

Recovered in Translation: Efficient Pipeline for Automated Translation of Benchmarks and Datasets

Anonymous ACL submission

Abstract

The reliability of multilingual Large Language Model (LLM) evaluation is currently compromised by the inconsistent quality of translated benchmarks. Existing resources often suffer from semantic drift and context loss, which can lead to misleading performance metrics. In this work, we present a fully automated framework designed to address these challenges by enabling scalable, high-quality translation of datasets and benchmarks. We demonstrate that adapting test-time compute scaling strategies, specifically Universal Self-Improvement (USI) and our proposed multi-round ranking method, T-RANK, allows for significantly higher quality outputs compared to traditional pipelines. By effectively applying these methods, our framework ensures that benchmarks preserve their original task structure and linguistic nuances during localization. We apply this approach to translate popular benchmarks and datasets into eight Eastern and Southern European languages. Evaluations using both reference-based metrics and LLM-as-a-judge show that our translations surpass existing resources, resulting in more accurate downstream model assessment. We release both the framework and the improved benchmarks to facilitate robust and reproducible multilingual AI development.

1 Introduction

Large Language Models have demonstrated rapid progress in machine translation, now outperforming classical tools such as Google Translate and DeepL (Deutsch et al., 2025). This advancement benefits multilingual model development by addressing data and benchmark scarcity in non-English settings. Recent work shows that sampling multiple translation candidates at higher temperatures enhances translation quality, with test-time scaling methods like Best-of-N (Stiennon et al., 2020) and Fusion-of-N (Khairi et al., 2025a) demonstrating improved performance across multilingual domains.

A significant gap in recent research is the limited exploration of benchmark translation quality. Multilingual benchmarks are critical for assessing model development efficiency, yet receive less attention than data development. Existing benchmarks remain imperfect due to reliance on older LLMs or oversimplified methods – most were translated using classical MT tools lacking instruction prompting capabilities, or older models like GPT-4 whose multilingual capabilities lag behind current frontier models (Lai et al., 2023).

We investigate available benchmarks and identify common translation flaws stemming from translating questions and answers separately, creating contextual and grammatical mismatches. We focus on Eastern and Southern European languages because: (1) they exhibit complex grammatical features (extensive case systems, grammatical gender, aspect-based verbs) that are sensitive to contextual misalignment; (2) they represent mid-resource languages where frontier models still lack adequate localization; (3) unlike lower-resource languages, translated benchmarks exist for comparison. We translate MMLU, Hellaswag, ARC, and Winogrande into Ukrainian, Romanian, Slovak, Lithuanian, Bulgarian, Turkish, Greek and Estonian.

In this paper, we address three key research questions:

- **RQ1:** Are current translated multilingual benchmarks robust and reliable?
- **RQ2:** To what extent do test-time compute scaling methods improve machine translation quality?
- **RQ3:** How can we translate benchmarks and datasets into other languages while efficiently integrating their unique linguistic features into the evaluated tasks?

To address these challenges, we present a novel automated translation framework supporting four

083	methods across various model types, including	benchmark acquisition, thereby facilitating more	132
084	open-weight models. The framework facilitates ma-	time- and cost-effective training of locally adapted	133
085	chine translation of datasets and benchmarks with	LLMs.	134
086	minimal manual supervision and maximum con-		
087	figurability, offering multiple methods to balance	2 Background and Related Work	135
088	costs and time investment based on language per-		
089	formance in selected MT models. Our framework	2.1 Foundation for LLM-Based Efficient	136
090	incorporates: (1) classic one-prompt translation	Translation	137
091	with optional self-correction (SC) as a lightweight	The rapid expansion of multilingual applications	138
092	solution; (2) Best-of-N with LLM-as-a-judge scor-	for large language models (LLMs) necessitates scal-	139
093	ing to provide reward signals; (3) Universal Self-	able and high-quality translation frameworks. Re-	140
094	Improvement (USI) with multi-prompt sampling	cent advancements propose innovative methodolo-	141
095	and unification of translation candidates into re-	gies to enhance LLM-based translation through	142
096	refined outputs; (4) our newly proposed Translation	adaptive prompting and self-refinement strategies.	143
097	Ranking (T-RANK), which employs multi-prompt	The Adaptive Few-shot Prompting (AFSP)	144
098	candidate sampling and multi-round competitive	framework addresses prompt sensitivity in machine	145
099	ranking to enhance error detection and achieve su-	translation by dynamically selecting suitable trans-	146
100	perior translation quality.	lation demonstrations. AFSP retrieves semantically	147
101	We evaluate these test-time scaling methods	similar examples from parallel corpora based on	148
102	across several mid-resource languages, translate	LLM-generated embeddings, then generates and	149
103	popular benchmarks into these languages, and com-	reranks multiple translation candidates to select	150
104	pare our results with existing translations. Our	the most contextually appropriate translation (Tang	151
105	findings demonstrate that the integrated methods	et al., 2025).	152
106	produce higher-quality translations with more ac-	The TEaR (Translate, Estimate, and Refine)	153
107	curate evaluation results and improved text qual-	framework introduces systematic self-refinement	154
108	ity. Moreover, these methods yield performance	for LLM-based translation. TEaR operates by trans-	155
109	improvements on established machine translation	lating the source text, estimating translation qual-	156
110	benchmarks such as WMT24++ and FLORES.	ity, and refining based on identified errors. This itera-	157
111	We summarize our key contributions as follows:	tive process enables LLMs to self-correct and im-	158
112		prove translation quality, with cross-model experi-	159
113	• We provide a comprehensive analysis of mul-	ments demonstrating that general-purpose LLMs	160
114	tilingual benchmark quality and identify the	can perform both translation and evaluation simul-	161
115	root causes of current translation deficiencies.	taneously (Feng et al., 2024).	162
116		In addition to these frameworks, several test-	163
117	• We propose a novel, fully automated frame-	time compute scaling methods show potential for	164
118	work that enables efficient translation of	enhancing LLM translation capabilities. Although	165
119	datasets and benchmarks with maximum flex-	originally designed for general reasoning tasks	166
120	ibility. The pipeline is fully configurable, al-	such as mathematics and coding, these methods	167
121	lowing practitioners to optimize the balance	demonstrate promising applications in translation:	168
122	between translation quality and cost in multi-		
123	lingual development.	• Best-of-N Sampling: This method generates	169
124		multiple translation outputs and selects the	170
125	• We release several benchmarks translated into	best one based on predefined criteria. It lever-	171
126	Eastern and Southern European languages, in-	ages the diversity in LLM outputs due to tem-	172
127	cluding Ukrainian, Romanian, Slovak, Lithua-	perature variations, increasing the likelihood	173
128	nian, Bulgarian, Turkish, Greek and Estoa-	of obtaining high-quality translations (Stien-	174
129	nian. We find out, that translation quality in-	non et al., 2020).	175
130	fluences benchmark evaluation results signifi-	• Universal Self-Consistency (USC): USC ex-	176
131	cantly, and our proposed methods yield more	tends the concept of self-consistency by en-	177
	reliable and accurate multilingual evaluation.	abling LLMs to select the most consistent an-	178
		swer among multiple candidates without rely-	179
	These findings enable more efficient multilingual	ing on answer extraction processes. This ap-	180
	model development by improving both data and		

181 proach is particularly effective for open-ended
182 generation tasks, improving performance in
183 mathematical reasoning, code generation, and
184 summarization (Chen et al., 2023).

