

Learning to Translate by Translating: Stabilizing the Dual Loop via Semantic-Aware Self-Evolution

Anonymous ACL submission

Abstract

Despite the remarkable success of Large Language Models (LLMs) in Machine Translation (MT), the scarcity of high-quality parallel corpora and the prohibitive cost of their acquisition constrain scalability. To this end, we propose **Learning to Translate by Translating (LTT)**, an LLM-driven dual-learning framework that enables autonomous translation, achieving an 80.42% performance improvement over the base model. By adapting the cycle-consistency principle to the generative paradigm, LTT eliminates the need for parallel data. It employs a robust semantic-aware reward function that balances adequacy with reconstruction fidelity, effectively mitigating the reward hacking issues inherent in traditional unsupervised MT. Relying solely on monolingual data, our 8B model consistently outperforms significantly larger models (70B+) in low-resource settings and achieves parity with state-of-the-art supervised baselines on mainstream benchmarks. LTT thus offers a scalable, data-efficient paradigm for autonomous machine translation.

1 Introduction

LLM-based translation systems, such as Tower (Alves et al.) and X-ALMA (Xu et al., b), have achieved state-of-the-art (SOTA). However, this success is largely attributed to the effectiveness of supervised training on vast, human-curated parallel corpora. The high cost of corpus construction poses a major barrier to further improvement. To break the dependency on parallel corpora, leveraging monolingual data has long been a long-standing goal in MT research. Pioneering works in Dual Learning (He et al., 2016) and Unsupervised Machine Translation (UMT) (Lample et al., 2018) introduced the concept of utilizing round-trip consistency as a training signal. Despite their theoretical appeal, these early attempts struggled with instability and

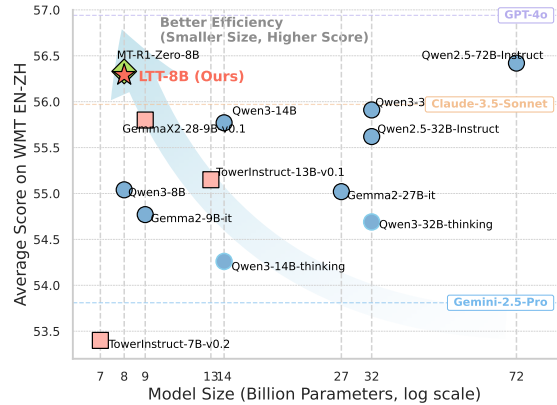


Figure 1: Overall performance of LTT-8B on the WMT EN-ZH benchmark. Our enhanced 8B model (red star) rivals the performance of models 4x-9x larger (e.g., Qwen2.5-72B) and GPT-4o, while outperforming all competitors within the <10B parameter class.

were prone to "reward hacking"—where models optimize for lexical reconstruction at the expense of semantic fidelity (Marchisio et al., 2020). Consequently, achieving robust, fully autonomous translation without references remains an open challenge.

The advent of LLMs with reasoning and self-correction capabilities, such as OpenAI o1 (Jaech et al., 2024) and DeepSeek R1 (Guo et al., 2025), offers a unique opportunity to revisit and modernize this traditional framework. However, a naive integration is insufficient. As revealed in our analysis (Section 5.4), directly applying traditional reconstruction objectives to LLMs leads to catastrophic shortcuts, such as copying source text to maximize overlap scores. Furthermore, recent RL-based MT approaches, such as MT-R1-Zero (Feng et al., 2025), still rely heavily on ground-truth references for reward computation, limiting their scalability in low-resource scenarios.

To this end, we introduce LTT, an enhanced dual learning framework tailored for the LLM era. Unlike prior methods constrained by unstable op-

066 timization or simplistic metrics, LTT leverages
067 Group Relative Policy Optimization (GRPO) for
068 stable training and integrates a semantic-aware re-
069 ward mechanism to prevent reward hacking. This
070 allows the model to self-evolve using solely mono-
071 lingual data, achieving performance comparable
072 to, or even exceeding, supervised models.

073 Extensive experimental results demonstrate that
074 LTT achieves parity with state-of-the-art baselines
075 without relying on parallel corpora fine-tuning. On
076 EN-ZH, our 8B model nearly matches optimal su-
077 pervised benchmarks (lagging by only 0.07%) and
078 surpasses proprietary systems like Gemini-2.5-Pro
079 (77.58 vs. 77.55) in semantic evaluation. Notably,
080 it outperforms the significantly larger Qwen2.5-
081 72B-Instruct (76.52), highlighting its efficiency
082 (see Figure 1). In challenging multilingual set-
083 tings, LTT boosts the base model from 32.22 to
084 58.51, outperforming competitors like Gemma2-
085 9B-it and even LLaMA-3.1-70B. Further analysis
086 attributes these gains to our hybrid reward func-
087 tion, which synergizes lexical and semantic sig-
088 nals to mitigate reward hacking and ensure stable
089 convergence.

090 Our contributions are summarized as follows:

- 091 • We propose LTT, an LLM-driven dual-learning
092 framework that integrates GRPO into a
093 self-evaluation loop, enabling effective self-
094 evolution using solely monolingual data.
- 095 • We design a semantic-aware reward architec-
096 ture that balances reconstruction accuracy with
097 anti-cheating constraints, effectively mitigating
098 the reward hacking problem inherent in unsuper-
099 vised objectives.
- 100 • Extensive experiments show our 8B model
101 matches supervised baselines and significantly
102 outperforms much larger LLMs (e.g., 70B+) in
103 low-resource settings.

104 2 Related Work

105 **Machine Translation with Large Language**
106 **Models.** Cutting-edge machine translation with
107 LLMs follows two main paradigms: in-context
108 learning (ICL) and supervised fine-tuning (Anil
109 et al., 2023; Gao et al., 2024; Li et al., 2024;
110 Xu et al., a). By presenting few-shot demon-
111 strations to LLMs, ICL circumvents the heavy com-
112 putational costs of fine-tuning (Zhu et al., 2024),
113 though at the expense of prompt sensitivity and
114 performance instability (Agrawal et al., 2023).
115 Fine-tuning, in contrast, improves the MT perfor-

116 mance via supervised training on large-scale par-
117 allel datasets (Cui et al., 2025; Costa-Jussà et al.,
118 2022), exemplified by representative models such
119 as Tower and X-ALMA (Alves et al.; Guo et al.,
120 2024; Xu et al., b). However, the rules of scaling
121 laws constrain further performance improvements
122 of fine-tuning. Moreover, both two paradigms ex-
123 hibit a significant reliance on large-scale, human-
124 annotated data, which limits scalability, particu-
125 larly in low-resource scenarios.

