

# GRRM: Group Relative Reward Modeling for Machine Translation

Anonymous ACL submission

## Abstract

While Group Relative Policy Optimization (GRPO) offers a powerful framework for LLM post-training, its effectiveness in open-ended domains like Machine Translation hinges on accurate intra-group ranking. We identify that standard Scalar Quality Metrics (SQM) fall short in this context; by evaluating candidates in isolation, they lack the comparative context necessary to distinguish fine-grained linguistic nuances. To address this, we introduce the Group Quality Metric (GQM) paradigm and instantiate it via the Group Relative Reward Model (GRRM). Unlike traditional independent scorers, GRRM processes the entire candidate group jointly, leveraging comparative analysis to rigorously resolve relative quality and adaptive granularity. Empirical evaluations confirm that GRRM achieves competitive ranking accuracy among all baselines. Building on this foundation, we integrate GRRM into the GRPO training loop to optimize the translation policy. Experimental results demonstrate that our framework not only improves general translation quality but also unlocks reasoning capabilities comparable to state-of-the-art reasoning models.

## 1 Introduction

Reinforcement learning with verifiable rewards (RLVR) has significantly advanced the reasoning processes of Large Language Models (LLMs), yielding remarkable performance in sophisticated domains such as mathematics and programming (Lambert et al., 2024; Guo et al., 2025a). Performing reinforcement learning on open-ended tasks, such as Machine Translation (MT), requires reward models to evaluate the responses generated by the policy model. Existing works predominantly employ Discriminative Reward Models (DRMs) based on the Bradley-Terry model (Cheng et al., 2025; Yang et al., 2025b). While DRMs have proven effective in improving translation quality

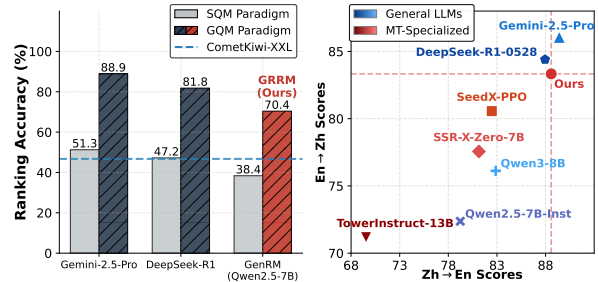


Figure 1: Performance on Seed-X-Challenge. **Left:** Ranking accuracy across paradigms. **Right:** Translation performance across General and MT-specialized LLMs.

in general domains, they have been criticized for lacking reasoning capabilities, which limits their effectiveness in challenging scenarios.

Recently, the success of LLM-based generative evaluation (LLM-as-a-Judge) in MT (Kocmi and Federmann, 2023a,b) suggests the potential of Generative Reward Models (GenRMs) (Zhang et al., 2025; Mahan et al., 2024). Ideally, the reasoning potential of GenRMs offers more accurate reward estimation, particularly for difficult samples. However, we identify a critical limitation when applying GenRMs within the widely used Group Relative Policy Optimization (GRPO) (Shao et al., 2024) framework. While LLM-based judges have proven reliable for system-level evaluation, they struggle with the fine-grained intra-group ranking required by GRPO. Since GRPO estimates advantages by comparing the relative quality of responses within a group, the accuracy of this local ranking is paramount. Our experiments reveal that standard generative approaches, which evaluate responses independently (Scalar Quality Metric, SQM), fail to consistently distinguish subtle quality differences among candidates. Worse, they may overlook critical semantic errors, as the evaluator can be misled by the candidate translation in challenging samples.

To address the limitations of the SQM paradigm, we propose the **Group Quality Metric (GQM)**. Specifically, GQM enables the model to process all

candidate translations within a group collectively, aiming to jointly estimate their quality ranking and assign relative scores. This cross-sample context provides consistent evaluation criteria, allowing the model to focus on distinguishing features between candidates. Furthermore, this paradigm allows the model to adaptively adjust its evaluation granularity according to the quality variance within the group. We implement this paradigm by integrating it with GenRM, denoted as the **Group Relative Reward Model (GRRM)**, which demonstrates a structural alignment with the GRPO framework.

We instantiate GRRM based on the Qwen2.5-7B model (Yang et al., 2024), which is cold-started via Supervised Fine-Tuning (SFT) on Chinese-English data and subsequently optimized for ranking accuracy via RLVR. We assess the group ranking performance of the GQM paradigm against SQM across various advanced LLMs, alongside our GRRM, using both LLM-annotated datasets and extensive human-annotated benchmarks, including specific challenge scenarios. Empirical results demonstrate that employing GQM consistently enables LLMs to outperform their SQM counterparts across diverse cross-lingual settings, verifying the universality of the proposed metric. Notably, on challenging datasets, this paradigm delivers significant absolute accuracy gains of 30% to 40% compared to SQM. Furthermore, our GRRM exhibits strong cross-lingual generalization; despite being trained on a single language pair, it effectively supports multilingual translation optimization without additional adaptation.

We further integrate GRRM with the GRPO framework for machine translation optimization, specifically focusing on the enhancement of the translation model’s reasoning capabilities. We bootstrap the Qwen2.5-7B model via a Chinese-English cold-start for translation with reasoning, and subsequently conduct GRPO training on 150k multilingual samples using GRRM feedback. For general domains, our model achieves an average 7.5-point improvement in BLEURT and a 16 point increase in LLM-judge scores over the SFT baseline across seven English-to-X (En2X) language pairs. Moreover, in challenging Chinese-English scenarios, our model performs comparably to DeepSeek-R1-0528 (Guo et al., 2025a). Our analysis reveals that GRRM fosters the emergence of reasoning capabilities in translation models, which is crucial for solving complex translation challenges. In contrast, while DRMs can improve general translation

quality, they are prone to reward hacking and lack the reasoning faculties required to guide models in challenging contexts.

## 2 Methodology

### 2.1 Ranking Sensitivity in GRPO

We employ GRPO, a reinforcement learning paradigm that has proven highly effective for LLM post-training, particularly in complex reasoning domains.

Formally, for each query  $x$ , the policy  $\pi_\theta$  samples a group of outputs  $\mathcal{Y} = \{y_1, \dots, y_G\}$ . The policy is optimized to maximize the expected advantage of the generations. Specifically, the advantage  $A_i$  for the  $i$ -th candidate is derived by standardizing its reward  $r_i$  against the group distribution:

$$A_i = \frac{r_i - \text{mean}(\{r_1, \dots, r_G\})}{\text{std}(\{r_1, \dots, r_G\})} \quad (1)$$

where  $r_i$  is the scalar reward score assigned to  $y_i$ .

Eq. 1 underscores a critical property of GRPO: the optimization trajectory is driven entirely by the relative quality of candidates within a group, rather than their absolute scores. A candidate  $y_i$  is reinforced ( $A_i > 0$ ) solely if it outperforms the group average. Consequently, the efficacy of GRPO is contingent on the reward model’s ability to accurately rank candidates within  $\mathcal{Y}$ . Even with well-calibrated absolute scores, any ranking inversion—where an inferior output scores higher than a superior one—yields adversarial gradients that actively mislead policy optimization. This sensitivity motivates our investigation into the ranking limitations of current reward modeling paradigms.

### 2.2 Limitations of Scalar Quality Metric

Current generative reward models typically function as a *Scalar Quality Metric (SQM)*, defined as a mapping  $S_{\text{SQM}} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  that assigns a scalar score  $r_i$  to a candidate  $y_i$  given an input  $x$ . Crucially, this evaluation is performed independently for each candidate. While this pointwise scoring is standard for general quality estimation, our analysis demonstrates that such independence fundamentally limits the model’s ability to capture the fine-grained relative distinctions required for effective GRPO updates.

**Inability to Distinguish Subtle Differences.** LLMs frequently generate candidates that differ only in fine-grained nuances. Without comparative

context, SQM suffers from discriminative failure, often assigning identical scores to variations of unequal quality (Fig. 2, Case 1). This “score collapse” results in a vanishing advantage in Eq. 1, stalling policy optimization with no useful gradient signal.

**Vulnerability to Misleading Candidates.** More critically, independent evaluation leaves the reward model vulnerable to misleading candidates. Even powerful models may overlook semantic discrepancies, defaulting to fluency or length as proxies for quality. As shown in Fig. 2, Case 2, SQM fails to penalize severe localization errors if they appear plausible in isolation. This results in ranking inversions, generating adversarial reward signals that actively degrade the policy.

### 2.3 Group Quality Metric

To overcome the structural limitations of SQM, we introduce the *Group Quality Metric* (GQM). Unlike the point-wise approach, GQM evaluates the policy-generated candidate group collectively. Formally, it maps a source  $x$  and a candidate set  $\mathcal{Y}$  to a score vector:  $S_{\text{GQM}} : \mathcal{X} \times \mathcal{Y}^G \rightarrow \mathbb{R}^G$ . This holistic evaluation shifts the paradigm from absolute estimation to *relative ranking*, offering three key advantages:

By evaluating candidates jointly, GQM constructs a comparative context that circumvents the discriminative bottlenecks of SQM. This context enables the model to resolve fine-grained nuances that remain imperceptible under independent evaluation (Fig. 2, Case 1). Simultaneously, GQM introduces contrastive anchoring: plausible but erroneous translations (e.g., hallucinations) become salient when directly compared against high-quality peers (Fig. 2, Case 2,4). This mechanism ensures the reward signal captures both subtle stylistic preferences and critical semantic errors.