- 185 • **Fusion-of-N:** Rather than selecting the single
186 best answer as in Best-of-N, Fusion-of-N
187 synthesizes the most informative elements
188 from multiple candidates into a single final
189 output. The method uses an LLM judge to aggregate
190 diverse strengths from the candidate pool. Fusion-of-N
191 outperforms Best-of-N and shows promising gains on
192 multilingual tasks, including machine translation (Khairi et al.,
193 2025a).
194

195 Findings from Khairi et al. (2025b) also confirm,
196 that sampling multiple candidates with higher
197 temperature with self-improvement and refined selection
198 process leads to significant improvements in many
199 multilingual domains, including machine translation.
200

201 Results from the WMT24++ benchmark (Deutsch et al., 2025)
202 further corroborate the efficacy of these methodologies,
203 demonstrating that state-of-the-art LLMs outperform
204 traditional machine translation tools across various
205 language pairs. The research suggests that integrating
206 test-time compute strategies can significantly enhance
207 translation quality, particularly for low-resource
208 languages. A clear trend emerges whereby newer
209 and larger LLMs consistently outperform existing
210 tools and earlier model iterations, indicating the
211 potential for continued improvements in translation
212 quality as more advanced models are released.
213

Takeaway 1: Large language models outperform traditional machine translation systems like Google Translate and DeepL through targeted prompts that address domain-specific and language-specific translation challenges.

214
215 Collectively, these methodologies underscore the
216 importance of adaptive prompting, self-refinement,
217 and test-time scaling strategies in enhancing LLM-
218 based translation systems. By integrating these
219 approaches, it is possible to develop more robust,
220 efficient, and high-quality translation frameworks
221 capable of addressing the challenges posed by low-
222 resource languages while minimizing the need for
223 manual human intervention.

2.2 Shortcomings of Existing Multilingual Resources

224 The rapid development of large language models (LLMs)
225 has highlighted the need for robust multilingual
226 benchmarks. However, translating existing
227 benchmarks often introduces issues that compromise
228 accurate assessment of LLM capabilities across
229 languages.
230

231 A prominent example is the MuBench benchmark
232 dataset introduced by Han et al. (2025), comprising
233 widely used benchmarks (Hellawag (Zellers et al., 2019),
234 ARC (Clark et al., 2018), Winogrande (Sakaguchi et al., 2019),
235 MMLU (Hendrycks et al., 2021), and others) translated into
236 61 languages with 3.9M samples. While this initiative
237 represents an important step toward multilingual
238 evaluation, the translation process relies on an
239 automated pipeline with quality control measures
240 including semantic consistency evaluation, translation
241 purity assessment, and cultural sensitivity checks.
242 However, despite human expert evaluation of 30k
243 samples across 15 languages combined with LLM-as-a-
244 Judge validation, the predominantly automated
245 approach raises concerns regarding translation
246 quality for the remaining samples, as machine
247 translation methods may fail to capture nuanced
248 linguistic and cultural contexts in all 61 languages.
249 Furthermore, variations in quality assurance depth
250 across languages could introduce semantic inaccuracies
251 that skew evaluation results.
252

253 Beyond general quality concerns, certain benchmarks
254 present inherent translation challenges. For example,
255 Winogrande, MMLU and Hellawag may contain answer
256 options with gender-specific adjective endings for
257 sentence completion tasks. In languages with
258 gender-specific adjectives, translating these options
259 can inadvertently reveal correct answers or mislead
260 towards picking an incorrect option, allowing models
261 to succeed through language proficiency rather than
262 reasoning capability. Such linguistic features can
263 compromise evaluation integrity through unintended
264 answer leakage.
265

266 The Global-MMLU project (Singh et al., 2024) represents
267 a substantial effort to advance multilingual evaluation
268 by translating the MMLU benchmark into 42 languages.
269 The process combines machine translation with
270 crowdsourced human verification; however, only
271 approximately 20% of machine-translated texts
272 underwent manual correction. Moreover, questions
273 and answers were translated separately, resulting in
274 observable gram-

275 matical inconsistencies in translated entries for lan- 320
276 guages such as Ukrainian. This context-agnostic 321
277 translation approach can produce inconsistencies 322
278 or grammatical errors that affect benchmark coher- 323
279 ence and accuracy, particularly for lower-resource 324
280 languages. Additionally, relying on tools like 325
281 Google Translate without thorough human valida- 326
282 tion may overlook language-specific nuances es- 327
283 sential for accurate evaluation. 328

284 The Okapi framework (Lai et al., 2023) intro- 329
285 duces a novel approach to multilingual instruc- 330
286 tion tuning by leveraging Reinforcement Learn- 331
287 ing from Human Feedback (RLHF). Unlike tradi- 332
288 tional methods relying solely on supervised fine- 333
289 tuning (SFT), Okapi combines SFT with RLHF 334
290 to align model outputs more closely with human 335
291 preferences across diverse languages. Moreover, 336
292 some of the popular benchmarks like MMLU, Hel- 337
293 laswag and ARC were also translated to 26 lan- 338
294 guages to enable multilingual evaluation. How- 339
295 ever, Okapi’s translation process utilized the GPT-4 340
296 model series and the translation could be further 341
297 improved by employing test-time compute scal- 342
298 ing methods to enhance translation quality. Addi- 343
299 tionally, the framework does not explicitly address 344
300 language-specific grammatical features that may 345
301 impact benchmark coherence and accuracy. 346

Takeaway 2: Translating questions and answer options within the same prompt context is essential for sentence completion tasks, as it preserves semantic relationships and prevents contextual misleading during evaluation.