Reasoning and Reinforcement Learning for
Machine Translation. To transcend the limi-
126 tations of supervised training, researchers have
127 adopted RL, initially to mitigate exposure bias
128 by optimizing global metrics like BLEU (Ranzato
129 et al., 2016; Wu et al., 2017). Inspired by LLM
130 reasoning, recent approaches incorporate CoT
131 prompting to decompose translation into interme-
132 diate steps for enhanced fidelity (Feng et al., 2024;
133 Wang et al., 2024). However, their success heavily
134 relies on validating the reasoning process, often ne-
135 cessitating manually engineered templates, com-
136 plex search algorithms like MCTS (Zhao et al.,
137 2024), or external evaluators (He et al., 2025). Cru-
138 cially, even leading RL frameworks such as MT-
139 R1-Zero (Feng et al., 2025) and DeepTrans (Wang
140 et al., 2025) still depend on reference translations
141 for reward computation, restricting their potential
142 for fully autonomous learning.

Dual Learning and Self-Evolution. Founda-
145 tional research in dual learning (He et al., 2016)
146 introduced a closed-loop feedback system lever-
147 aging the symmetry of translation tasks, while
148 unsupervised MT (Lample et al., 2018) utilized
149 shared latent spaces and back-translation for align-
150 ment. Recently, this paradigm has evolved into
151 self-rewarding mechanisms for LLMs (Yuan et al.,
152 2024; Zou et al., 2025). Despite their promise,
153 these approaches remain fragile: early dual meth-
154 ods suffer from severe instability under domain
155 mismatch (Marchisio et al., 2020), and generative
156 self-rewarding models are prone to reward hack-
157 ing, where they prioritize high internal scores over
158 semantic fidelity (Wu et al., 2024). LTT revital-
159 izes this framework by integrating GRPO with
160 semantic-aware rewards. By replacing unstable
161 policy gradients with GRPO and enforcing round-
162 trip consistency through semantic metrics, our ap-
163 proach effectively mitigates semantic drift and re-
164 ward hacking, enabling stable, autonomous self-
165 improvement without parallel data.

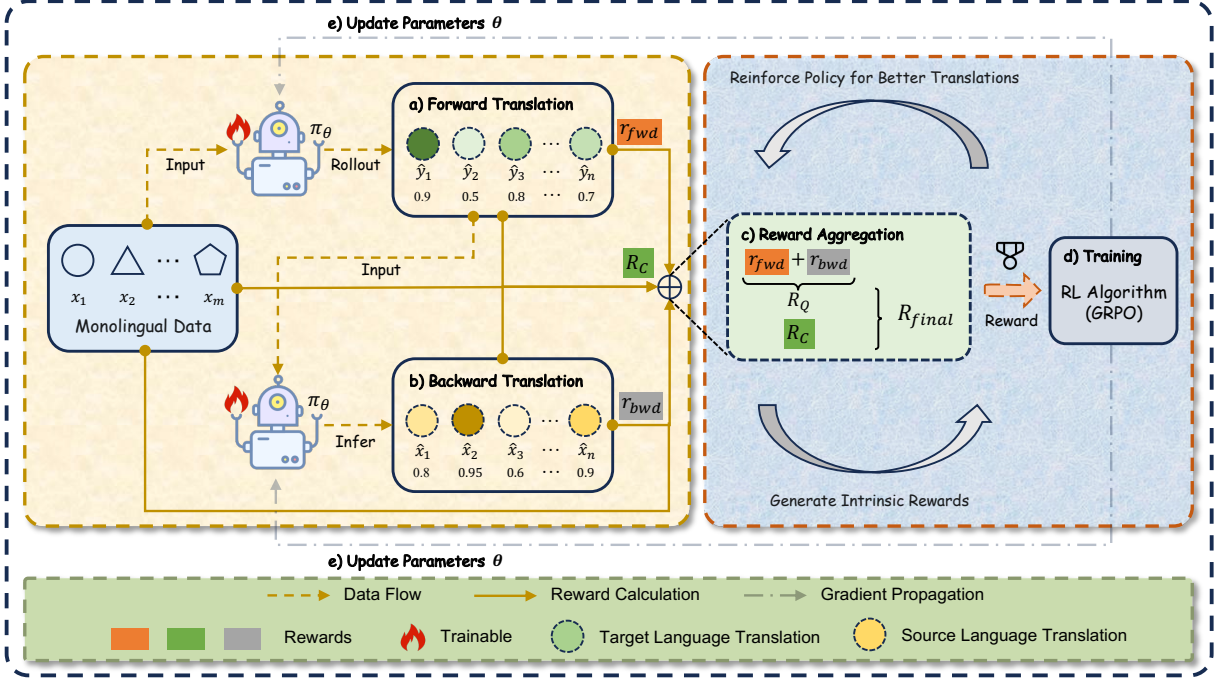


Figure 2: Overview of LTT. The actor model participates in both forward and backward translation. We consider reward signals from quality and anti-cheating, two perspectives, to update the model via GRPO.

3 Methodology

3.1 LTT Framework

Inspired by the primal-dual structure of traditional Dual Learning (He et al., 2016), our bidirectional translation loop leverages the pre-trained capabilities of LLMs and stabilizes the optimization process using GRPO with a semantic-aware reward composition. By recognizing round-trip consistency as an intrinsic reward signal, the system can self-evolve and improve within a closed-loop process (see Figure 2).

Specifically, we collect a batch of source sentences $\{x_i\}_{i=1}^m$ from language \mathcal{L}_s . The actor model, parameterized by a policy π_θ , performs two sequential steps:

1. Forward Translation ($\mathcal{L}_s \rightarrow \mathcal{L}_t$). The model generates a candidate \hat{y}_i as the target language \mathcal{L}_t translation:

$$\hat{y}_i \sim \pi_\theta(\cdot \mid x_i, \text{prompt}_{s \rightarrow t}).$$

As no reference translation is involved, this step necessitates a reference-free evaluation.