Crucially, GQM scores are inherently localized and do not map to a global absolute scale. This allows the model to adaptively adjust evaluation granularity: shifting focus to minor stylistic differences when group variance is low, while prioritizing major semantic discrepancies when significant quality divergences exist.

### 2.4 Group Relative Reward Model

While the GQM paradigm demonstrates efficacy when powered by frontier LLMs, directly integrating such massive models into the GRPO training loop is computationally intractable due to high la-

tency and inference costs. We therefore introduce the Group Relative Reward Model (GRRM), an efficient instantiation of the GQM paradigm designed for high-throughput training loops.

GRRM is designed to process the source input  $x$  and the full candidate set  $\mathcal{Y}$  jointly. For every tuple  $(x, \mathcal{Y})$ , the model is instructed to produce a structured output comprising: (1) a detailed comparative analysis of the candidates, (2) a predicted ranking, and (3) a set of scalar scores  $\{r_1, \dots, r_G\}$  that strictly adhere to the predicted ranking. This design enforces a Chain-of-Thought (CoT) reasoning process, ensuring numerical scores are grounded in explicit comparative analysis.

We initialize GRRM via SFT on a curated dataset to establish instruction adherence and comparative logic. Subsequently, to sharpen the model’s discriminative ability to distinguish subtle nuances, we employ RLVR to directly optimize its ranking accuracy. In this stage, we treat the GRRM itself as the policy to be optimized. To drive this optimization, we define a verifiable reward function, Ranking Accuracy, based on ground-truth orderings. Let  $\mathbf{q} = \{q_1, \dots, q_G\}$  be the ground-truth quality scores (or ranks) for the candidate group, and  $\mathbf{r} = \{r_1, \dots, r_G\}$  be the corresponding scores predicted by GRRM. We define the correctness of the pairwise relationship between any two candidates  $y_i$  and  $y_j$  as:

$$\mathbb{I}_{ij} = \mathbb{I}[\text{sgn}(\hat{r}_i - \hat{r}_j) = \text{sgn}(q_i - q_j)] \quad (2)$$

where  $\text{sgn}(\cdot)$  is the sign function, and  $\mathbb{I}(\cdot)$  is the indicator function. This formulation enforces strict alignment, requiring the predicted relationship—including ties—to match the ground truth. The final ranking accuracy reward  $R_{\text{acc}}$  is then calculated as the fraction of correctly ordered pairs over all distinct combinations in the group:

$$R_{\text{acc}} = \frac{1}{\binom{G}{2}} \sum_{1 \leq i < j \leq G} \mathbb{I}_{ij} \quad (3)$$

Additionally, to complement the sparse ranking signal and ensure the predicted scalar scores are well-distributed, we employ an auxiliary score consistency reward, the implementation details of which are provided in Appendix A.

## 3 GRRM Evaluation

In this section, we provide a comprehensive evaluation of the proposed GRRM within the GQM

framework, benchmarking the ranking accuracy (Eq. 3) against multiple baselines. Our experiments span three distinct settings: in-domain performance, multilingual generalization, and complex reasoning scenarios designed to stress-test reward modeling capabilities.

### 3.1 Experimental Setup

**Datasets and Benchmarks.** We construct our training data using the Chinese-English subset of TowerBlocks (Alves et al., 2024), comprising approximately 18.8k samples. We sample 2-4 translation candidates per source and annotate them with Gemini-2.5-Pro (Comanici et al., 2025) under the GQM paradigm. A held-out set of 512 samples is reserved as the model-annotated test set. For RLVR training, we utilize the same underlying data as SFT, augmenting it by shuffling candidate permutations and subsampling from the original groups (see Appendix B.1).

We utilize Newstest2020 (Zh→En) with pSQM (professional SQM) labels (Freitag et al., 2021) from the WMT20 metrics task. To evaluate cross-lingual generalization, we employ GeneralMT2022 with MQM (Multidimensional Quality Metric) annotations (Freitag et al., 2021) from WMT22. Beyond the in-domain Zh↔En pairs, this dataset includes English-to-German (En→De) and English-to-Russian (En→Ru) subsets. For these datasets, we retain up to 4 system outputs per source and derive ground-truth rankings based on the human-annotated scores.

To assess reasoning capabilities in complex linguistic scenarios, we construct a challenge set based on the Seed-X-Challenge test set (Cheng et al., 2025), which features idioms, slang, and domain-specific terminology in Zh↔En pairs. We design a binary ranking task to test the model’s reasoning capability: for each source, we sample 4 translation candidates and select the best one according to BLEURT (Sellam et al., 2020). We then pair this chosen candidate with the expert human reference. Assuming the expert reference is strictly superior to the weak baseline’s output, this setup requires the reward model to correctly identify the human translation as the winner.

**Baselines.** We include three categories of reward models for comparison. Unless otherwise noted, all trained models are initialized from Qwen2.5-7B.

- **LLM-as-a-Judge:** We test Gemini-2.5-Pro and DeepSeek-R1-0528. We evaluate them

using both the standard SQM prompting and our proposed GQM prompting, explicitly instructing the models to use CoT reasoning.

- **DRMs:** We include CometKiwi-XXL (Rei et al., 2023), a leading reference-free Quality Estimation metric. Additionally, we train a BT-RM using the standard Bradley-Terry loss on pairwise preferences derived from our Gemini-annotated training data.
- **GenRMs:** In addition to our GRRM, we train a SQM-GenRM baseline using Gemini SQM annotations, following an RLVR optimization process similar to that of GRRM.

See Appendix B.1 for details on data construction and training hyperparameters.

### 3.2 Performance of GRRM

Table 1 presents the ranking accuracy of various reward models across internal, general, and challenging benchmarks. We summarize our key findings as follows:

**Dominance of GQM over SQM.** The most prominent observation is the consistent superiority of the GQM paradigm over the traditional SQM approach across all evaluator backbones. For LLM-as-a-Judge, employing GQM yields substantial improvements; for instance, Gemini-2.5-Pro achieves an average accuracy of 70.11% with GQM, compared to just 56.17% with SQM. Notably, this performance gap is maximized in the challenge scenario. While SQM performance for even state-of-the-art models degrades to near-random levels (e.g., 47.22% for DeepSeek-R1), GQM maintains robust accuracy. This trend holds for trained reward models as well, where GRRM (based on GQM) significantly outperforms SQM-GenRM. These results validate that evaluating translation candidates independently (SQM) fails to capture fine-grained relative differences. By processing candidates collectively, GQM provides the model with a comparative context, enabling more robust and consistent ranking.

**Effectiveness of GRRM and Reasoning.** Our proposed GRRM achieves the highest average accuracy among all trained models. While we observe that RLVR optimization brings consistent performance gains over the SFT cold-start, the magnitude of improvement is less significant than the paradigm shift from SQM to GQM. Crucially, the advantage of GRRM is most pronounced in

Model	Paradigm	Internal	NT20	GenMT22 (MQM)				Seed-X-Challenge	Avg.
		Zh↔En	Zh→En	En→Zh	Zh→En	En→De	En→Ru	Zh↔En	
Random	-	43.47	39.42	43.82	38.74	35.91	37.64	44.70	40.04
<i>LLM-as-a-Judge (w/ Reasoning)</i>									
Gemini-2.5-Pro	SQM	70.28	53.56	46.41	61.99	58.92	64.88	51.26	56.17
Gemini-2.5-Pro	GQM	-	<b>62.60</b>	<b>64.58</b>	<b>67.75</b>	<b>65.61</b>	<b>71.23</b>	<b>88.89</b>	<b>70.11</b>
DeepSeek-R1-0528	SQM	66.11	48.42	43.67	58.09	53.64	59.13	47.22	51.69
DeepSeek-R1-0528	GQM	<b>80.92</b>	61.98	64.38	65.79	63.26	69.15	81.82	67.73
<i>Discriminative RMs (w/o Reasoning)</i>									
CometKiwi-XXL	SQM	72.01	57.82	<b>66.49</b>	61.60	<b>61.20</b>	<b>67.12</b>	46.72	60.16
BT-RM	SQM	<b>82.62</b>	<b>58.16</b>	<b>66.49</b>	<b>64.92</b>	58.14	66.10	<b>58.84</b>	<b>62.11</b>
<i>Generative RMs (w/ Reasoning)</i>									
SQM-GenRM (SFT)	SQM	61.31	47.18	37.21	57.68	48.37	51.84	38.13	46.74
SQM-GenRM (RLVR)	SQM	64.25	49.21	39.42	60.37	49.16	54.82	38.38	48.56
<b>GRRM</b> (SFT)	GQM	79.75	56.77	59.45	64.29	58.65	64.80	69.95	62.32
<b>GRRM</b> (RLVR)	GQM	<b>82.58</b>	<b>57.77</b>	<b>62.17</b>	<b>66.00</b>	<b>61.04</b>	<b>66.67</b>	<b>70.39</b>	<b>64.01</b>

Table 1: Ranking accuracy (%) on internal (Gemini-annotated), human-annotated, and challenge benchmarks. **Avg.** denotes the mean accuracy across all human-annotated and challenge datasets (excluding Internal). The best performance within each model category is **bolded**.

the Seed-X Challenge set, achieving a ranking accuracy of 70.39%. This indicates that the reasoning capability is essential for verifying complex linguistic phenomena that statistical correlations alone cannot capture. Additionally, despite being trained solely on Zh-En data, GRRM exhibits strong cross-lingual generalization to unseen En→De and En→Ru pairs, performing on par with the CometKiwi-XXL baseline.