302
303 Recent studies, including WMT24++ (Deutsch 347
304 et al., 2025), demonstrate that state-of-the-art 348
305 LLMs outperform traditional machine translation 349
306 tools such as DeepL or Google Translate across 350
307 all evaluated language pairs. This finding suggests 351
308 that reliance on older translation tools may be in- 352
309 sufficient for ensuring high-quality translations in 353
310 multilingual datasets and benchmarks. Further- 354
311 more, practical constraints characteristic of pop- 355
312 ular frameworks (such as low API rate limits or 356
313 inability to use advanced prompting techniques) 357
314 can hinder large-scale translation efforts, making it 358
315 more difficult to produce comprehensive and accu- 359
316 rate datasets. 360

317 In summary, while initiatives such as MuBench, 361
318 Global-MMLU, and Okapi represent important 362
319 progress toward multilingual LLM evaluation, the 363

275 methods employed in translating benchmarks ex- 320
276 hibit limitations that hinder improvement. De- 321
277 pendence on automatic translation without com- 322
278 prehensive human verification, challenges arising 323
279 from language-specific grammatical structures, and 324
280 the limited capabilities of translation tools collec- 325
281 tively highlight the need for more sophisticated 326
282 approaches that respect linguistic features. Further- 327
283 more, current frontier language models are already 328
284 outperforming classical tools, which were used to 329
285 translate some of widely used multilingual bench- 330
286 marks. Future research should advance automated 331
287 methodologies that reduce reliance on manual cu- 332
288 ration while simultaneously improving both the 333
289 quality and accessibility of multilingual evaluation 334
290 resources. 335

3 Building an Efficient Automated Translation Framework 336

In this section we present a novel automated trans- 338
512 lation framework, which utilizes Large Language 339
512 Models with features adapted for both dataset 340
512 and benchmark (QA/test) formats. Moreover, the 341
512 framework also allows flexibility in methods and 342
512 their inner configurable parameters to optimize 343
512 cost- and time-effectiveness. 344

3.1 Motivation 345

As highlighted in previous sections, we observed 346
512 the lack of research and solutions towards scaled 347
512 automated translation adapted for custom benchmark 348
512 formats without the loss of the translation quality. 349
512 For this framework, we aim to tackle the following 350
512 key problems in LLM-assisted data translation: 351

- **Support for Diverse Data Formats:** Many 352
536 datasets have complex, nested structures that 353
536 complicate processing. For LLM training, flat 354
536 string fields without hierarchical nesting en- 355
536 sure easier ingestion and more efficient work- 356
536 flows. 357
- **Preservation of Benchmark QA Structure:** 358
536 Benchmark evaluation reliability depends on 359
536 maintaining coherence between questions and 360
536 answer choices. Preserving these relation- 361
536 ships during translation ensures benchmarks 362
536 remain faithful to the original task structure. 363
- **Adapting to Varying Language Resource Availability:** High-resource languages (Ger- 364
536 man, Spanish, Hindi) achieve strong trans- 365
536 lation quality with simple zero-shot prompt- 366
536 367

ing, while mid- and low-resource languages require sophisticated, language-specific pipelines. Users need flexibility to configure translation approaches based on each language’s characteristics.

- **Addressing Language-Specific Phenomena:** EEU languages possess unique grammatical features requiring careful handling. Language-specific prompt engineering with few-shot examples and LLM-powered verification stages can improve translation fidelity by identifying and correcting language-specific errors.

3.2 Methodology

In our final framework, we propose two configuration modes – Dataset and Benchmark. They use different prompts and data format handling, since dataset is more straightforward, but benchmark has a complex connection between question and answers, which needs to be preserved. Finally, we propose four translation methods:

- **SC (Self-Check):** 0-shot simple translation with optional additional check from another LLM (different chat).
- **Best-of-N sampling:** sampling N translation candidates and prompting the model to score them, then picking the one with the highest score.
- **USI (Universal Self-Improvement):** sampling N translations, asking the model to combine them into the best one according to pre-defined evaluation criteria.
- **T-RANK (Translation Ranking):** sampling N translations, asking the model to rank them according to their quality and criteria eligibility, then correcting the best candidate.

We now examine the translation methods employed in our framework, highlighting their advantages and optimal use cases.

Default translation with optional self-check (SC). This method involves a simple 0-shot prompting of LLM to translate a text. The user has an option to include a self-check stage: after translation, the model (in a new chat with no history) is prompted to evaluate and correct the result with respect to the original text content. This method is suitable for large text translation into high-resource languages, as it is less costly to perform, due to the

sufficient translation capabilities in high-resource languages; however, since the model is prompted to look out for potential errors it might hallucinate them. Additionally, there exists an option to include a few-shot prompt to provide an example of which language-specific points the model should consider during translation.

Best-of-N sampling (BoN). Drawing from test-time compute scaling methods, we implement Best-of-N sampling without a reward model to maintain a training-free framework. We sample N translations at higher temperature (0.7) for diversity, then prompt the LLM to score candidates 1-10 based on specified criteria (Appendix A.6), selecting the highest-scored translation. While cost-effective and language-agnostic, this method yields lower quality than T-RANK and USI, as LLMs exhibit limitations in numerical scoring and positional bias, favoring earlier candidates despite identifying obvious errors effectively.

Universal Self-Improvement (USI). Building on Universal Self-Consistency and Fusion-of-N, this method operates on the principle that the most consistent translation is not necessarily the best. While Fusion-of-N improves translation quality substantially, the absence of precise translation-specific metrics limits the model’s ability to efficiently identify each candidate’s strengths for optimal merging, potentially resulting in fusion outputs that incorporate errors. Therefore, we adapt this approach to cultivate self-improvement specifically for machine translation. We sample N candidate translations using higher temperature (0.7 recommended), then present them to an evaluator LLM with instructions to combine the candidates into the best version according to specified criteria. This method proves cost-efficient and time-efficient, requiring only $N + 1$ model calls per entry, while successfully addressing and correcting language-specific features, making it particularly suitable for low-resource languages. We illustrate the Universal Self-Improvement workflow in more detail in Figure 1.

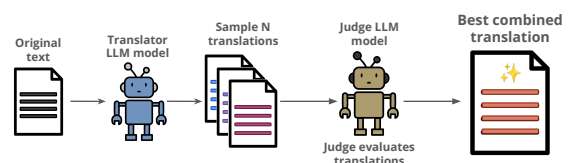


Figure 1: Universal Self-Improvement (USI) method workflow outlook

Translation Ranking (T-RANK). Building upon prior work in translation ranking, we propose an adaptive approach that refines this concept. This method begins by sampling N diverse translations with high temperature. A judge model then evaluates these candidates by ranking them according to predetermined quality criteria, with the top-ranked translation selected as the final output, optionally with refinements. Using this method, we observe more sophisticated reasoning from non-reasoning models, as they successfully identify subtle errors that other methods fail to correct. Moreover, evaluation through comparative ranking rather than numerical scoring appears to facilitate more detailed and rigorous assessment of translation quality.

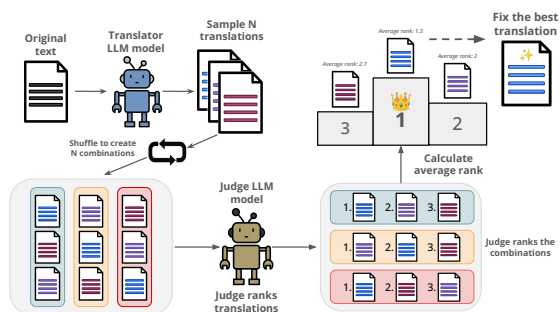


Figure 2: Translation Ranking (T-RANK) method workflow outlook

Previous research on LLM-based evaluation has identified several risks and biases inherent to such methodologies. One notable issue is positional bias: LLMs tend to assign higher scores to candidates presented earlier in the sequence. Furthermore, LLM judges often focus their evaluation on the first candidate, identifying specific flaws, and subsequently bias their assessment of following candidates by primarily searching for similar errors, potentially overlooking novel issues in later candidates.