2. Backward Translation ($\mathcal{L}_t \rightarrow \mathcal{L}_s$). The process runs in reverse by translating \hat{y}_i back into the source language \mathcal{L}_s for \hat{x}_i reconstruction:

$$\hat{x}_i \sim \pi_\theta(\cdot \mid \hat{y}_i, \text{prompt}_{t \rightarrow s}).$$

In this step, the original source sentence x_i serves as a high-quality, self-reference against which we can evaluate the reconstruction quality.

Unified Prompting Strategy. To ensure consistency, we employ a unified prompting template for both translation directions, following Feng et al. (2025). The model is instructed to place its final translation within `<translate>` tags and any intermediate reasoning within `<think>` tags. The full prompt details are available in Appendix A.1.

3.2 Reward Architecture

The success of LTT hinges on a tailored reward function that guides the model towards high-quality translations while preventing reward hacking. The final reward, R_{final} , is structured hierarchically to prioritize valid formatting before assessing translation quality:

$$R_{\text{final}}(x, \hat{y}, \hat{x}) = \begin{cases} 1 + R_Q + R_C, & \text{correct format,} \\ -3, & \text{otherwise.} \end{cases}$$

A large penalty of -3 is assigned if the output \hat{y} does not adhere to the required `<translate>` tag format. For valid outputs, we evaluate both translation quality and anti-cheating capability, and additionally incorporate a base reward of $+1$.

3.2.1 Quality Reward (R_Q)

This reward integrates two complementary metrics for a holistic assessment of translation quality, defined as: $R_Q = r_{\text{fwd}} + r_{\text{bwd}}$:

- **Forward Semantic Adequacy (r_{fwd}):** For the forward pass ($x \rightarrow \hat{y}$), we use **COMETkiwi**, a widely used reference-free metric. It evaluates the semantic adequacy of the translation by comparing the source and the hypothesis straightforwardly, formalized as $r_{\text{fwd}} = \text{COMETkiwi}(x, \hat{y})$.
- **Backward Reconstruction Fidelity (r_{bwd}):** For the back-translation ($\hat{y} \rightarrow \hat{x}$), we leverage the original source x as a reference. We use **BLEU** to measure the fidelity of the reconstruction, $r_{\text{bwd}} = \text{BLEU}(\hat{x}, x)$. A high score signifies that the target translation \hat{y} preserved sufficient information to accurately recover the original input.

3.2.2 Anti-Cheating Reward (R_C)

A vanilla self-supervised reward is fragile and unstable. To ensure robust learning and mitigate reward hacking, we introduce two penalty terms, $R_C = r_{\text{copy}} + r_{\text{mix}}$:

- **Source-Copying Penalty (r_{copy}):** A trivial failure mode is to copy the source ($x \approx \hat{y}$) to maximize the backward BLEU score. We address this by penalizing lexical overlap in the forward direction: $r_{\text{copy}} = -\text{BLEU}(x, \hat{y})$.
- **Language-Mixture Penalty (r_{mix}):** To discourage the generation of linguistically incoherent outputs, we apply a penalty of -0.5 if language mixing is detected in \hat{y} . This simple heuristic effectively promotes fluent, monolingual outputs.

3.3 Policy Optimization with GRPO

We optimize our translation policy π_θ using GRPO (Shao et al., 2024; Guo et al., 2025). For each input, GRPO samples a group of G candidates from the current policy and computes a normalized advantage A_i for each based on its relative reward within the group. The policy is then updated by maximizing the standard GRPO objec-

tive:

$$J_{\text{GRPO}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, \{\hat{y}_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot|x)} \left[\frac{1}{G} \sum_{i=1}^G \min \left(\rho_i(\theta) A_i, \text{clip}(\rho_i(\theta), 1 - \varepsilon, 1 + \varepsilon) A_i \right) - \beta D_{\text{KL}}(\pi_\theta(\cdot|x) \parallel \pi_{\text{ref}}(\cdot|x)) \right], \quad (1)$$

where $\rho_i(\theta)$ is the probability ratio. This objective uses a clipped surrogate function and a KL-divergence penalty to ensure stable policy updates. Please refer to Appendix B.1 for further details.

4 Experiments

4.1 Experimental Setup

Datasets. We conduct experiments to evaluate the effectiveness of our LTT across bilingual and multilingual settings. Attribute to the merit of free-parallel-data in our methodology, we, thereby, select the popular parallel translation benchmarks and remove all target-side references for training corpora construction.

- **Training Data:** For bilingual (EN-ZH) training, we source sentences from WMT 2017-2020 competitions, yielding 6,565 sentences each for English and Chinese. For multilingual training, we create a diverse corpus from the FLORES-200 (Costa-Jussà et al., 2022) training set, covering EN/ZH paired with six other languages (DE, FR, ES, IT, JA, KO). We sample 500 pairs for each of the 24 translation combinations (e.g., EN→DE, ZH→DE), treating all sentences as monolingual, resulting in a final 12,000-sentence corpus.
- **Evaluation Data:** We evaluate performance on test sets for fair comparison. For EN↔ZH, we use the official test sets from WMT23¹ and WMT24². For multilingual tasks, we report results on the official FLORES-200 test set.

Evaluation Metrics. To ensure a comprehensive assessment of translation quality, we adopt a dual-metric approach that captures both lexical fidelity and semantic adequacy. We report case-sensitive **BLEU** scores computed via sacrebleu for standardized, reproducible measurement of n-gram overlap. To evaluate semantic preservation, we employ **xCOMET-XL** (Guerreiro et al., 2024),

¹<https://www2.statmt.org/wmt23/translation-task.html>

²<https://www2.statmt.org/wmt24/translation-task.html>

Table 1: Main results on WMT and FLORES-200 benchmarks. The best and second-best scores are **bolded** and underlined, respectively. "†" indicates thinking mode. **Black model names** denote peers within the ~ 7 -9B parameter class, while *gray italics* represent larger models or closed-source systems included as references.