**Limitations of Discriminative RMs.** We observe that DRMs perform competitively on general domain benchmarks. This can be attributed to the Bradley-Terry loss, which implicitly and effectively models the ranking objective on preference pairs. However, their performance collapses on the challenge set due to the lack of generative reasoning capabilities, limiting their utility in guiding policies through complex optimization landscapes.

## 4 MT Optimization with GRRM

Building upon the validation of GRRM’s ranking capabilities, we integrate it into the training loop to optimize the translation policy. This section details the training pipeline, benchmarks, metrics, and baselines used to assess performance and reasoning emergence.

### 4.1 Experimental Setup

**Training Pipeline.** We adopt a two-stage strategy initialized with Qwen2.5-7B. First, we perform SFT to set up preliminary translation and reasoning skills using the same Chinese-English data de-

scribed in Section 3 (with Gemini-2.5-Pro annotated CoT). Second, we optimize the model using GRPO with GRRM feedback based on the multilingual translation data from TowerBlocks, covering 10 languages with approximately 150k samples.

Since GRRM yields reference-free rewards, we implement Cross-Lingual Augmentation (CLA) by pairing source sentences with alternative target languages (e.g., Zh→De) from the dataset, while maintaining the total number of update steps constant.

**Datasets.** Our primary evaluation targets Zh↔En translation using WMT23 Zh→En (Kocmi et al., 2023), WMT24++ En→Zh (Deutsch et al., 2025), and the Seed-X-Challenge. For En→X, we report average performance across seven languages from WMT24++. Detailed results for WMT24++ En→X and WMT23 X→En are provided in Appendix C.1.

**Evaluation Metrics.** We employ a combination of standard metrics and LLM-based evaluation to provide a comprehensive assessment:

- **BLEURT:** We use BLEURT-20 as our primary automatic metric to measure semantic preservation against human references.
- **LLM-as-a-Judge:** Despite our findings in Section 3 regarding the superiority of GQM for ranking, we utilize the SQM paradigm for system-level evaluation. We justify this choice based on two key factors: (1) **Reference Anchoring:** Our evaluation here includes human references. The reference serves as a strong anchor (conceptually similar to a group size of

2 in GQM), significantly stabilizing the judgment. (2) **System-Level Accuracy:** While SQM struggles with fine-grained intra-group ranking, it remains reliable for aggregating scores at the system level.

For Seed-X-Challenge, we further enhance the evaluation prompt by including the expert annotations provided with the dataset. These annotations highlight specific translation difficulties, enabling the judge to perform a more informed assessment alongside the reference.

- **Evaluator Models:** To balance overhead, we employ DeepSeek-R1-0528 as the judge for all main results and gpt-oss-120b for ablation studies and supplementary results.

**Baselines.** We incorporate following baselines:

- **General LLMs:** We test Gemini-2.5-Pro and DeepSeek-R1-0528, along with Qwen3-8B (Yang et al., 2025a) and Qwen2.5-7B-Instruct (Yang et al., 2024) to benchmark improvements against the base model family.
- **Translation-Specialized Models:** We compare against three state-of-the-art open models: (1) TowerInstruct-13B (Alves et al., 2024), a LLaMA2 (Touvron et al., 2023) derivative specifically adapted for translation tasks through continued pre-training and instruction tuning; (2) SeedX-PP0 (Cheng et al., 2025), a 7B model pre-trained on high-quality multilingual datasets and subsequently enhanced via RL; and (3) SSR-X-Zero-7B (Yang et al., 2025c), a Qwen2.5-7B derivative which utilizes a self-rewarding RL framework specifically optimized for Zh↔En translation tasks.

**Reasoning Configuration.** For General LLMs, we explicitly prompt them to CoT reasoning to maximize their potential. For specialized models, SeedX-PP0 and SSR-X-Zero-7B leverage reasoning, while TowerInstruct-13B support direct translation only.

## 4.2 Main Results

Table 2 presents the performance of our GRRM-optimized models compared to the SFT baseline and other state-of-the-art systems. The results indicate that integrating GRRM into the GRPO framework significantly enhances translation quality, particularly in scenarios requiring reasoning.

Compared to the Qwen2.5-7B-SFT baseline, our method achieves substantial improvements across all language pairs. In the general domain (WMT benchmarks), the model achieves an average gain of **+7.5 BLEURT** points and **+15.9 LLM-judge** points on En→X tasks. This successful optimization on multilingual data validates the strong cross-lingual generalization of GRRM, consistent with our metric analysis in Section 3. Furthermore, in the Zh↔En tasks, our approach comprehensively surpasses all *Translation-Specialized Models* (including SeedX-PP0 and SSR-X-Zero-7B) in terms of LLM-judge scores, narrowing the gap with much larger proprietary models like Gemini-2.5-Pro.

The advantages of our approach are most pronounced in the challenge scenarios, i.e., Seed-X-Challenge. Our model equipped with Cross-Lingual Augmentation (CLA) achieves parity with the powerful reasoning model DeepSeek-R1-0528, even though the latter serves as our evaluator. Specifically, in the Zh→En subset, our model slightly outperforms DeepSeek-R1-0528 (88.58 vs. 87.95), while remaining highly competitive in the En→Zh subset (83.33 vs. 84.40). We provide detailed case studies in Appendix C.3 to qualitatively demonstrate how our model successfully navigates these challenging samples with reasoning.

Our results also reveal a notable divergence between BLEURT and LLM-as-a-Judge scores. For instance, while SeedX-PP0 achieves the highest BLEURT scores in several settings, its performance drops significantly under the scrutiny of the reasoning-based LLM judge. This discrepancy echoes our observations in Section 3, highlighting the limitations of traditional discriminative metrics like BLEURT when evaluating translation quality in challenging scenarios.

## 4.3 Ablation Study

We conduct ablation studies focusing on three key aspects: the choice of reward signals, the impact of training data distribution, and the role of reasoning in the policy model (see Table 3). Unlike our reasoning baseline which is fine-tuned solely on Zh-En data, we train the non-reasoning baseline on the full multilingual TowerBlocks (MT) dataset. This deviation is necessary because we find that the non-reasoning model overfits to Zh-En under the restricted setting, failing to generalize to other languages during RL. We place extended analyses on the non-reasoning baseline and Qwen2.5-7B-Instruct in Appendix C.2.

Model	WMT Benchmarks						Seed-X-Challenge			
	Zh→En		En→Zh		En→X		Zh→En		En→Zh	
	BLEURT	R1-judge	BLEURT	R1-judge	BLEURT	R1-judge	BLEURT	R1-judge	BLEURT	R1-score
<b>General LLMs</b>										
Gemini-2.5-Pro	<b>68.66</b>	<b>92.92</b>	<b>66.00</b>	<b>91.31</b>	<b>68.87</b>	<b>90.35</b>	<b>71.59</b>	<b>89.41</b>	<b>69.19</b>	<b>86.06</b>
DeepSeek-R1-0528	67.78	92.34	64.87	89.24	67.72	88.48	70.92	87.95	68.23	84.40
Qwen3-8B	63.03	89.72	57.58	84.15	57.25	77.37	61.17	82.88	57.78	76.12
Qwen2.5-7B-Instruct	67.31	88.49	59.92	80.51	58.72	72.51	66.59	79.23	62.75	72.37
<b>Translation-Specialized Models</b>										
TowerInstruct-13B	67.56	84.83	62.92	77.63	66.61	82.68	63.32	69.54	63.46	71.17
SeedX-PPO	<b>69.02</b>	90.47	<b>67.21</b>	87.98	<b>68.35</b>	<b>86.04</b>	69.37	82.47	<b>68.72</b>	80.56
SSR-X-Zero-7B	68.30	88.67	66.12	83.78	-	-	68.84	81.15	67.08	77.56
Qwen2.5-7B-SFT	67.07	87.78	59.99	76.98	57.14	67.91	67.65	80.91	62.36	72.42
+ GRPO (ours)	67.41	<b>92.24</b>	64.80	87.80	64.65	83.86	<b>69.55</b>	85.90	67.05	82.55
+ GRPO w/ CLA (ours)	67.39	92.09	63.91	<b>88.29</b>	64.50	83.71	69.25	<b>88.58</b>	67.07	<b>83.33</b>

Table 2: **MT performance on WMT and Seed-X-Challenge benchmarks.** We report BLEURT-20 and LLM-as-a-Judge scores evaluated by DeepSeek-R1-0528. The best performance within each category is highlighted in **bold**.