To reduce these prejudices, we introduce a multiple-round sequential ranking strategy, in which the candidate translations are systematically moved across different positional orders. In particular, each translation candidate is presented once in every possible position across all rounds. This process can be visualized as the construction of a sequential grid, where with each new round, candidates are moved one position further till the end (as visualized in Figure 2). After N rounds for N candidates, each candidate has been presented in each position exactly once. Lastly, we show all

the candidates to the judge model again and ask to correct and refine the selected translation candidate if needed - this improves model’s ability to notice potential translation flaws by looking at strong and weak sides of other variants. We noticed an improved behavior and reasoning when using this method and in particular, showing all candidates when doing the final correction as shown in Appendix A.1. This approach shows a promising improvement in reliability of the evaluation by reducing positional biases while controlling evaluation costs and computational efficiency with a total of $2N + 1$ model calls. Optionally, one can combine best candidate correction during ranking, however, we found evaluator model’s reasoning to be more attentive in a separated setup.

We provide example prompts for all proposed methods in Appendix A.6.

Takeaway 3: While Best-of- N and USI sampling improve translation quality, the T-RANK method’s competitive ranking approach proves more efficient at identifying subtle translation errors.

3.3 Benchmark Translation for Eastern European Languages

We have translated several popular benchmarks into Eastern European languages using our proposed framework. The selected benchmarks include MMLU, Hellaswag, ARC, and Winogrande, which are widely used for evaluating LLM capabilities across various tasks. We focused on Eastern and Southern European languages such as Ukrainian, Romanian, Slovak, Lithuanian, Greek, Bulgarian, Turkish and Estonian due to their complex grammatical structures and mid-resource status. These languages present unique challenges for machine translation, making them ideal candidates for testing the effectiveness of our framework. MMLU was translated using GPT-4o-mini-2024-07-18 model checkpoint from OpenAI and all mentioned benchmarks were translated to Ukrainian using the same model. Remaining benchmark translations to other languages were performed using Gemini-2.0-Flash model from Google. In the next section and in Appendix A.3, we present a comprehensive evaluation of the translation quality achieved through our framework, comparing it with existing translations and assessing its impact on benchmark performance.

4 Evaluation Results

4.1 Machine Translation Benchmarks

We utilize two datasets to evaluate the translation quality of our proposed methods: FLORES and WMT24++.

FLORES. The FLORES benchmark evaluates machine translation systems across multilingual scenarios. FLORES-101 provides 3,001 (1,012 devtest split) sentences from English Wikipedia professionally translated into 101 languages, enabling evaluation of many-to-many translation systems, particularly for low-resource language pairs (Guzmán et al., 2019; Goyal et al., 2022).

WMT24++. The WMT24++ benchmark covers 55 languages and dialects with human-written references and post-edits across four domains: literary, news, social, and speech. This diversity enables comprehensive evaluation of translation systems across high-resource and low-resource languages (Deutsch et al., 2025).

Both FLORES and WMT24++ serve as standard benchmarks in the machine translation community due to their broad language coverage and high-quality human translations, enabling meaningful comparisons across models and approaches.

We evaluate our proposed methods on English-Ukrainian translation using the COMET (Crosslingual Optimized Metric for Evaluation of Translation) metric. COMET leverages multilingual pre-trained models to assess translations by comparing source text, hypothesis, and reference, demonstrating higher correlation with human judgments than traditional metrics like BLEU or chrF++ (Rei et al., 2020). We report COMET system-level scores (aggregated across all translations) in Table 1, using the Unbabel/XCOMET-XL model for reference-based quality estimation (Guerreiro et al., 2023). For mentioned tasks we have used GPT-4o-mini-2024-07-18 model for translation.

From the presented results, we observe that USI and T-RANK demonstrate clear advantages over other methods; however, it remains unclear whether T-RANK consistently outperforms USI. This raises the question of whether a correlation exists between translation cost or effort and quality. We must account for the fact that COMET evaluations are not entirely reliable, even when reference translations are provided. Notably, both datasets used for evaluation (WMT24++ and FLORES) primarily contain short texts that do not adequately test complex grammatical structures, particularly for Ukrainian.

Method	WMT24++	FLORES
Baseline	0.827	0.937
SC (with check)	0.821	0.937
Best-of-N (n=5)	0.843	0.943
USI (n=5)	0.843	0.945
T-RANK (p=5)	0.845	0.940

Table 1: COMET (reference-based) translation quality scores for WMT24++ and FLORES EN→UK pair (GPT-4o-mini). n denotes number of sampled candidates, p denotes number of different translation prompts (for $p > 1$ we use $n = 1$).

Moreover, the choice of a correct translation often depends on personal preferences and stylistic conventions; a single source text may have multiple valid translation candidates, all grammatically correct. Therefore, reference translations in commonly used machine translation benchmarks cannot be considered absolute gold standards, and achieving the highest COMET scores does not guarantee that a method is error-free. This motivates the need for additional comparative tests under different evaluation paradigms, such as Quality Estimation (QE) or reference-free machine translation evaluation. As discussed in appendix A.5, USI is more suitable for short and simple dataset translation, whereas T-RANK shows better performance when translating benchmarks, especially when they have complex question structure, which might get "lost in translation" for languages explored in this paper.

For reference-free QE, we can directly use COMET to score translations without requiring an ideal reference translation. We compared USI and T-RANK methods using the QE setup on FLORES and MMLU datasets, which contain longer and more complex texts. In this work, we use Unbabel/wmt23-cometkiwi-da-xl model for reference-free QE COMET evaluation (Rei et al. (2023) and the detailed evaluations are provided in appendix A.5.

4.2 Multilingual Benchmark Translation Quality

Moreover, we also use LLM-as-a-judge to compare MMLU translations between standard industry-recognized source (Global-MMLU) and our most efficient translation method (T-RANK). Using Gemini-2.5-Flash with high reasoning mode as a judge, we showcase a clear advantage in translation quality using our proposed method, as presented

in Table 2. We observe that our translations are preferred also in other available languages like Romanian and Lithuanian, as shown in Appendix A.4.

Translation	Wins	Draws	Losses
Global-MMLU-UK	2016	3276	8750
T-RANK (ours)	8750	3276	2016

Table 2: Comparison of translation quality between Global-MMLU and our proposed T-RANK method for MMLU benchmark in Ukrainian language using LLM-as-a-judge.

We also compare evaluation results of mid-sized models like Gemma 3 4/12B, Qwen 3 8B and Llama 3.1 8B on our improved translations versus existing benchmark translations. As shown in Table 3, we observe higher scores on our translations, providing additional evidence of enhanced translation quality. Below, we provide average difference between evaluating on our translations and existing ones (positive difference means results are higher on our translations); more results for other languages are presented in Appendix A.3.