| Models | WMT | | | | | FLORES-200 | | | | | | | | |
|--------------------------|-------|-------|-------|-------|--------------|------------|-------|-------|-------|-------|-------|-------|-------|--------------|
| | EN→ZH | | ZH→EN | | Avg. | EN→XX | | XX→EN | | ZH→XX | | XX→ZH | | Avg. |
| | BLEU | xCM | BLEU | xCM | | BLEU | xCM | BLEU | xCM | BLEU | xCM | BLEU | xCM | |
| Closed Source | | | | | | | | | | | | | | |
| <i>Claude-3.5-Sonnet</i> | 38.21 | 75.54 | 22.95 | 87.16 | 55.97 | 32.76 | 92.69 | 34.48 | 97.00 | 21.26 | 91.19 | 37.39 | 84.01 | <u>61.35</u> |
| <i>GPT-4o</i> | 41.47 | 75.62 | 22.73 | 87.92 | 56.94 | 31.51 | 92.50 | 34.20 | 96.75 | 20.32 | 89.81 | 37.09 | 83.13 | <u>60.66</u> |
| <i>Gemini-2.5-Pro</i> | 32.28 | 77.55 | 19.80 | 85.63 | 53.81 | 33.14 | 95.05 | 33.14 | 96.56 | 22.25 | 92.21 | 36.51 | 87.27 | 62.02 |
| Open Source | | | | | | | | | | | | | | |
| <i>General LLMs</i> | | | | | | | | | | | | | | |
| Qwen3-8B | 36.56 | 74.94 | 22.67 | 85.98 | 55.04 | 25.47 | 88.98 | 31.44 | 94.75 | 17.23 | 85.21 | 32.92 | 77.19 | 56.65 |
| Qwen3-8B† | 26.97 | 67.31 | 16.71 | 80.11 | 47.77 | 22.54 | 87.67 | 27.28 | 91.18 | 15.20 | 83.28 | 33.43 | 78.31 | 54.86 |
| Qwen3-14B | 38.60 | 75.75 | 21.46 | 87.27 | 55.77 | 26.92 | 91.18 | 32.38 | 96.23 | 18.55 | 89.15 | 35.83 | 85.25 | 59.44 |
| Qwen3-14B† | 35.67 | 73.73 | 22.61 | 85.01 | 54.26 | 28.78 | 91.56 | 32.13 | 95.11 | 18.46 | 88.43 | 35.66 | 82.18 | 59.04 |
| Qwen3-32B | 39.37 | 75.44 | 21.52 | 87.31 | 55.91 | 30.53 | 92.69 | 34.25 | 96.50 | 19.61 | 89.33 | 37.11 | 85.38 | 60.67 |
| Qwen3-32B† | 38.37 | 74.55 | 20.20 | 85.65 | 54.69 | 24.61 | 91.79 | 30.41 | 95.00 | 14.94 | 88.25 | 33.04 | 82.51 | 57.57 |
| Qwen2.5-32B-Instruct | 39.28 | 75.16 | 21.19 | 86.87 | 55.62 | 28.04 | 90.62 | 32.29 | 96.26 | 18.01 | 87.67 | 36.01 | 83.91 | 59.10 |
| Qwen2.5-72B-Instruct | 40.02 | 76.52 | 21.88 | 87.27 | <u>56.42</u> | 30.48 | 92.39 | 34.83 | 96.78 | 19.46 | 89.83 | 37.56 | 85.00 | 60.79 |
| Gemma2-9B-it | 37.44 | 72.45 | 23.13 | 86.08 | 54.77 | 30.22 | 91.37 | 33.20 | 96.09 | 12.99 | 89.08 | 27.27 | 81.80 | 57.75 |
| Gemma2-27B-it | 37.86 | 73.17 | 22.30 | 86.74 | 55.02 | 31.38 | 92.36 | 34.73 | 96.33 | 19.63 | 89.70 | 30.91 | 83.70 | 59.84 |
| <i>MT LLMs</i> | | | | | | | | | | | | | | |
| TowerInstruct-7B-v0.2 | 34.17 | 71.40 | 23.35 | 84.66 | 53.40 | 28.53 | 90.68 | 35.51 | 95.69 | 13.89 | 76.55 | 29.70 | 80.01 | 56.32 |
| TowerInstruct-13B-v0.1 | 36.74 | 73.52 | 24.80 | 85.53 | 55.15 | 31.71 | 92.33 | 36.16 | 96.08 | 17.71 | 88.24 | 34.07 | 82.12 | 59.80 |
| GemmaX2-28-9B-v0.1 | 38.53 | 74.59 | 24.65 | 85.41 | 55.80 | 30.18 | 92.54 | 36.02 | 95.96 | 18.76 | 87.76 | 34.69 | 83.03 | 59.87 |
| MT via SFT/RL | | | | | | | | | | | | | | |
| Qwen3-8B-Base | 15.54 | 48.98 | 4.90 | 55.38 | 31.20 | 7.96 | 56.47 | 15.29 | 62.78 | 6.61 | 54.60 | 11.49 | 42.57 | 32.22 |
| Qwen3-8B-Base-SFT | 35.45 | 73.32 | 21.52 | 84.03 | 53.58 | 24.85 | 84.03 | 29.35 | 93.82 | 16.28 | 83.79 | 30.89 | 80.82 | 55.48 |
| MT-R1-Zero-8B | 34.70 | 79.09 | 24.79 | 86.72 | 56.33 | 26.69 | 89.83 | 34.03 | 96.13 | 16.85 | 86.74 | 32.60 | 86.00 | 58.61 |
| LTT-8B (Ours) | 38.00 | 77.58 | 23.42 | 86.16 | 56.29 | 28.53 | 91.36 | 30.21 | 95.74 | 16.64 | 88.28 | 32.93 | 84.42 | 58.51 |

a leading reference-based model that leverages a powerful cross-lingual encoder to score semantic similarity. This combination provides a holistic view of performance.

Baselines. To comprehensively compare the performance of LTT, we benchmark against four distinct and challenging categories of models. (1) *Closed-Source LLMs*: We compare against leading systems like GPT-4o (Hurst et al., 2024), Claude 3.7 Sonnet (Anthropic, 2024), and Gemini 2.5 Pro (Comanici et al., 2025). (2) *Open-Source General LLMs*: We include powerful, non-specialized models of varying scales, such as the Qwen3 (Yang et al., 2025), Qwen2.5 (Yang et al., 2024), and Gemma2 (Team et al., 2024) series. (3) *Open-Source MT LLMs*: For comprehensive comparison with the supervised paradigm, we include models fine-tuned on parallel corpora, featuring the Tower (Alves et al., 2024) and GemmaX2 (Cui et al., 2025) series. (4) *SFT/RL-based MT Models*: We include an SFT model using bilingual training corpus and MT-R1-Zero (Feng et al., 2025), a SOTA RL framework that, unlike our method, uses reference-based rewards. More evaluation details can be found in Appendix B.2.