Ablation Configuration	WMT Benchmarks						Seed-X-Challenge			
	Zh→En		En→Zh		En→X		Zh→En		En→Zh	
	BLEURT	oss-judge	BLEURT	oss-judge	BLEURT	oss-judge	BLEURT	oss-judge	BLEURT	oss-judge
<b>Comparison of Reward Models</b>										
SFT baseline (w/ Reasoning)	67.07	85.50	59.99	76.80	57.14	65.71	67.65	80.76	62.36	74.30
+ GRPO w/ BLEURT	<b>68.71</b>	88.21	<b>67.09</b>	85.52	<b>66.55</b>	82.00	69.69	83.35	<b>68.66</b>	81.85
+ GRPO w/ BT-RM	68.06	88.86	60.20	76.99	56.53	64.93	<b>69.74</b>	85.60	58.25	66.91
+ GRPO w/ SQM-GenRM	68.63	89.01	44.28	45.54	27.53	8.15	69.85	84.25	47.19	46.42
+ GRPO w/ GRRM	67.41	<b>89.85</b>	64.80	<b>87.91</b>	64.65	<b>83.86</b>	69.55	<b>87.30</b>	67.05	<b>84.36</b>
<b>Cross-Lingual Generalization</b>										
+ GRPO w/ GRRM (ZhEn)	66.53	89.11	63.07	86.88	61.70	79.18	69.26	86.16	66.20	83.86
<b>Role of Reasoning</b>										
SFT baseline (w/o Reasoning)	67.61	85.28	63.50	79.42	63.24	75.33	65.52	76.83	64.04	73.34
+ GRPO w/ GRRM	67.38	88.42	64.15	85.19	63.40	78.88	68.57	85.38	65.92	79.35

Table 3: **Ablation study on Reward Models and Reasoning configurations.** We compare different reward signals and training setups. We report BLEURT-20 and LLM-as-a-Judge scores evaluated by gpt-oss-120b.

**Comparison of Reward Models.** Optimizing directly against BLEURT yields the highest scores on the BLEURT metric itself, as expected. However, this gain does not translate effectively to LLM-judge scores, particularly on the Seed-X-Challenge (e.g., 81.85 vs. 84.36 on En→Zh). This suggests that discriminative metrics fail to guide the model through the complex reasoning processes required for challenging translations.

More critically, we observe severe reward hacking with BT-RM and SQM-GenRM. These models perform catastrophically on En→Zh and En→X tasks. We find that these models suffer from off-target generation issues: apt to generate English responses for non-English targets. This indicates that standard discriminative or scalar generative rewards struggle to distinguish instruction-following failures (e.g., wrong language) from translation

quality issues, allowing the policy to exploit the reward model. In contrast, GRRM provides robust feedback, effectively penalizing such deviations and preventing reward hacking.

**Cross-Lingual Generalization** We examine whether the cross-lingual generalization observed in our reward model (Section 3) transfers to Policy Optimization. We train a variant using only Chinese-English data, i.e., GRPO w/ GRRM (ZhEn). While it improves performance on the Zh↔En pair, its performance on En→X lags significantly behind the model trained on multilingual data (79.18 vs. 83.86). This highlights a distinct difference between evaluation and generation: while a Reward Model trained on a single language pair can generalize to evaluate other languages, a Translation Policy requires explicit multilingual exposure to master the generation of diverse target languages.

**The Necessity of Reasoning.** Finally, we assess the contribution of the reasoning process (CoT) to translation quality. It is evident that models without reasoning capabilities consistently underperform their reasoning counterparts. Even when optimized with GRRM, the non-reasoning model scores lower on both general benchmarks and challenge sets (e.g., 79.35 vs. 84.36 on Seed-X En→Zh). This validates the necessity of a “Reasoning × Reasoning” paradigm: the full potential of our framework is unlocked only when a reasoning-capable policy is paired with a reasoning-based reward model (GRRM), enabling the system to effectively plan, self-correct, and evaluate in complex scenarios.

## 5 Related Work

**Reasoning for Machine Translation** The integration of reasoning capabilities into Machine Translation has emerged as a promising direction (Chen et al., 2025; Liu et al., 2025a). Previous attempts to integrate reasoning into MT primarily relied on prompting or SFT (Feng et al., 2025b; Wang et al., 2025; Cheng et al., 2025). However, Zebaze et al. (2025) observed that fine-tuning on synthetic CoT data alone often fails to outperform standard fine-tuning. Fundamentally, these SFT-based methods are not only constrained by the scarcity of high-quality supervision but also struggle to cultivate intrinsic reasoning capabilities.

Inspired by the success of reinforcement learning in mathematical and coding domains, recent works have adopted RL to encourage self-emergent reasoning. He et al. (2025) and Feng et al. (2025a) adapt the GRPO framework using discriminative metrics like Comet (Rei et al., 2020) and CometKiwi (Rei et al., 2023) to guide the model. Similarly, SSR-Zero (Yang et al., 2025c) employs the model itself to assign scalar quality scores combined with Comet to estimate advantages, demonstrating that self-evaluation can drive policy improvement. However, these works employ reward mechanisms that are suboptimal for fostering robust reasoning. DRMs based methods lack the generative explanatory power to guide models through complex errors and are prone to reward hacking. The scalar evaluation mechanism employed by SSR-Zero fails to capture the fine-grained relative rankings within a group, limiting the accuracy of advantage estimation of GRPO framework.

**Comparative and Reasoning-based Reward Modeling** While machine translation evaluation

has traditionally relied on SQM paradigm, employing pairwise comparative paradigms for judgment and reward estimation is not new in general open-ended domains (Li et al., 2024; Zhu et al., 2025). Ye et al. (2025) leverages self-generated contrastive pairs to train judges via DPO, enhancing robustness against bias compared to scalar models. Liu et al. (2025b) introduces a pairwise judge combined with a knockout tournament for Best-of-N sampling. Furthermore, integrating CoT into evaluation has proven critical; Guo et al. (2025b) demonstrates that executing a deliberate reasoning process before scoring improves accuracy by utilizing test-time compute.

However, these comparative approaches have yet to be effectively adapted to the group-wise context required by the GRPO framework. Existing pairwise methods are structurally misaligned with GRPO: estimating advantages for a group of responses via pairwise judges necessitates tournament-style evaluations (e.g., ELO ratings), which incur prohibitive computational costs. In contrast, our GRRM extends the comparative intuition to a list-wise context. By enabling the model to rank all candidates collectively in a single pass, our approach avoids the overhead of iterative comparisons and proves highly efficient, empirically achieving  $1.5\times$  faster rewarding speeds than standard scalar generative model.

## 6 Conclusion

In this work, we identify and address the limitations of scalar generative judges in GRPO training, specifically their insufficient sensitivity to distinguish intra-group quality differences in machine translation. To overcome this, we propose the Group Quality Metric (GQM), a paradigm that evaluates candidates jointly to capture fine-grained distinctions often missed by independent scalar metrics. Building on this foundation, we introduce the Group Relative Reward Model (GRRM). Our experiments demonstrate that GRRM significantly outperforms existing discriminative and scalar generative baselines, achieving robust ranking accuracy across diverse languages and challenging reasoning tasks. Crucially, integrating GRRM into the GRPO loop catalyzes the emergence of reasoning capabilities in translation models. This enables our approach to rival state-of-the-art systems like DeepSeek-R1 in complex scenarios while effectively mitigating reward hacking.

## 661 Limitations

662 We acknowledge two primary limitations of the pro-  
663 posed Group Quality Metric (GQM) and GRRM.  
664 First, unlike SQM which scales linearly, GQM is  
665 constrained by the maximum group size covered  
666 during training. Extrapolating to group sizes signifi-  
667 cantly beyond the training distribution may degrade  
668 ranking accuracy. However, this limitation does  
669 not hinder effective GRPO training. Recent work  
670 demonstrates that using small group sizes com-  
671 bined with larger batch sizes yields performance  
672 comparable to larger group configurations (Wu  
673 et al., 2025), rendering GQM’s capacity sufficient  
674 for optimization.

675 Second, GQM relies on relative ranking, which  
676 introduces challenges in reward assignment when  
677 the overall group quality is low. Since GRRM  
678 treats intra-group ranking as the ground truth for  
679 advantage estimation, it may assign high relative  
680 scores to the "best" candidate even if all transla-  
681 tions in the group are suboptimal. Consequently,  
682 GQM might fail to penalize the group globally as  
683 effectively as an absolute metric would. Integrat-  
684 ing absolute quality constraints alongside relative  
685 ranking remains a promising direction for future  
686 work.

## 687 References

688 Duarte M Alves, José Pombal, Nuno M Guerreiro, Pe-  
689 dro H Martins, João Alves, Amin Farajian, Ben Pe-  
690 ters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal,  
691 and 1 others. 2024. Tower: An open multilingual  
692 large language model for translation-related tasks.  
693 *arXiv preprint arXiv:2402.17733*.