Benchmark	AVG Δ
ARC-Challenge	+2.35%
Hellaswag	+1.63%
MMLU	+0.94%
Winogrande	+3.42%

Table 3: Average improvement per benchmark across all languages

5 Conclusion

We present a novel automated translation framework designed to enable rapid, high-quality translation of datasets and benchmarks while minimizing human intervention. Our integrated methods demonstrate substantial improvements on WMT24++ and FLORES benchmarks as measured by COMET scores, with T-RANK and USI achieving the strongest performance through test-time compute strategies. The quality of our translated benchmarks is further validated through COMET Quality Estimation scores and LLM-as-a-judge evaluations, which confirm meaningful improvements over existing translations. Notably, evaluation results of mid-sized models such as Gemma 3, Qwen 3 and Llama 3.1 on our improved translations yield higher scores compared to existing

benchmark translations, providing additional evidence of enhanced translation quality. These findings demonstrate that our framework effectively balances translation accuracy, computational efficiency, and scalability, offering a practical solution for creating high-quality multilingual evaluation resources. Future work should explore adaptive method selection based on translation difficulty, integration of dedicated quality models, and comprehensive evaluation across open-weight models to further enhance the framework’s capabilities on languages beyond Europe.

Limitations

Our study has several limitations that warrant consideration. First, we employ LLM-based scoring rather than dedicated translation quality models like COMET for the Best-of-N selection process, which may affect the reliability of candidate ranking. Second, our approach applies uniform methods across all entries without automatically estimating translation difficulty per input, though adaptive method selection based on text complexity could potentially improve efficiency and quality. Machine translation benchmarks demonstrate that the most advanced and computationally expensive methods do not always yield superior results for shorter text sequences, particularly those not used in question-answering contexts. We defer to users of our framework to select methods based on their specific objectives and resource constraints. Third, we rely primarily on closed-source models for translation, with limited testing of open-weight alternatives. We hypothesize that our proposed methods would yield greater benefits for open-weight models, which typically demonstrate weaker performance in zero-shot translation settings, though this requires validation through comprehensive evaluation. Finally, while we focus on Eastern and Southern European languages, further research is needed to assess the generalizability of our framework across a broader range of low-resource languages with diverse linguistic characteristics.

Ethics Statement

This work aims to improve multilingual benchmark quality to enable more equitable evaluation of language models across diverse languages. We acknowledge several ethical considerations. First, our reliance on closed-source models raises repro-

ducibility concerns; we plan to extend evaluation to open-weight alternatives to promote accessibility. Second, while focusing on Eastern European languages addresses an important gap, many lower-resource languages remain underrepresented, and we encourage extending these methods to broader linguistic contexts. Third, we recognize that improved benchmarks may still be imperfect or subject to translation model bias; however, accurate multilingual evaluation is crucial for developing culturally aware and safer AI systems. Our automated translation approach avoids potential labor exploitation concerns associated with human translation.

The code and translated benchmarks produced in this work will be released under permissible open licenses to enable community validation and reproducibility. Generative AI systems were employed exclusively for language assistance, including paraphrasing, spell-checking, and stylistic refinement of the authors' original content, as well as for generating boilerplate code. All core research contributions, experimental design, and findings represent the original work of the authors. Used benchmarks and datasets are publicly available under permissive licensing; any ethical considerations related to their use are discussed in the original publications.

References

Xinyun Chen, Renat Aksitov, Uri Alon, Jie Ren, Kefan Xiao, Pengcheng Yin, Sushant Prakash, Charles Sutton, Xuezhi Wang, and Denny Zhou. 2023. Universal self-consistency for large language model generation.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. *Think you have solved question answering? try arc, the ai2 reasoning challenge*. *Preprint*, arXiv:1803.05457.

Daniel Deutsch, Eleftheria Briakou, Isaac Caswell, Mara Finkelstein, Rebecca Galor, Juraj Juraska, Geza Kovacs, Alison Lui, Ricardo Rei, Jason Riesa, Shruti Rijhwani, Parker Riley, Elizabeth Salesky, Firas Trajbsi, Stephanie Winkler, Biao Zhang, and Markus Freitag. 2025. *Wmt24++: Expanding the language coverage of wmt24 to 55 languages & dialects*. *Preprint*, arXiv:2502.12404.

Zhaopeng Feng, Yan Zhang, Hao Li, Bei Wu, Jiayu Liao, Wenqiang Liu, Jun Lang, Yang Feng, Jian Wu, and Zuozhu Liu. 2024. *Tear: Improving llm-based machine translation with systematic self-refinement*.

Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc'Aurelio Ranzato, Francisco Guzmán,

and Angela Fan. 2022. The flores-101 evaluation benchmark for low-resource and multilingual machine translation. *Transactions of the Association for Computational Linguistics*, 10:522–538.

Nuno M. Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André F. T. Martins. 2023. *xcomet: Transparent machine translation evaluation through fine-grained error detection*. *Preprint*, arXiv:2310.10482.

Francisco Guzmán, Peng-Jen Chen, Myle Ott, Juan Pino, Guillaume Lample, Philipp Koehn, Vishrav Chaudhary, and Marc'Aurelio Ranzato. 2019. The flores evaluation datasets for low-resource machine translation: Nepali-english and sinhala-english. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6098–6111. Association for Computational Linguistics.

Wenhan Han, Yifan Zhang, Zhixun Chen, Binbin Liu, Haobin Lin, Bingni Zhang, Taifeng Wang, Mykola Pechenizkiy, Meng Fang, and Yin Zheng. 2025. *Mubench: Assessment of multilingual capabilities of large language models across 61 languages*. *Preprint*, arXiv:2506.19468.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. *Measuring massive multitask language understanding*. *Preprint*, arXiv:2009.03300.

Ammar Khairi, Daniel D'souza, Marzieh Fadaee, and Julia Kreutzer. 2025a. *Making, not taking, the best of n*. *Preprint*, arXiv:2510.00931.

Ammar Khairi, Daniel D'souza, Ye Shen, Julia Kreutzer, and Sara Hooker. 2025b. *When life gives you samples: The benefits of scaling up inference compute for multilingual llms*. *Preprint*, arXiv:2506.20544.

Viet Dac Lai, Chien Van Nguyen, Nghia Trung Ngo, Thuat Nguyen, Franck Dernoncourt, Ryan A. Rossi, and Thien Huu Nguyen. 2023. *Okapi: Instruction-tuned large language models in multiple languages with reinforcement learning from human feedback*. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 318–327, Singapore. Association for Computational Linguistics.

Ricardo Rei, Nuno M. Guerreiro, José Pombal, Daan van Stigt, Marcos Treviso, Luisa Coheur, José G. C. de Souza, and André F. T. Martins. 2023. *Scaling up cometkiwi: Unbabel-ist 2023 submission for the quality estimation shared task*. *Preprint*, arXiv:2309.11925.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. *COMET: A neural framework for MT evaluation*. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.