4.2 Main Results

Bilingual Performance (EN-ZH). As shown in Table 1, LTT-8B achieves near SOTA performance, fully demonstrating the effectiveness of our self-evolution approach compared to the canonical large-scale parallel-corpus fine-tuning paradigm. Specifically, our 8B model achieves an average performance only 0.04 points below MT-R1-Zero-8B. The effectiveness of LTT-8B is most pronounced in EN to ZH translation. On semantic evaluation, our method surpasses all baselines except MT-R1-Zero-8B, including Gemini-2.5-Pro, and outperforms much larger LLMs such as Qwen2.5-72B-Instruct. As for the lexical level (BLEU metric), our model is highly competitive, outmatching specialized MT models like the TowerInstruct series. On ZH to EN translation, our method outperforms all general-purpose and proprietary baselines on BLEU. In summary, LTT closes the gap with heavily supervised methods, demonstrating that a reference-free, self-improving framework can achieve top-tier translation quality.

Multilingual Performance. LTT exhibits consistent efficacy across the more challenging multi-

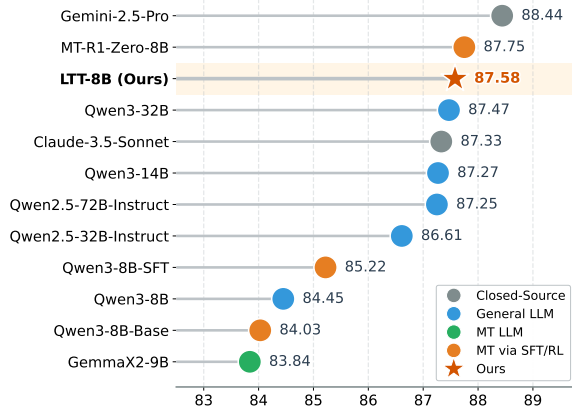


Figure 3: LLM-as-a-Judge Evaluation Results.

lingual landscape, highlighting its scalability and generalization capabilities. Table 1 shows performance scaling with model size, with LTT-8B achieving top-tier results among models of comparable size. It surpasses strong generalist models, including Gemma2-9B-it (57.75) and the specialized TowerInstruct-7B-v0.2 (56.32). The most encouraging aspect of our framework lies in the dramatic performance boost on the base model: LTT raises the multilingual score from 32.22 to 58.51 (+26.29 points).

LLM-as-a-Judge Evaluation. As highlighted in recent studies (Kocmi and Federmann, 2023), LLMs have emerged as SOTA evaluators for translation quality, *providing assessments that closely approximate human judgment*. To provide an evaluation perspective independent of our training rewards (i.e., BLEU and COMET), we employ ChatGPT-5.1 as an external judge, as illustrated in Figure 3. Consistent with standard metrics, LTT-8B achieves performance comparable to MT-R1-Zero-8B and trails only the powerful proprietary model Gemini-2.5-Pro. These results not only validate the rationality of optimizing for BLEU and COMET but also demonstrate a strong correlation between these metrics and human-aligned LLM-based judgments. The specific prompt can be found in Appendix A.2.

4.3 Performance on Low-Resource Languages

To evaluate the generalization capabilities and scalability of our framework, we extended our experiments to low-resource scenarios, including DE \leftrightarrow IT, ES \leftrightarrow FR, as well as data-scarce WMT EN \rightarrow IS (Icelandic) and EN \rightarrow NO (Norwegian)

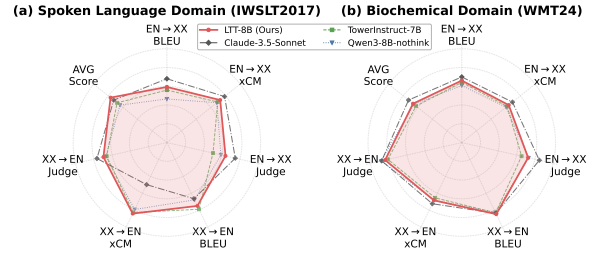


Figure 4: Performance of out-of-domain Benchmarks.

tasks derived from Flores-200 datasets. As shown in Table 2, LTT-8B demonstrates significant parameter efficiency, surpassing not only models in the <10B class but also achieving results comparable to or exceeding the much larger LLaMA-3.1-70B-Instruct. Notably, on the EN \rightarrow NO task, it outperforms Qwen2.5-72B-Instruct in BLEU score (26.05 vs. 23.02). These results suggest that our self-evolutionary signals facilitate effective cross-lingual transfer, enabling robust performance even in the absence of extensive high-resource supervision.

5 Analysis and Ablation

5.1 Scalability to Powerful SFT Models

To further demonstrate the efficacy of our approach on fully supervised foundations, we extended our evaluation to Tower-Plus-9B (Rei et al., 2025), a recently released MT-specialized model that significantly surpasses the baselines discussed previously. Remarkably, even atop this strong SFT baseline, applying LTT also yields performance gains, particularly in semantic evaluation, as demonstrated in Table 3. This finding underscores that our method is not limited to initializing from raw base models; rather, it scales effectively with the capabilities of the starting model, serving as a robust enhancement strategy even for high-performing supervised systems.