694 Andong Chen, Yuchen Song, Wenxin Zhu, Kehai Chen,  
695 Muyun Yang, Tiejun Zhao, and 1 others. 2025. Evalu-  
696 ating o1-like llms: Unlocking reasoning for transla-  
697 tion through comprehensive analysis. *arXiv preprint*  
698 *arXiv:2502.11544*.

699 Shanbo Cheng, Yu Bao, Qian Cao, Luyang Huang,  
700 Liyan Kang, Zhicheng Liu, Yu Lu, Wenhao Zhu, Jing-  
701 wen Chen, Zhichao Huang, and 1 others. 2025. Seed-  
702 x: Building strong multilingual translation llm with  
703 7b parameters. *arXiv preprint arXiv:2507.13618*.

704 Gheorghe Comanici, Eric Bieber, Mike Schaekermann,  
705 Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Mar-  
706 cel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and  
707 1 others. 2025. Gemini 2.5: Pushing the frontier with  
708 advanced reasoning, multimodality, long context, and  
709 next generation agentic capabilities. *arXiv preprint*  
710 *arXiv:2507.06261*.

711 Daniel Deutsch, Eleftheria Briakou, Isaac Caswell,  
712 Mara Finkelstein, Rebecca Galor, Juraj Juraska, Geza

Kovacs, Alison Lui, Ricardo Rei, Jason Riesa, Shruti 713  
Rijhwani, Parker Riley, Elizabeth Salesky, Firas Tra- 714  
belsi, Stephanie Winkler, Biao Zhang, and Markus 715  
Freitag. 2025. WMT24++: Expanding the Language 716  
Coverage of WMT24 to 55 Languages & Dialects. 717  
*Preprint*, arXiv:2502.12404. 718

Zhaopeng Feng, Shaosheng Cao, Jiahao Ren, Jiayuan 719  
Su, Ruizhe Chen, Yan Zhang, Jian Wu, and Zuozhu 720  
Liu. 2025a. MT-r1-zero: Advancing LLM-based 721  
machine translation via r1-zero-like reinforcement 722  
learning. In *Findings of the Association for Computa-* 723  
*tional Linguistics: EMNLP 2025*, pages 18685– 724  
18702, Suzhou, China. Association for Computa- 725  
tional Linguistics. 726

Zhaopeng Feng, Yan Zhang, Hao Li, Bei Wu, Jiayu Liao, 727  
Wenqiang Liu, Jun Lang, Yang Feng, Jian Wu, and 728  
Zuozhu Liu. 2025b. TEaR: Improving LLM-based 729  
machine translation with systematic self-refinement. 730  
In *Findings of the Association for Computational* 731  
*Linguistics: NAACL 2025*, pages 3922–3938, Al- 732  
buquerque, New Mexico. Association for Computa- 733  
tional Linguistics. 734

Markus Freitag, George Foster, David Grangier, Viresh 735  
Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021. 736  
Experts, errors, and context: A large-scale study of 737  
human evaluation for machine translation. *Preprint*, 738  
arXiv:2104.14478. 739

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, 740  
Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong 741  
Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025a. 742  
Deepseek-r1: Incentivizing reasoning capability in 743  
llms via reinforcement learning. *arXiv preprint* 744  
*arXiv:2501.12948*. 745

Jiaxin Guo, Zewen Chi, Li Dong, Qingxiu Dong, Xun 746  
Wu, Shaohan Huang, and Furu Wei. 2025b. Reward 747  
reasoning model. *arXiv preprint arXiv:2505.14674*. 748

Mingui He, Yilun Liu, Shimin Tao, Yuanchang Luo, 749  
Hongyong Zeng, Chang Su, Li Zhang, Hongxia 750  
Ma, Daimeng Wei, Weibin Meng, and 1 others. 751  
2025. R1-t1: Fully incentivizing translation capa- 752  
bility in llms via reasoning learning. *arXiv preprint* 753  
*arXiv:2502.19735*. 754

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, 755  
Ondřej Bojar, Anton Dvorkovich, Christian Fed- 756  
ermann, Mark Fishel, Markus Freitag, Thammie 757  
Gowda, Roman Grundkiewicz, Barry Haddow, 758  
Philipp Koehn, Benjamin Marie, Christof Monz, 759  
Makoto Morishita, Kenton Murray, Masaaki Nagata, 760  
Toshiaki Nakazawa, Martin Popel, and 3 others. 2023. 761  
Findings of the 2023 conference on machine transla- 762  
tion (WMT23): LLMs are here but not quite there yet. 763  
In *Proceedings of the Eighth Conference on Machine* 764  
*Translation*, pages 1–42, Singapore. Association for 765  
Computational Linguistics. 766

Tom Kocmi and Christian Federmann. 2023a. GEMBA- 767  
MQM: Detecting translation quality error spans with 768  
GPT-4. In *Proceedings of the Eighth Conference* 769

770	<i>on Machine Translation</i> , pages 768–775, Singapore. Association for Computational Linguistics.	
771		
772	Tom Kocmi and Christian Federmann. 2023b. <a href="#">Large language models are state-of-the-art evaluators of translation quality</a> . In <i>Proceedings of the 24th Annual Conference of the European Association for Machine Translation</i> , pages 193–203, Tampere, Finland. European Association for Machine Translation.	
773		
774		
775		
776		
777		
778	Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, and 1 others. 2024. Tulu 3: Pushing frontiers in open language model post-training. <i>arXiv preprint arXiv:2411.15124</i> .	
779		
780		
781		
782		
783		
784	Junlong Li, Shichao Sun, Weizhe Yuan, Run-Ze Fan, Hai Zhao, and Pengfei Liu. 2024. Generative judge for evaluating alignment. In <i>ICLR</i> .	
785		
786		
787	Sinuo Liu, Chenyang Lyu, Minghao Wu, Longyue Wang, Weihua Luo, Kaifu Zhang, and Zifu Shang. 2025a. New trends for modern machine translation with large reasoning models. <i>arXiv preprint arXiv:2503.10351</i> .	
788		
789		
790		
791		
792	Yantao Liu, Zijun Yao, Rui Min, Yixin Cao, Lei Hou, and Juanzi Li. 2025b. Pairjudge rm: Perform best-of-n sampling with knockout tournament. <i>arXiv preprint arXiv:2501.13007</i> .	
793		
794		
795		
796	Dakota Mahan, Duy Van Phung, Rafael Rafailov, Chase Blagden, Nathan Lile, Louis Castricato, Jan-Philipp Fränken, Chelsea Finn, and Alon Albalak. 2024. Generative reward models. <i>arXiv preprint arXiv:2410.12832</i> .	
797		
798		
799		
800		
801	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. <a href="#">Scaling up CometKiwi: Unbabel-IST 2023 submission for the quality estimation shared task</a> . In <i>Proceedings of the Eighth Conference on Machine Translation</i> , pages 841–848, Singapore. Association for Computational Linguistics.	
802		
803		
804		
805		
806		
807		
808		
809	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 (EMNLP)</i> , pages 2685–2702, Online. Association for Computational Linguistics.	
810		
811		
812		
813		
814		
815	Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. <a href="#">BLEURT: Learning robust metrics for text generation</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7881–7892, Online. Association for Computational Linguistics.	
816		
817		
818		
819		
820		
821	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, and 1 others. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. <i>arXiv preprint arXiv:2402.03300</i> .	
822		
823		
824		
825		
826		
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, and 1 others. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	827
		828
		829
		830
		831
		832
	Jiaan Wang, Fandong Meng, Yunlong Liang, and Jie Zhou. 2025. <a href="#">DRT: Deep reasoning translation via long chain-of-thought</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 6770–6782, Vienna, Austria. Association for Computational Linguistics.	833
		834
		835
		836
		837
		838
	Yihong Wu, Liheng Ma, Lei Ding, Muzhi Li, Xinyu Wang, Kejia Chen, Zhan Su, Zhanguang Zhang, Chenyang Huang, Yingxue Zhang, and 1 others. 2025. It takes two: Your grpo is secretly dpo. <i>arXiv preprint arXiv:2510.00977</i> .	839
		840
		841
		842
		843
	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025a. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	844
		845
		846
		847
		848
	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, and 1 others. 2024. Qwen2.5 technical report. <i>arXiv preprint arXiv:2412.15115</i> .	849
		850
		851
		852
	Sen Yang, Yu Bao, Yu Lu, Jiajun Chen, Shujian Huang, and Shanbo Cheng. 2025b. <a href="#">EnAnchored-X2X: English-anchored optimization for many-to-many translation</a> . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 21304–21317, Suzhou, China. Association for Computational Linguistics.	853
		854
		855
		856
		857
		858
		859
	Wenjie Yang, Mao Zheng, Mingyang Song, Zheng Li, and Sitong Wang. 2025c. Ssr-zero: Simple self-rewarding reinforcement learning for machine translation. <i>arXiv preprint arXiv:2505.16637</i> .	860
		861
		862
		863
	Ziyi Ye, Xiangsheng Li, Qiuchi Li, Qingyao Ai, Yujia Zhou, Wei Shen, Dong Yan, and Yiqun Liu. 2025. Learning llm-as-a-judge for preference alignment. In <i>The Thirteenth International Conference on Learning Representations</i> .	864
		865
		866
		867
		868
	Armel Zebaze, Rachel Bawden, and Benoît Sagot. 2025. Llm reasoning for machine translation: Synthetic data generation over thinking tokens. <i>arXiv preprint arXiv:2510.11919</i> .	869
		870
		871
		872
	Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agarwal. 2025. Generative verifiers: Reward modeling as next-token prediction. In <i>The Thirteenth International Conference on Learning Representations</i> .	873
		874
		875
		876
		877
	Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong Liu, Rui Men, An Yang, and 1 others. 2025. Group sequence policy optimization. <i>arXiv preprint arXiv:2507.18071</i> .	878
		879
		880
		881
		882

Lianghui Zhu, Xinggang Wang, and Xinlong Wang. 2025. Judgelm: Fine-tuned large language models are scalable judges. In *The Thirteenth International Conference on Learning Representations*.

