Rad-

ford et al., 2019; Brown et al., 2020; Chowdhery et al., 2022; Ouyang et al., 2022). Nowadays, optimising LLMs directly for specific tasks becomes infeasible yet unnecessary since they generalize to downstream tasks without explicit supervision. With more parameters and training data, LLMs may offer stronger performance than dedicated translation or post-editing models. The method we use to elicit a response from GPT is zero-shot hard prompting (Brown et al., 2020), which means affixing a description to the original task input to form a query to the model. Researchers have benchmarked LLMs’ capability to translate (Vilar et al., 2023; Zhang et al., 2023; Jiao et al., 2023; HENDY et al., 2023), and to evaluate translations (Kocmi and Federmann, 2023; Lu et al., 2023; Xu et al., 2023).

Recent findings show that GPT produces less literal translations, especially for out-of-English translations (Raunak et al., 2023a), which to some extent stands in contrast with our evaluation outcome. Concurrent with our study, Raunak et al. (2023b) formalized post-editing as a chain-of-thought process (Wei et al., 2022) with GPT-4 and showed promising results. Different from their focus, our work features the iterative refinement process as a means to mitigate translationese. The improvement, especially for into-English, may be attributed to the abundant English pre-training data available for LLMs. To the best of our knowledge, although the concept of iterative refinement is not new, ours is the pioneering paper in applying such strategies to LLMs for translation.