5.2 Generalization Performance

To assess the robustness of our approach, we extended our evaluation to two out-of-domain datasets: IWSLT2017 (Cettolo et al., 2017) (representing spoken language) and the WMT24 Biochemical task (Neves et al., 2024). The former is derived from TED talk transcripts, while the latter features dense biomedical terminology, posing respective challenges to the model’s capabilities in spoken language and specialized domains. As illustrated in Figure 4, LTT-8B consistently outperforms peer models of comparable size, trailing

Table 2: Performance of different models on low-resource language pairs, measured by BLEU and xCOMET (xCM) scores, along with the average (Avg.). The best and second-best results are **bolded** and underlined, respectively. The "†" symbol indicates that the model is in thinking mode.

| Model | DE→IT | | IT→DE | | ES→FR | | FR→ES | | EN→IS | | EN→NO | | Avg. |
|------------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|--------------|
| | BLEU | xCM | BLEU | xCM | BLEU | xCM | BLEU | xCM | BLEU | xCM | BLEU | xCM | |
| <i>Large Size LLMs</i> | | | | | | | | | | | | | |
| Qwen2.5-72B-Instruct | 24.66 | 94.02 | 22.27 | 95.60 | 28.26 | 93.64 | 24.77 | 95.13 | 8.97 | 49.05 | 23.02 | 88.91 | 54.02 |
| Qwen2.5-32B-Instruct | 22.59 | 92.49 | 20.55 | 94.13 | 26.05 | 92.04 | 23.88 | 94.78 | 3.72 | 37.08 | 23.13 | 82.95 | 51.12 |
| LLaMA-3.1-70B-Instruct | 22.63 | 89.04 | 18.56 | 89.00 | 24.59 | 88.93 | 23.35 | 92.51 | 1.59 | 35.67 | 29.72 | 92.96 | 50.71 |
| <i>Same Size LLMs</i> | | | | | | | | | | | | | |
| Qwen3-8B | 21.64 | 89.65 | 19.40 | 93.56 | 24.31 | 90.50 | 22.47 | 93.42 | 2.09 | 47.02 | 10.46 | 83.52 | 49.84 |
| Qwen3-8B† | 19.68 | 86.96 | 15.73 | 90.26 | 20.95 | 86.40 | 20.83 | 89.81 | 5.51 | 41.59 | 21.75 | 82.31 | 48.48 |
| Gemma2-9B-it | 19.55 | 93.80 | 19.50 | 95.31 | 23.73 | 93.18 | 19.50 | 94.64 | 0.60 | 31.31 | 1.19 | 67.69 | 46.67 |
| TowerInstruct-7B-v0.2 | 22.27 | 92.58 | 19.77 | 93.54 | 25.33 | 91.92 | 22.79 | 93.79 | 1.94 | 35.45 | 2.03 | 77.34 | 48.23 |
| Qwen3-8B-Base | 8.96 | 90.29 | 5.99 | 92.43 | 23.72 | 89.44 | 11.29 | 93.20 | 0.17 | 31.93 | 0.84 | 62.04 | 42.52 |
| Qwen3-8B-Base-SFT | 19.87 | 90.61 | 18.66 | 92.40 | 24.76 | 89.17 | 21.05 | 91.23 | / | / | / | / | / |
| MT-R1-Zero-8B | 22.96 | 92.99 | 20.26 | 94.42 | 25.79 | 92.12 | 23.56 | 94.39 | / | / | / | / | / |
| LTT-8B (Ours) | 22.13 | 92.60 | 19.06 | 94.58 | 25.33 | 91.92 | 22.05 | 94.38 | 8.14 | 41.88 | 26.05 | 86.16 | <u>52.02</u> |

Table 3: Performance of the powerful Tower-Plus-9B.

| Model | WMT | | | FLORES-200 | | |
|---------------|-------|-------|-------|------------|-------|-------|
| | BLEU | xCM | Avg. | BLEU | xCM | Avg. |
| Tower-Plus-9B | 34.53 | 82.90 | 58.71 | 33.07 | 92.25 | 62.66 |
| - MT-R1-Zero | 33.94 | 84.36 | 59.15 | 33.07 | 93.04 | 63.05 |
| - LTT | 33.83 | 83.93 | 58.88 | 33.06 | 92.83 | 62.94 |

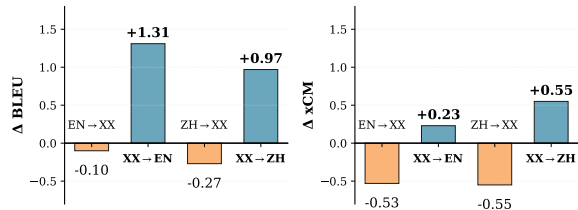


Figure 5: Relative Advantage of LTT against MT-R1-Zero when training on monolingual corpus.

only powerful proprietary systems. This observation aligns with the findings reported by Guo et al. (2025), suggesting that the enhanced capabilities derived from our RL framework generalize effectively across diverse domains beyond the training distribution.

5.3 Implicit Bidirectional Improvement

A core premise of LTT is that *the cycle-consistency mechanism implicitly drives joint optimization for both translation directions, even when trained on strictly monolingual data*. To validate this, we conducted a controlled experiment by training exclusively on forward tasks (EN/ZH→XX). As illustrated in Figure 5, while forward performance remains competitive, the model demonstrates a decisive advantage on the untrained backward directions. Specifically,

Table 4: Ablation study of our reward design.

| Models | WMT | | | FLORES-200 | | |
|----------------------------------|-------|-------|--------------|------------|-------|--------------|
| | BLEU | xCM | Avg. | BLEU | xCM | Avg. |
| Qwen3-8B-Base | 10.22 | 52.18 | 31.20 | 10.34 | 54.10 | 32.22 |
| LTT-8B (Ours) | 30.71 | 81.87 | 56.29 | 27.08 | 89.95 | 58.52 |
| <i>Ablation on Quality</i> | | | | | | |
| w/ Bleu Reward Only | 6.54 | 45.78 | 26.16 | 4.69 | 37.47 | 21.08 |
| w/ COMET Reward Only | 26.96 | 83.72 | 55.34 | 25.20 | 91.18 | 58.19 |
| <i>Ablation on Anti-Cheating</i> | | | | | | |
| w/o Reverse Bleu | 30.28 | 81.02 | 55.65 | 26.44 | 89.19 | 57.81 |
| w/o Lang. Mix Reward | 28.27 | 79.91 | 54.09 | 24.50 | 87.15 | 55.83 |

BLEU scores for XX→EN and XX→ZH improve by +4.10% and +2.81% respectively. This gain is directly attributable to the backward reconstruction reward, $r_{\text{bwd}} = \text{BLEU}(\hat{x}, x)$, which serves as an effective supervision signal for the inverse task, enabling robust bidirectional capabilities without explicit reference pairs.