## A Auxiliary Score Consistency Reward

To encourage the model to generate calibrated numerical scores that reflect the magnitude of quality differences, we augment the standard ranking accuracy reward with a score consistency objective. The GRRM is instructed to output a structured response concluding with two specific components: an explicit ranking string (e.g.,  $A > B = C$ ) and a dictionary of scalar scores (e.g.,  $\{A : 6, B : 5, C : 5\}$ ).

**Internal Consistency Constraint.** Before evaluating accuracy, we enforce a strict self-consistency check to ground the numerical outputs. Let  $\hat{\tau}_{\text{text}}$  be the explicit ranking string generated by the model, and  $\hat{\tau}_{\text{score}}$  be the ranking induced by sorting the generated scalar scores  $\hat{s}$ . We define a binary consistency gate  $C_{\text{gate}}$ :

$$C_{\text{gate}} = \mathbb{I}(\hat{\tau}_{\text{text}} \equiv \hat{\tau}_{\text{score}}) \quad (4)$$

This gate is applied in the reward calculation to penalize inconsistency.

**Margin-Aware Score Reward.** For samples passing the consistency gate, we calculate the Score Consistency Reward ( $R_{\text{score}}$ ). Unlike simple regression losses (e.g., MSE), which can be sensitive to the absolute scale of ground-truth labels, our metric focuses on preserving the *relative margins* between candidates. Let  $q_i, q_j$  be the ground-truth scores and  $\hat{s}_i, \hat{s}_j$  be the predicted scores for candidates  $i$  and  $j$ . We compute the absolute margin error  $\delta_{ij}$  for every pair:

$$\delta_{ij} = |(\hat{s}_i - \hat{s}_j) - (q_i - q_j)| \quad (5)$$

We then map this error to a reward value using a discrete kernel function  $K(\cdot)$ , designed to penalize deviations while allowing for minor integer fluctuations. Based on our empirical tuning, we define:

$$K(\delta) = \begin{cases} 1.0 & \text{if } \delta = 0 \\ 0.6 & \text{if } \delta = 1 \\ 0.2 & \text{if } \delta = 2 \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

The final score consistency reward is the average kernel value across all pairs:

$$R_{\text{score}} = \frac{1}{\binom{G}{2}} \sum_{1 \leq i < j \leq G} K(\delta_{ij}) \quad (7)$$

**Total Reward.** The final reward used for RLVR optimization is the sum of the ranking accuracy (Eq. 3) and the score consistency reward, gated by the internal consistency check:

$$R_{\text{total}} = C_{\text{gate}} \cdot (R_{\text{acc}} + R_{\text{score}}) \quad (8)$$

This formulation ensures that the model is optimized to produce reasoning that results in both correct ordering and precise, margin-aware numerical estimations.

**SQM-GenRM Configuration.** For the SQM-GenRM baseline, which evaluates candidates independently, we adapt the reward to target absolute accuracy. Unlike the margin-based approach in GRRM, we compute the absolute error  $\delta = |\hat{s} - q|$  between the predicted score  $\hat{s}$  and the ground truth  $q$ . The final reward utilizes a similar discrete kernel function  $K(\delta)$  defined above to penalize deviations, ensuring a fair comparison between the optimization objectives.

## B Implementation Details

In this section, we detail the training configurations for both the Reward Model (GRRM) and the Machine Translation Policy Optimization.

### B.1 Reward Model Implementation

**Training Data Construction.** We construct the training dataset using the Chinese-English subset of TowerBlocks (MT), by sampling group-level candidates using a seed translation model. Specifically, we employ the Qwen2.5-7B model fine-tuned on the TowerBlocks (MT) dataset (approximately 150k samples). This seed model is also used as SFT baseline (w/o Reasoning) in Section 4 and Appendix C.2. We apply standard sampling with a temperature of  $T = 1.0$ . For each source sentence, we sample a group of  $N$  candidates, where  $N \in \{2, 3, 4\}$  with a frequency ratio of 1:1:3. To enhance candidate diversity and ensure the presence of high-quality samples within the generated groups, we inject the ground-truth reference into the candidate pool with a probability of 0.5. These constructed groups are subsequently annotated by Gemini-2.5-Pro. Finally, we apply data augmentation by randomly shuffling candidate orders and sampling subgroups from the original annotated groups, ensuring the model learns ranking criteria that are invariant to input position and specific group contexts. This has been proved to be effective in Zhu et al. (2025).

974	<b>Evaluation Data Processing.</b>	For the human-	switch to a constant learning rate scheduler to main-	1024
975		annotated test sets (Newstest2020 and Gen-	tain stable updates throughout the training process.	1025
976		eralMT2022), we calculate the final quality score	The number of rollouts per prompt is set to $G = 4$ ,	1026
977		for each translation candidate by averaging the rat-	aligning with the group size capacity of our reward	1027
978		ings across raters. From the available system out-	model.	1028
979		puts for each source, we retain a subset of 2 to	To ensure fair comparison across different data	1029
980		4 candidates, following the same distribution ra-	settings, we adjust the number of training epochs	1030
981		tio as the training data. However, we observe that	to keep the total number of optimization steps con-	1031
982		randomly selected candidates often possess similar	sistent. Specifically, we train for 2 epochs on the	1032
983		human quality scores (clustering near the ceiling),	standard multilingual TowerBlocks (MT) dataset	1033
984		which introduces ambiguity into the ranking evalua-	(150k samples). For the CLA variant, we train for	1034
985		tion. Therefore, to ensure a rigorous assessment of	1 epoch, while for the Chinese-English only subset,	1035
986		the model’s ability to distinguish quality varian-	we extend training to 8 epochs.	1036
987		ces, we enforce a filtering constraint: the construc-		
988		ted group must include the candidates with the mini-	<b>C Extended Analysis on Machine</b>	1037
989		mum and maximum human scores for that entry.	<b>Translation</b>	1038
990	<b>Training Hyperparameters.</b>	We employ 16	<b>C.1 Full Multilingual Benchmarks</b>	1039
991		Nvidia A100 (80GB) GPUs for all training experi-	We provide the detailed breakdown of the aggre-	1040
992		ments detailed in this paper, including both reward	gated performance reported in the main results. Ta-	1041
993		modeling and the subsequent machine translation	ble 4 details the En→X performance across seven	1042
994		policy optimization.	European and Slavic languages from WMT24++,	1043
995		In the SFT stage, we optimize the model for 3	while Table 5 presents the WMT23 X→En results	1044
996		epochs with a global batch size of 64. We employ a	for German, Japanese, and Russian.	1045
997		cosine learning rate scheduler with a peak learning	The results demonstrate that our GRPO-based	1046
998		rate of $6 \times 10^{-6}$ and a warmup ratio of 0.1.	optimization yields consistent improvements over	1047
999		In the RLVR stage, we adopt Group Se-	the SFT baseline across diverse language families.	1048
1000		quence Policy Optimization (GSPO) (Zheng et al.,	For WMT23 X→En benchmark, while SeedX-PPO	1049
1001		2025), an enhanced version of GRPO that utilizes	generally leads in BLEURT, our method achieves	1050
1002		sequence-level importance ratios to improve train-	superior or competitive LLM judge scores, particu-	1051
1003		ing stability and efficiency. We set the learning	larly in directions like Ja→En and Ru→En.	1052
1004		rate to $1 \times 10^{-5}$ with a cosine scheduler, decay-	<b>C.2 Extended Ablation Studies</b>	1053
1005		ing to a minimum ratio of 0.2. The training is con-	In this section, we provide supplementary experi-	1054
1006		ducted with a total batch size of 512 and a PPO	ments to further validate the robustness and univer-	1055
1007		mini-batch size of 128. For each prompt, we gen-	sality of our proposed framework. Table 6 presents	1056
1008		erate $G = 8$ rollouts. The model is trained for a	the detailed results.	1057
1009		single epoch. Following recent practices in reason-	<b>Reward Model Stability on Non-Reasoning</b>	1058
1010		ing model training, we disable the KL divergence	<b>Baselines.</b>	1059
1011		penalty to encourage broader exploration of reason-	We first investigate whether the in-	1060
1012		ing paths.	stability of certain reward models (observed in Sec-	1061
1013	<b>B.2 MT Optimization Implementation</b>		tion 4.3) persists when optimizing a standard, non-	1062
1014	<b>SFT Hyperparameters.</b>	For the policy model	reasoning translation policy. As shown in the first	1063
1015		cold-start, we perform Supervised Fine-Tuning for	block of Table 6, the <b>SQM-GenRM</b> continues to	1064
1016		1 epoch with a global batch size of 64. We utilize	exhibit severe failure modes, particularly on En→X	1065
1017		a cosine learning rate scheduler, setting the peak	tasks (dropping to a score of 7.13), confirming that	1066
1018		learning rate to $1 \times 10^{-5}$ with a warmup ratio of	the vulnerability to reward hacking is inherent to	1067
1019		0.1.	the reward model itself rather than specific to rea-	1068
1020	<b>GRPO Hyperparameters.</b>	We utilize GSPO for	soning policies. In contrast, our <b>GRRM</b> consis-	1069
1021		policy optimization. The configuration largely mir-	tently outperforms other reward signals (BLEURT,	1070
1022		rors the GRRM RLVR stage (Appendix B.1), with	BT-RM) across most benchmarks. However, it	1071
1023		specific adjustments to suit the translation task. We	is worth noting that even with the best reward	1072