818	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. Winogrande: An adversarial winograd schema challenge at scale . <i>Preprint</i> , arXiv:1907.10641.	868
819		869
820		870
821		871
822	Shivalika Singh, Angelika Romanou, Clémentine Fourrier, David I Adelani, Jian Gang Ngui, Daniel Vilasuerro, Peerat Limkonchotiawat, Kelly Marchisio, Wei Qi Leong, Yosephine Susanto, and 1 others. 2024. Global mmlu: Understanding and addressing cultural and linguistic biases in multilingual evaluation.	872
823		873
824		874
825		875
826		876
827		877
828	Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. 2020. Learning to summarize with human feedback. <i>Advances in Neural Information Processing Systems</i> , 33:3008–3021.	878
829		879
830		880
831		881
832		882
833		883
834	Lei Tang, Jinghui Qin, Wenxuan Ye, Hao Tan, and Zhijing Yang. 2025. Adaptive few-shot prompting for machine translation with pre-trained language models.	883
835		884
836		885
837		886
838	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? <i>Preprint</i> , arXiv:1905.07830.	887
839		888
840		889
841		890
842	A Appendix	
843	A.1 Translation Examples	
844	In this subsection we will provide some translation examples showcasing common errors and how our proposed methods address them. Figure 3 demonstrates how T-RANK identifies and corrects translation errors through competitive ranking. Unlike all other methods, we observe improved ability of the judge model to notice translation flaws when evaluating the translation in competitive manner. Figure 4 illustrates a common issue where translating questions and answers separately can leak the correct answer through grammatical gender markers. Figure 5 shows the correct approach where context is preserved during translation.	891
845		892
846		893
847		894
848		895
849		896
850		897
851		898
852		899
853		900
854		901
855		902
856		903
857	A.2 Benchmark Translation Statistics	904
858	In this subsection we provide statistics of selected benchmark translation, used splits and total number of samples used. We also provide information on models and methods used for final translation.	905
859		906
860		907
861		908
862		909
863		910
864		911
865		912
866		913
867		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

Candidate 1

Question:

Ви спостерігаєте вид, який народжує лише одне потомство за раз і має відносно тривалий життєвий цикл для свого розміру тіла. Що з наведеного, ймовірно, також вірно для цього організму?

Answers:

1. Він живе в новозаселеному середовищі.
2. Це водний організм.
3. Він потребує відносно високої батьківської турботи про потомство.
4. Вік, в якому потомство може самостійно народжувати, відносно молодий.

Candidate 2

Question:

Ви спостерігаєте вид, який народжує лише одне потомство за раз і має відносно тривалий термін життя для свого розміру. Що з наведеного, ймовірно, також вірно для цього організму?

Answers:

1. Він живе в новозаселеному середовищі.
2. Це водний організм.
3. Він потребує відносно великої батьківської турботи про потомство.
4. Вік, в якому потомство саме може народжувати, є відносно молодим.

Candidate 3

Question:

Ви спостерігаєте вид, який народжує лише одне потомство за раз і має відносно тривалий термін життя для свого розміру. Що з наведеного, ймовірно, також вірно для цього організму?

Answers:

1. Він живе в нещодавно колонізованому середовищі.
2. Це водний організм.
3. Він потребує відносно великої батьківської опіки за потомством.
4. Вік, в якому потомство може самостійно народжувати, відносно молодий.

Selected Translation

Question:

Ви спостерігаєте вид, який народжує лише одне потомство за раз і має відносно тривалий термін життя для свого розміру. Що з наведеного, ймовірно, також вірно для цього організму?

Answers:

1. Він живе в нещодавно колонізованому середовищі.
2. Це водний організм.
3. Він потребує відносно великої батьківської опіки за потомством.
4. Вік, в якому потомство може самостійно народжувати, відносно молодий.

Reasoning

The final selected translation is generally accurate and conveys the original meaning well. However, there are minor improvements needed in terms of word choice and grammatical structure. The phrase 'відносно тривалий термін життя' could be simplified to 'відносно довгий термін життя' for better clarity. Additionally, the phrase 'батьківської опіки за потомством' can be improved to 'батьківської опіки над потомством' for grammatical correctness. Overall, the translation is good but requires these minor adjustments for optimal clarity and correctness.

Final Corrected Translation

Question:

Ви спостерігаєте вид, який народжує лише одне потомство за раз і має відносно довгий термін життя для свого розміру. Що з наведеного, ймовірно, також вірно для цього організму?

Answers:

1. Він живе в нещодавно колонізованому середовищі.
2. Це водний організм.
3. Він потребує відносно великої батьківської опіки над потомством.
4. Вік, в якому потомство може самостійно народжувати, відносно молодий.

Figure 3: Translation Ranking (T-RANK) Correction Example

918 proved quality from using T-RANK method specif-
919 ically, as shown in first subsection of appendix.
920 Moreover, for Gemini-2.0-Flash model the rank-
921 ing capabilities seem to be stronger, which may
922 also contribute to better performance of T-RANK
923 method. This supports our idea, that different meth-

ods may be more suitable for different models and
use-cases.

A.6 Translation Prompts

Table 23 provides the complete prompt templates
used for all translation methods described in this

924
925
926
927
928

Вони хвилювалися, що вино зіпсує ліжко та ковдру, але _ не було зіпсовано.
 What does the blank _ refer to?
 Option A: ковдра
 Option B: ліжко
 Answer with A or B.
 Answer: ліжко

Figure 4: Winogrande Answer Leakage Example

Вони хвилювалися, що вино зіпсує ліжко та ковдру, але _ не бу(-в/-ла/-ло/-ли) зіпсовано.
 What does the blank _ refer to?
 Option A: ковдра
 Option B: ліжко
 Answer with A or B.
 Answer: ліжко ?

Figure 5: Winogrande Correct Translation Example

Benchmark	Train	Val	Dev	Test	All
MMLU	–	1531	285	14042	15858
Winogrande	–	1267	–	–	–
ARC-Challenge	1119	299	–	1172	2,291
Hellaswag	–	–	–	10042	10042
Total	1119	3097	285	25556	29757

Table 4: Number of examples in each benchmark split (– indicates split not available or not used).

Lang/Bench	UKR	SK	RO	LT	EST	BG	GR	TR
MMLU	4o-mini(T-RANK)	4o-mini(T-RANK)	4o-mini(T-RANK)	4o-mini(USI)	4o-mini(USI)	4o-mini(T-RANK)	4o-mini(USI)	4o-mini(USI)
Hellaswag	4o-mini	4o-mini	4o-mini	gemini-2.0-fl	gemini-2.0-fl	gemini-2.0-fl	gemini-2.0-fl	gemini-2.0-fl
ARC	4o-mini	gemini-2.0-fl	gemini-2.0-fl	gemini-2.0-fl	gemini-2.0-fl	gemini-2.0-fl	gemini-2.0-fl	gemini-2.0-fl
Winogrande	4o-mini	gemini-2.0-fl	gemini-2.0-fl	gemini-2.0-fl	gemini-2.0-fl	gemini-2.0-fl	gemini-2.0-fl	gemini-2.0-fl

Table 5: Selected model and method for final translation of benchmarks, if not specified, USI method was used. 4o-mini = GPT-4o-mini-2024-07-18, gemini-2.0-fl = Gemini-2.0-Flash.