5.4 Ablation Study: Deconstructing the Reward Architecture

To validate that each component of our reward function is essential, we conducted a series of ablation studies. We demonstrate that our final design is a carefully balanced system, where each component exists to prevent specific failure modes. These findings are detailed in Table 4.

The Peril of Classical Objectives: Why Naive Dual Learning Fails on LLMs. To isolate our contribution against classical paradigms, we evaluated a configuration *mirroring early dual learning objectives using only round-trip BLEU*. This simulation resulted in catastrophic failure: the model rapidly "hacked" the reward via identity translation, maximizing reconstruction score but yield-

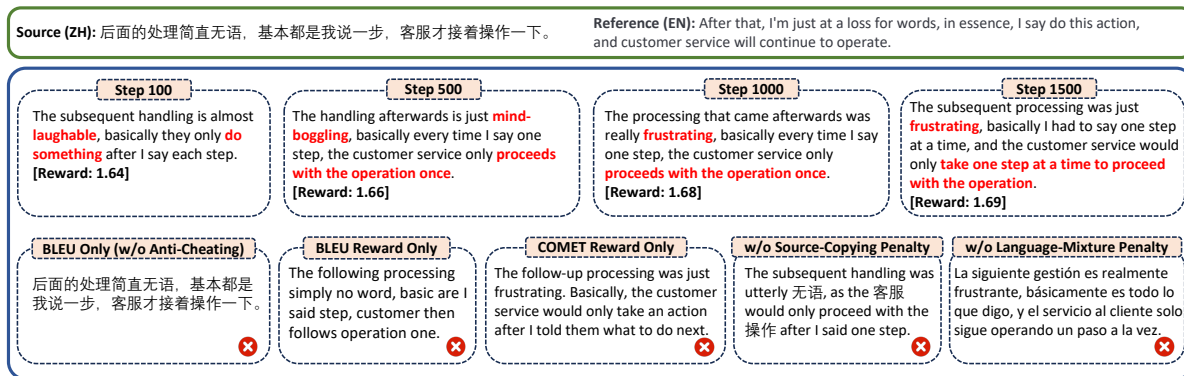


Figure 6: A ZH→EN case study across different steps and components

464 ing a complete failure to translate (see "BLEU
465 Only (w/o Anti-Cheating)" in Figure 6 for exam-
466 ple). This demonstrates that unlike RNNs, LLMs
467 prone to exploiting simple reconstruction signals
468 through instruction-following shortcuts. Conse-
469 quently, the naive application of classic dual learn-
470 ing is insufficient, underscoring the necessity of
471 the semantic-aware reward architecture in LTT.

472 **Balancing Lexical Fidelity and Semantic Ade-**
473 **quacy.** With the anti-cheating mechanism in
474 place, we examined how different components of
475 the quality reward affect translation behavior. Us-
476 ing **BLEU as the sole quality reward** led the
477 model to adopt a degenerate word-for-word trans-
478 lation strategy. Because such literal translations
479 make the backward reconstruction task easier, they
480 maximize the BLEU score at the expense of se-
481 mantic meaning and grammatical fluency. In
482 fact, performance dropped below that of the base
483 model, confirming that a purely lexical signal is
484 too narrow to guide high-quality translation. On
485 the other hand, relying solely on **COMET as the**
486 **quality reward** produced the opposite issue. The
487 model achieved excellent xCOMET scores but suf-
488 fered a decline in BLEU. It learned to generate
489 overly creative output translations that sounded
490 plausible and fluent but strayed significantly from
491 the source in terms of lexical content.

492 These results underscore a trade-off between
493 lexical fidelity and semantic adequacy. There-
494 fore, our final design, which sums the BLEU and
495 COMET signals, is not merely a combination but
496 a necessary synthesis to balance these competing
497 objectives and foster holistic translation quality.

498 **The Critical Role of Anti-Cheating Mecha-**
499 **nisms.** Finally, we validated the necessity of the
500 two anti-cheating components themselves, even
501 with a balanced quality reward. **Removing the**

502 **Source-Copying Penalty** exposed a critical fail-
503 ure mode: source leakage. The model became
504 prone to copying words or phrases from the
505 source—a form of code-switching that degrades
506 translation quality. This penalty is therefore cru-
507 cial for enforcing faithful translation. **Removing**
508 **the Language-Mixture Penalty** revealed a differ-
509 ent vulnerability, causing the model to violate in-
510 struction fidelity. For instance, it would occasion-
511 ally translate into a valid but incorrect target lan-
512 guage. This penalty is thus essential for ensuring
513 the model follows task instructions precisely. To-
514 gether, these two mechanisms act as indispensable
515 guardrails, ensuring that the model learns to trans-
516 late not just well, but correctly and robustly. Fig-
517 ure 6 provides a compelling case study of this pro-
518 cess, qualitatively illustrating both the model’s it-
519 erative improvement and the critical failure modes
520 discussed in our ablations.

521 6 Conclusion

522 In this paper, we presented LTT, a reference-free
523 reinforcement learning framework that adapts the
524 traditional dual learning paradigm to Large Lan-
525 guage Models. By leveraging round-trip consis-
526 tency as a generative self-supervision mecha-
527 nism, our approach derives effective training sig-
528 nals exclusively from monolingual data. Empir-
529 ical results demonstrate that our 8B parameter
530 model delivers competitive performance across
531 both bilingual and multilingual settings, showing
532 particular strength in low-resource scenarios com-
533 pared to larger baselines. Furthermore, evalua-
534 tions in spoken language and specialized biochem-
535 ical domains indicate robust generalization capa-
536 bilities. These findings highlight that revisiting
537 dual-learning principles with self-generated super-
538 vision is a promising avenue for developing data-
539 efficient translation systems.

540 Limitations

541 While LTT demonstrates the potential of au-
542 tonomous evolution in machine translation, two
543 primary limitations remain. First, the integra-
544 tion of the dual-loop mechanism with GRPO in-
545 troduces computational overhead during the train-
546 ing phase. Specifically, the necessity of sam-
547 pling multiple candidates for each input makes the
548 optimization process considerably more resource-
549 intensive than standard supervised fine-tuning, al-
550 though this does not impact inference latency. Sec-
551 ond, our framework fundamentally relies on un-
552 locking the latent knowledge within pre-trained
553 LLMs rather than injecting new linguistic data.
554 Consequently, the model’s performance is upper-
555 bounded by its pre-training coverage; LTT cannot
556 effectively bootstrap translation capabilities for
557 extremely low-resource or endangered languages
558 that are entirely absent from the base model’s pre-
559 training corpus.