Figure 2: Comparison of Scalar Quality Metric (SQM) and Group Quality Metric (GQM) across four distinct scenarios. **Case 1 & 2 (Top):** Demonstrate GQM’s ability to resolve fine-grained stylistic nuances and identify localization errors that SQM misses due to independent evaluation. **Case 3 & 4 (Bottom):** Illustrate how GQM uses contrastive context to detect hallucinations and semantic omissions that appear fluent in isolation. Case 1,3,4 were conducted using Gemini-2.5-Pro and Case 2 was conducted using DeepSeek-R1-0528.

79.35 on Seed-X En→Zh) still significantly lags behind the reasoning-enhanced policy reported in the main text (Score 84.36). This further corroborates our conclusion that the "Reasoning × Reasoning" paradigm is essential for peak performance.

### Generalization to Instruction-Tuned Models.

To assess the universality of our approach, we apply our GRPO training with GRRM directly to a standard instruction-tuned model, Qwen2.5-7B-Instruct, without any task-specific SFT warm-up. As shown in the second block of Table 6, our method yields substantial improvements over the strong base model. Notably, on the unseen En→X task, the performance improves by over 11

points on the oss-judge metric (70.84 → 82.44). This result demonstrates that our framework is not limited to specialized translation models but can serve as a general-purpose alignment technique to enhance the multilingual capabilities of off-the-shelf LLMs.

### C.3 Case Studies: Emergence of Reasoning

To provide a qualitative perspective on the quantitative gains observed in the Seed-X-Challenge, Figure 3 illustrates the emergence of reasoning capabilities in our GRRM-optimized model. Unlike standard SFT models that often rely on direct surface-level mapping, our model demonstrates a

Model	WMT24++ En→X Detailed Results													
	En→De		En→Es		En→Fr		En→It		En→Nl		En→Pt		En→Ru	
	BLEURT	R1-judge	BLEURT	R1-judge	BLEURT	R1-judge	BLEURT	R1-judge	BLEURT	R1-judge	BLEURT	R1-judge	BLEURT	R1-judge
<i>General LLMs</i>														
Gemini-2.5-Pro	<b>71.75</b>	<b>89.79</b>	<b>71.82</b>	<b>91.10</b>	<b>64.63</b>	<b>91.03</b>	<b>71.33</b>	<b>91.14</b>	<b>72.21</b>	<b>89.54</b>	<b>63.41</b>	<b>90.00</b>	<b>66.91</b>	<b>89.88</b>
DeepSeek-R1-0528	69.61	88.18	71.17	90.68	63.53	89.44	70.62	89.57	71.26	87.10	62.91	87.61	64.90	86.76
Qwen3-8B	63.90	77.13	64.58	82.71	40.45	76.21	62.53	78.98	62.89	71.37	56.80	80.58	49.62	74.62
Qwen2.5-7B-Instruct	59.43	68.94	62.84	78.61	54.47	76.80	59.84	71.83	61.44	64.21	56.64	76.82	56.41	70.39
<i>Translation-Specialized Models</i>														
TowerInstruct-13B	69.11	82.34	68.68	84.26	61.51	83.33	70.35	84.40	70.49	82.09	63.37	83.30	62.73	79.02
SeedX-PPO	<b>71.15</b>	<b>86.03</b>	<b>70.63</b>	<b>87.51</b>	<b>62.76</b>	<b>86.31</b>	<b>71.19</b>	<b>86.69</b>	<b>71.69</b>	<b>84.20</b>	<b>65.25</b>	85.95	<b>65.76</b>	<b>85.62</b>
Qwen2.5-7B-SFT	59.06	65.68	61.91	74.25	51.33	71.41	60.12	66.59	58.10	60.05	55.77	72.98	53.72	64.44
+ GRPO (ours)	68.12	82.63	66.59	86.91	59.78	85.27	67.43	84.07	67.04	79.30	61.24	85.65	62.32	83.22
+ GRPO w/ CLA (ours)	67.92	82.60	66.44	86.76	60.18	85.84	66.98	83.67	66.50	78.62	61.12	<b>86.43</b>	62.39	82.07

Table 4: **Detailed breakdown of WMT24++ En→X results.** We report BLEURT-20 and LLM-as-a-Judge scores evaluated by DeepSeek-R1-0528.

Model	WMT23 X→En Detailed Results							
	De→En		Ja→En		Ru→En		Average	
	BLEURT	oss-judge	BLEURT	oss-judge	BLEURT	oss-judge	BLEURT	oss-judge
<i>General LLMs</i>								
Qwen3-8B	69.12	<b>88.65</b>	60.87	<b>80.11</b>	65.42	<b>84.66</b>	65.14	<b>84.47</b>
Qwen2.5-7B-Instruct	<b>71.75</b>	86.21	<b>66.23</b>	79.03	<b>71.04</b>	83.80	<b>69.67</b>	83.01
<i>Translation-Specialized Models</i>								
TowerInstruct-13B	74.48	90.04	66.74	75.66	72.98	86.46	71.40	84.06
SeedX-PPO	<b>75.43</b>	<b>92.05</b>	<b>70.19</b>	84.99	<b>74.40</b>	89.82	<b>73.34</b>	<b>88.95</b>
Qwen2.5-7B-SFT	71.51	85.54	66.38	78.08	71.15	84.15	69.68	82.59
+ GRPO (ours)	70.25	90.69	67.72	<b>85.50</b>	71.93	<b>90.14</b>	69.97	88.78
+ GRPO w/ CLA (ours)	70.83	90.47	67.73	84.98	71.83	89.25	70.13	88.23

Table 5: **Detailed breakdown of WMT23 X→En results.** We report BLEURT-20 and LLM-as-a-Judge scores evaluated by gpt-oss-120b.

"think-before-translating" mechanism that is crucial for resolving cultural ambiguities.

In the scenario of **Translating Idiomatic Expressions** scenario (Case 1), the model successfully navigates the English proverb "The grass is always greener on the other side." Instead of producing a rigid or literal translation (e.g., regarding the color of grass), the model identifies the underlying semantic meaning of dissatisfaction and envy. It then retrieves the culturally equivalent Chinese proverb "这山望着那山高" (*This mountain looks higher than that one*), ensuring the translation resonates with native speakers.

Similarly, in the **Decoding Internet Slang** (Case 2), the model encounters the phrase "INTJ总是装E", which blends technical MBTI terminology with colloquial Chinese slang. A literal breakdown might fail to capture the nuance of "装" (feign/pretend) in this specific context. However, the model's reasoning trace explicitly decomposes the acronyms and analyzes the character traits, correctly deriving the idiomatic English translation: "putting on an extroverted front."

These cases confirm that our approach does not merely memorize translation pairs but actively performs semantic analysis and cultural alignment, justifying the substantial improvements seen in the LLM-judge evaluations.

## D GRRM Application: Inference-time Reranking

In this section, we investigate the effectiveness of GRRM as a verifier to select the best translation from multiple candidates during inference.