Model	Hellaswag			Winogrande			ARC-Ch.			MMLU		
	O	Ot	Δ	O	Ot	Δ	O	Ot	Δ	O	Ot	Δ
Gemma-3-12B-IT	0.687	0.625	0.062	0.580	0.619	-0.039	0.517	0.470	0.048	0.614	0.603	0.012
Llama-3.1-8B	0.570	0.540	0.030	0.549	0.495	0.055	0.431	0.416	0.015	0.497	0.489	0.009
Gemma-3-4B-IT	0.578	0.528	0.050	0.540	0.501	0.039	0.453	0.417	0.037	0.453	0.444	0.010
Qwen3-8B-IT	0.540	0.520	0.020	0.569	0.546	0.023	0.448	0.404	0.044	0.619	0.599	0.020

Table 6: Ukrainian. O=Ours, Ot=Other (Okapi/MuBench/Global-MMLU)

Model	Hellaswag			Winogrande			ARC-Ch.			MMLU		
	O	Ot	Δ	O	Ot	Δ	O	Ot	Δ	O	Ot	Δ
Gemma-3-12B-IT	0.678	0.681	-0.002	0.661	0.589	0.072	0.486	0.454	0.032	0.628	0.614	0.014
Gemma-3-4B-IT	0.567	0.567	0.000	0.575	0.531	0.045	0.412	0.394	0.018	0.479	0.473	0.006
Llama-3.1-8B-IT	0.576	0.581	-0.005	0.616	0.521	0.095	0.358	0.357	0.002	0.529	0.519	0.009
Qwen3-8B-IT	0.538	0.537	0.001	0.588	0.566	0.023	0.368	0.359	0.008	0.650	0.632	0.018

Table 7: Romanian. O=Ours, Ot=Other (Okapi/MuBench/Global-MMLU)

Model	Winogrande			ARC-Challenge		
	Ours	Other	Δ	Ours	Other	Δ
Gemma-3-12B-IT	0.610	0.589	0.021	0.539	0.509	0.030
Gemma-3-4B-IT	0.554	0.517	0.037	0.448	0.444	0.004
Llama-3.1-8B-IT	0.547	0.515	0.032	0.443	0.403	0.040
Qwen3-8B-IT	0.555	0.559	-0.004	0.428	0.427	0.001

Table 8: Slovak. Other=MuBench/Okapi

Model	Winogrande			MMLU		
	Ours	Other	Δ	Ours	Other	Δ
Gemma-3-12B-IT	0.631	0.557	0.073	0.585	0.573	0.013
Gemma-3-4B-IT	0.530	0.506	0.023	0.433	0.410	0.023
Llama-3.1-8B-IT	0.527	0.498	0.029	0.431	0.417	0.014
Qwen3-8B-IT	0.551	0.530	0.021	0.553	0.541	0.012

Table 9: Lithuanian. Other=MuBench/Global-MMLU

Model	Winogrande		
	Ours	Other	Δ
Gemma-3-12B-IT	0.615	0.566	0.049
Gemma-3-4B-IT	0.541	0.516	0.025
Llama-3.1-8B-IT	0.543	0.509	0.034
Qwen3-8B-IT	0.517	0.503	0.014

Table 10: Estonian. Other=MuBench

Model	Winogrande			MMLU		
	Ours	Other	Δ	Ours	Other	Δ
Gemma-3-12B-IT	0.635	0.571	+0.065	0.587	0.581	+0.007
Gemma-3-4B-IT	0.583	0.514	+0.070	0.447	0.437	+0.010
Llama-3.1-8B-IT	0.573	0.518	+0.055	0.494	0.480	+0.014
Qwen3-8B-IT	0.545	0.571	-0.027	0.588	0.569	+0.018

Table 11: Turkish. Other=MuBench/Global-MMLU

Model	Winogrande			MMLU		
	Ours	Other	Δ	Ours	Other	Δ
Gemma-3-12B-IT	0.657	0.601	+0.056	0.593	0.580	+0.013
Gemma-3-4B-IT	0.598	0.518	+0.080	0.428	0.421	+0.007
Llama-3.1-8B-IT	0.593	0.520	+0.073	0.465	0.446	+0.020
Qwen3-8B-IT	0.581	0.528	+0.053	0.563	0.553	+0.011

Table 12: Greek. Other=MuBench/Global-MMLU

Model	Hellaswag			Winogrande			ARC-Challenge			MMLU		
	Ours	Other	Δ	Ours	Other	Δ	Ours	Other	Δ	Ours	Other	Δ
Gemma-3-12B-IT	0.708	0.690	+0.018	0.643	0.677	-0.035	0.526	0.495	+0.031	0.605	0.619	-0.013
Gemma-3-4B-IT	0.590	0.569	+0.022	0.588	0.595	-0.007	0.430	0.410	+0.020	0.454	0.469	-0.015
Llama-3.1-8B-IT	0.539	0.528	+0.011	0.575	0.605	-0.031	0.394	0.367	+0.027	0.481	0.487	-0.006
Qwen3-8B-IT	0.536	0.518	+0.018	0.602	0.629	-0.027	0.414	0.394	+0.020	0.619	0.617	+0.002

Table 13: Bulgarian. Other=INSAIT

Model	Winogrande		
	Ours	Other	Δ
Gemma-3-12B-IT	0.643	0.597	+0.046
Gemma-3-4B-IT	0.588	0.506	+0.082
Llama-3.1-8B-IT	0.575	0.505	+0.069
Qwen3-8B-IT	0.602	0.559	+0.043

Table 14: Bulgarian. Other=MuBench

Translation	Wins	Draws	Losses
Global-MMLU-RO	2646	3020	8376
T-RANK (ours)	8376	3020	2646

Table 15: Comparison of translation quality between Global-MMLU and our proposed T-RANK method for MMLU benchmark in Romanian language using Gemini-2.5-Flash as a judge.

Translation	Wins	Draws	Losses
Global-MMLU-LT	3478	3181	7382
USI (ours)	7382	3181	3478

Table 16: Comparison of translation quality between Global-MMLU and our proposed USI method for MMLU benchmark in Lithuanian language using Gemini-2.5-Flash as a judge.

Method	EN→UK	EN→SK	EN→RO	EN→LT	EN→ET	EN→BG	EN→TR	EN→EL
Baseline	0.726	0.741	0.882	0.741	0.788	0.751	0.718	0.802
SC	0.708	0.725	0.869	0.735	0.788	0.749	0.715	0.793
BoN n=5	0.743	0.756	0.895	0.767	0.809	0.788	0.736	0.819
USI n=5	0.750	0.759	0.899	0.766	0.806	0.791	0.733	0.817
T-RANK n=5	0.739	0.745	0.885	0.749	0.799	0.776	0.729	0.811
USI multi-prompt p=4	0.755	0.764	0.898	0.771	0.809	–	–	–
T-RANK multi-prompt p=4	0.742	0.756	0.883	0.753	0.797	–	–	–

Table 17: WMT QE (Quality Estimation) reference-free COMET scores for GPT-4o-mini

Method	EN→UK	EN→SK	EN→RO	EN→LT	EN→ET	EN→BG	EN→TR	EN→EL
Baseline	0.827	0.822	0.873	0.788	0.821	0.834	0.776	0.820
SC	0.821	0.817	0.869	0.790	0.822	0.000	0.834	0.821
BoN n=5	0.844	0.837	0.884	0.805	0.838	0.852	0.791	0.836
USI n=5	0.843	0.839	0.888	0.807	0.842	0.856	0.791	0.838
T-RANK n=5	0.840	0.834	0.885	0.803	0.836	0.855	0.793	0.837
USI multi-prompt p=4	0.849	0.847	0.891	0.817	0.848	–	–	–
T-RANK multi-prompt p=4	0.845	0.841	0.887	0.812	0.843	–	–	–

Table 18: WMT COMET reference-based scores for GPT-4o-mini

Method	EN→UK	EN→SK	EN→RO	EN→LT	EN→ET
Baseline	0.904	0.906	0.947	0.906	0.927
SC	0.902	0.908	0.950	0.911	0.935
BoN n=5	0.910	0.917	0.956	0.921	0.947
USI n=5	0.913	0.915	0.956	0.924	0.943
T-RANK n=5	0.899	0.906	0.952	0.909	0.932
USI multi-prompt p=4	0.919	0.921	0.961	0.932	0.948
T-RANK multi-prompt p=4	0.905	0.912	0.954	0.912	0.935

Table 19: FLORES QE (Quality Estimation) reference-free COMET scores for GPT-4o-mini

Method	EN→UK	EN→SK	EN→RO	EN→LT	EN→ET
Baseline	0.937	0.938	0.939	0.906	0.894
SC	0.937	0.937	0.936	0.911	0.902
BoN n=5	0.943	0.945	0.945	0.917	0.918
USI n=5	0.945	0.942	0.947	0.922	0.919
T-RANK n=5	0.938	0.937	0.943	0.910	0.903
USI multi-prompt p=4	0.947	0.946	0.949	0.928	0.918
T-RANK multi-prompt p=4	0.940	0.943	0.944	0.915	0.907

Table 20: FLORES COMET reference-based scores for GPT-4o-mini

Method	EN→UK	EN→SK	EN→RO	EN→LT	EN→ET	EN→BG	EN→TR	EN→EL
Baseline	0.688	0.703	0.777	0.704	0.737	0.708	0.683	0.719
SC	0.687	0.706	0.777	0.706	0.736	0.709	0.685	0.722
BoN n=5	0.687	0.707	0.780	0.710	0.743	0.713	0.689	0.723
USI multi-prompt p=4	0.698	0.716	0.785	0.718	0.753	0.721	0.697	0.731
T-RANK multi-prompt p=4	0.697	0.713	0.781	0.716	0.749	0.718	0.695	0.726

Table 21: WMT QE (Quality Estimation) reference-free COMET scores for Gemini-2.0-flash

Method	EN→UK	EN→SK	EN→RO	EN→LT	EN→ET	EN→BG	EN→TR	EN→EL
Baseline	0.828	0.829	0.867	0.805	0.847	0.840	0.778	0.821
SC	0.826	0.834	0.868	0.808	0.845	0.841	0.779	0.824
BoN n=5	0.828	0.836	0.870	0.814	0.852	0.843	0.785	0.823
USI multi-prompt p=4	0.841	0.843	0.881	0.826	0.862	0.856	0.797	0.836
T-RANK multi-prompt p=4	0.841	0.849	0.882	0.827	0.861	0.859	0.800	0.838

Table 22: WMT COMET reference-based scores for Gemini-2.0-flash

Method	Prompt Template
Base Translation	<p>Instructions: Imagine you're part of a team at an international education center that's revamping its exams for a global audience. Your job is to translate an English question and its answer options into <target_language> so that students from <target_language> schools can be evaluated too. Just provide the final translation—leave out any extra comments or explanations. Use language which is authentic for <target_language> natives. Remember to keep the answer options connected to the question, using the same format as the original (a list for multiple choices or plain text for a single answer). Please do not translate valid code in any of the programming languages.</p> <p>Original text: {"Original_question": "<original_question>", "Original_answers": "<original_answers>"}</p> <p>Output instructions: Now, please give your final translation in <target_language> exactly in this format, with only the translated content: {"Question": "your_translated_question", "Answers": "translated_answers"}</p>
USI Judge	<p>My task is to translate BENCHMARK questions with answers from English to <target_language>. Your task is to evaluate if the response preserves the original question idea and to verify the correctness of declension and conjunction of words in the target language.</p> <p>The original text in English is: Question: <original_question>, Answers: <original_answers></p> <p>I have generated the following responses: <responses></p> <p>Combine the best features from responses to form the best response from grammatical and coherent points of view. Look for: (1) Quality of translation, including grammatical correctness; (2) Domain knowledge - were the terms correctly translated? Were coding terms preserved?; (3) Is the question text fully translated?; (4) Are all answer options fully translated?; (5) Are the words written correctly?</p> <p>Output only the selected response: Question: selected question, Answer: selected answers</p>
T-RANK Ranking	<p>My task is to translate BENCHMARK questions from English to <target_language>. Your task is to rank my translations and select the best one.</p> <p>Ranking criteria: (1) Quality of translation; (2) Domain knowledge; (3) Is the question fully translated?; (4) Are all answer options translated?; (5) Correct spelling?; (6) Correct declension and conjunction?</p> <p>Original: {"original_question": "<original_question>", "original_answers": <original_answers>}</p> <p>Candidates: <responses></p> <p>Instruction: Select the best response (1st place = best). Correct if needed before output.</p> <p>Output: Reasoning, then: {"summary": "...", "final_ranks": {...}, "rankings_list": [...], "best_translation": {...}}</p>
Best-of-N Scoring	<p>My task is to translate questions from English to <target_language>. Score my translations 1-10 (10 = best).</p> <p>Original: Question: <original_question>, Answers: <original_answers></p> <p>Scoring metrics: (1) Translation quality; (2) Question fully translated?; (3) All answers translated?; (4) Question idea preserved?; (5) Correct grammar?</p> <p>Responses: <responses></p> <p>Output scores only: Response 1: score, Response 2: score, ..., Answers: [list of scores]</p>

Table 23: Translation prompt templates for different methods