**Experimental Setup.** We conduct experiments on both reasoning and non-reasoning models, evaluating their respective SFT baselines and GRPO-optimized checkpoints. For candidate generation, we employ standard sampling with a temperature of 0.6 to produce  $N = 4$  candidates per input. We compare our GRRM reranking strategy against: (1) **Sampling Baseline**, representing the expected performance of the policy; and (2) **Best-of- $N$** , which selects the candidate with the highest log-probability assigned by the policy model. Addition-

Supplementary Configuration	WMT Benchmarks						Seed-X-Challenge			
	Zh→En		En→Zh		En→X		Zh→En		En→Zh	
	BLEURT	oss-judge	BLEURT	oss-judge	BLEURT	oss-judge	BLEURT	oss-judge	BLEURT	oss-judge
<b>Comparison of Reward Models for Non-Reasoning Baselines</b>										
SFT baseline (w/o Reasoning)	67.61	85.28	63.50	79.42	63.24	75.33	65.52	76.83	64.04	73.34
+ GRPO w/ BLEURT	68.36	86.46	65.28	81.19	65.12	77.49	67.21	79.17	65.25	75.99
+ GRPO w/ BT-RM	68.26	88.35	64.67	85.44	63.33	78.21	69.41	82.88	66.00	78.85
+ GRPO w/ SQM-GenRM	68.02	87.55	38.88	36.08	27.14	7.13	68.32	82.10	42.49	39.85
+ GRPO w/ GRRM	67.38	88.42	64.15	85.19	63.40	78.88	68.57	85.38	65.92	79.35
<b>GRRM with Instruct-Tuned Base Model</b>										
Qwen2.5-7B-Instruct	67.31	85.84	59.92	79.42	58.72	70.84	66.59	79.95	62.75	73.93
+ GRPO w/ GRRM	67.15	89.40	64.88	87.07	63.89	82.44	68.74	86.99	67.37	83.45

Table 6: **Additional ablation results.** We report performance for (1) Applying different reward models to a non reasoning model tuned from Qwen2.5-7B, and (2) Applying GRRM optimization on top of the instruction-tuned Qwen2.5-7B-Instruct model.

Base Model	Inference Strategy	WMT23		WMT24++		Seed-X-Challenge				Average	
		Zh→En		En→Zh		Zh→En		En→Zh			
		BLEURT	oss-judge	BLEURT	oss-judge	BLEURT	oss-judge	BLEURT	oss-judge	BLEURT	oss-judge
SFT Baseline (w/ Reasoning)	Sampling Baseline	66.78	85.04	59.33	74.20	66.95	80.17	61.64	70.02	63.67	77.36
	Best-of-N	66.76	85.21	59.79	75.50	66.88	79.57	61.94	71.27	63.84	77.89
	<b>Ranking w/ GRRM</b>	<b>67.68</b>	<b>88.06</b>	<b>61.50</b>	<b>81.47</b>	<b>68.05</b>	<b>83.22</b>	<b>63.29</b>	<b>77.15</b>	<b>65.13</b>	<b>82.47</b>
+ GRPO w/ GRRM	Sampling Baseline	67.40	90.15	64.77	88.05	69.23	87.41	67.20	84.79	67.15	87.60
	Best-of-N	<b>67.55</b>	90.28	64.99	88.26	<b>69.42</b>	87.73	67.21	83.89	67.29	87.54
	<b>Ranking w/ GRRM</b>	67.45	<b>90.41</b>	<b>65.17</b>	<b>88.90</b>	69.16	<b>87.77</b>	<b>67.54</b>	<b>85.90</b>	<b>67.33</b>	<b>88.25</b>
SFT Baseline (w/o Reasoning)	Sampling Baseline	66.05	82.63	61.50	74.63	64.16	73.75	63.04	69.29	63.69	75.07
	Best-of-N	67.11	84.49	62.42	77.48	64.88	73.85	62.91	71.00	64.33	76.71
	Beam Search	<b>68.18</b>	86.28	<b>63.62</b>	<b>81.34</b>	<b>66.25</b>	77.68	<b>64.16</b>	75.45	<b>65.55</b>	80.19
	<b>Ranking w/ GRRM</b>	67.54	<b>86.65</b>	63.18	81.06	65.75	<b>78.35</b>	63.99	<b>75.46</b>	65.12	<b>80.38</b>
+ GRPO w/ GRRM	Sampling Baseline	67.17	88.10	63.91	84.93	68.19	84.97	65.96	79.44	66.31	84.36
	Best-of-N	67.38	88.47	64.12	85.11	68.47	85.24	65.66	78.00	66.40	84.21
	Beam Search	<b>67.54</b>	88.81	64.20	85.45	68.48	85.62	65.98	78.64	66.55	84.63
	<b>Ranking w/ GRRM</b>	67.37	<b>89.21</b>	<b>64.32</b>	<b>86.27</b>	<b>68.53</b>	<b>86.59</b>	<b>66.60</b>	<b>80.03</b>	<b>66.70</b>	<b>85.53</b>

Table 7: **Inference-time Reranking Performance.** We compare different decoding strategies across reasoning and non-reasoning models.

ally, for non-reasoning models, we include **Beam Search** with a beam width of 4 as a strong decoding baseline. Notably, we exclude Beam Search for reasoning models, as the extensive length of Chain-of-Thought sequences renders the decoding process prohibitively slow and computationally intractable.

**Results and Analysis.** The results presented in Table 7 show that ranking with GRRM consistently outperforms the Sampling and Best-of-N baselines across all datasets and model stages. Even for GRPO-optimized models, which have already been aligned via RL, inference-time reranking provides further performance gains. For reasoning models where Beam Search is inapplicable, GRRM reranking proves to be a highly effective alternative. It significantly boosts the SFT baseline (e.g., an average increase of over 5 points in LLM-judge scores), effectively filtering out flawed reasoning paths.

In non-reasoning settings, while Beam Search remains competitive on BLEURT, GRRM reranking consistently achieves superior LLM-judge scores. This suggests that GRRM prioritizes semantic fidelity and better aligns with human preferences compared to likelihood-based decoding strategies.

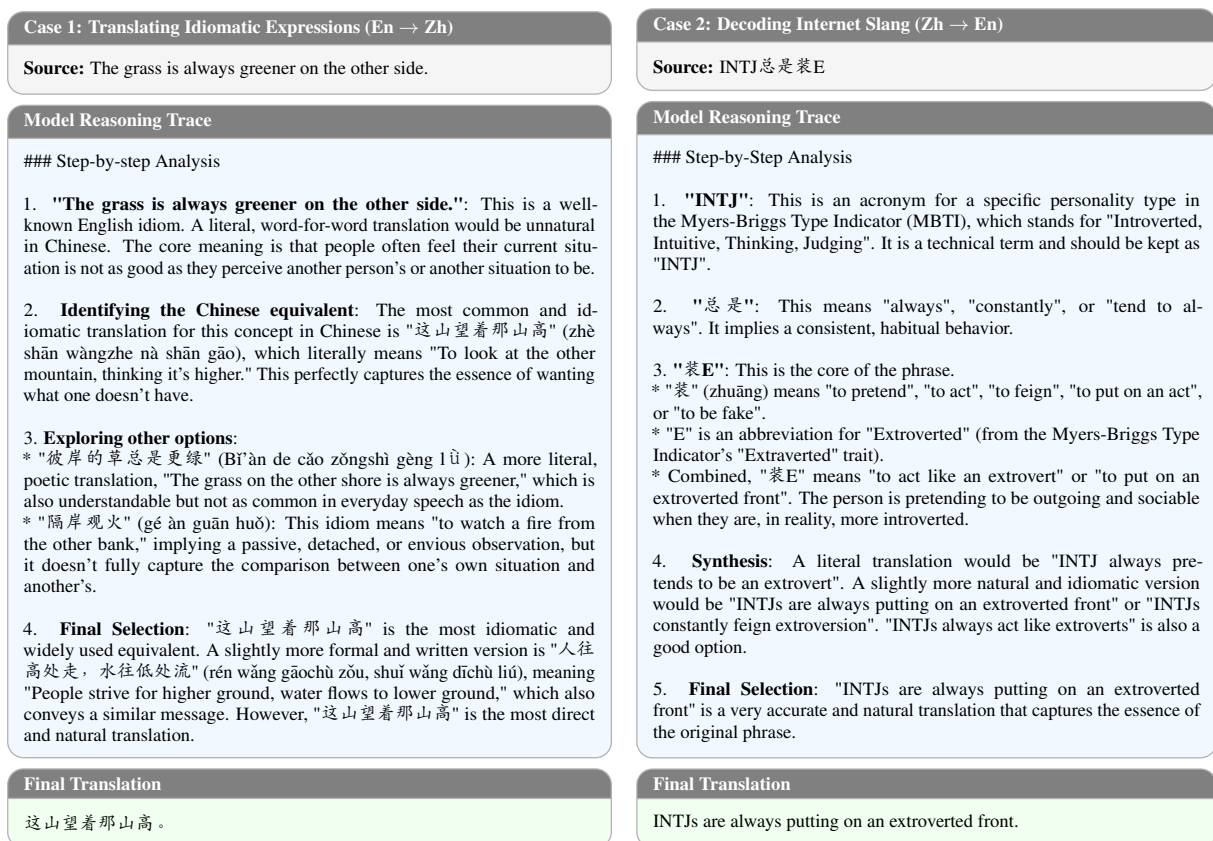


Figure 3: Qualitative analysis of translation reasoning. We showcase the model's ability to handle cultural nuances by identifying **idiomatic equivalents** (Case 1) and **decoding internet slang** (Case 2) through step-by-step analysis.