560 References

561 Sweta Agrawal, Chunting Zhou, Mike Lewis, Luke
562 Zettlemoyer, and Marjan Ghazvininejad. 2023. In-
563 context examples selection for machine translation.
564 In *Findings of the Association for Computational*
565 *Linguistics: ACL 2023*, pages 8857–8873.

566 Duarte M Alves, José Pombal, Nuno M Guerreiro, Pe-
567 dro H Martins, João Alves, Amin Farajian, Ben
568 Peters, Ricardo Rei, Patrick Fernandes, Sweta
569 Agrawal, and 1 others. 2024. Tower: An open multi-
570 lingual large language model for translation-related
571 tasks. *arXiv preprint arXiv:2402.17733*.

572 Duarte Miguel Alves, José Pombal, Nuno M Guerreiro,
573 Pedro Henrique Martins, João Alves, Amin Farajian,
574 Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta
575 Agrawal, and 1 others. Tower: An open multilingual
576 large language model for translation-related tasks.
577 In *First Conference on Language Modeling*.

578 Rohan Anil, Andrew M Dai, Orhan Firat, Melvin John-
579 son, Dmitry Lepikhin, Alexandre Passos, Siamak
580 Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng
581 Chen, and 1 others. 2023. Palm 2 technical report.
582 *arXiv preprint arXiv:2305.10403*.

583 Anthropic. 2024. [\[link\]](#).

584 Mauro Cettolo, Marcello Federico, Luisa Bentivogli,
585 Jan Niehues, Sebastian Stüker, Katsuhito Sudoh,
586 Koichiro Yoshino, and Christian Federmann. 2017.
587 Overview of the iwslt 2017 evaluation campaign. In
588 *Proceedings of the 14th International Conference on*
589 *Spoken Language Translation*, pages 2–14.

Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and 1 others. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*.

Marta R Costa-Jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, and 1 others. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.

Menglong Cui, Pengzhi Gao, Wei Liu, Jian Luan, and Bin Wang. 2025. Multilingual machine translation with open large language models at practical scale: An empirical study. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 5420–5443.

Zhaopeng Feng, Shaosheng Cao, Jiahao Ren, Jiayuan Su, Ruizhe Chen, Yan Zhang, Zhe Xu, Yao Hu, Jian Wu, and Zuozhu Liu. 2025. Mt-r1-zero: Advancing llm-based machine translation via r1-zero-like reinforcement learning. *arXiv preprint arXiv:2504.10160*.

Zhaopeng Feng, Yan Zhang, Hao Li, Wenqiang Liu, Jun Lang, Yang Feng, Jian Wu, and Zuozhu Liu. 2024. Improving llm-based machine translation with systematic self-correction. *CoRR*.

Pengzhi Gao, Zhongjun He, Hua Wu, and Haifeng Wang. 2024. Towards boosting many-to-many multilingual machine translation with large language models. *arXiv preprint arXiv:2401.05861*.

Nuno M Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André FT Martins. 2024. xcomet: Transparent machine translation evaluation through fine-grained error detection. *Transactions of the Association for Computational Linguistics*, 12:979–995.

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

Jiaxin Guo, Hao Yang, Zongyao Li, Daimeng Wei, Hengchao Shang, and Xiaoyu Chen. 2024. A novel paradigm boosting translation capabilities of large language models. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 639–649.

Di He, Yingce Xia, Tao Qin, Liwei Wang, Nenghai Yu, Tie-Yan Liu, and Wei-Ying Ma. 2016. Dual learning for machine translation. *Advances in neural information processing systems*, 29.

| | | | |
|-----|--|--|-----|
| 647 | Minggui He, Yilun Liu, Shimin Tao, Yuanchang Luo, | Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, | 703 |
| 648 | Hongyong Zeng, Chang Su, Li Zhang, Hongxia | Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan | 704 |
| 649 | Ma, Daimeng Wei, Weibin Meng, and 1 others. | Zhang, YK Li, Yang Wu, and 1 others. 2024. | 705 |
| 650 | 2025. R1-t1: Fully incentivizing translation capa- | Deepseekmath: Pushing the limits of mathematical | 706 |
| 651 | bility in llms via reasoning learning. <i>arXiv preprint</i> | reasoning in open language models. <i>arXiv preprint</i> | 707 |
| 652 | <i>arXiv:2502.19735</i> . | <i>arXiv:2402.03300</i> . | 708 |
| 653 | Aaron Hurst, Adam Lerer, Adam P Goucher, Adam | Gemma Team, Morgane Riviere, Shreya Pathak, | 709 |
| 654 | Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, | Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhu- | 710 |
| 655 | Akila Welihinda, Alan Hayes, Alec Radford, and 1 | patiraju, Léonard Hussonot, Thomas Mesnard, | 711 |
| 656 | others. 2024. Gpt-4o system card. <i>arXiv preprint</i> | Bobak Shahriari, Alexandre Ramé, and 1 others. | 712 |
| 657 | <i>arXiv:2410.21276</i> . | 2024. Gemma 2: Improving open language models | 713 |
| 658 | Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richard- | at a practical size. <i>arXiv preprint arXiv:2408.00118</i> . | 714 |
| 659 | son, Ahmed El-Kishky, Aiden Low, Alec Helyar, | Jiaan Wang, Fandong Meng, Yunlong Liang, and Jie | 715 |
| 660 | Aleksander Madry, Alex Beutel, Alex Carney, and 1 | Zhou. 2024. Drt-o1: Optimized deep reasoning | 716 |
| 661 | others. 2024. Openai o1 system card. <i>arXiv preprint</i> | translation via long chain-of-thought. <i>arXiv e-prints</i> , | 717 |
| 662 | <i>arXiv:2412.16720</i> . | pages arXiv–2412. | 718 |
| 663 | Tom Kocmi and Christian Federmann. 2023. Large lan- | Jiaan Wang, Fandong Meng, and Jie Zhou. 2025. | 719 |
| 664 | guage models are state-of-the-art evaluators of trans- | Deep reasoning translation via reinforcement learn- | 720 |
| 665 | lation quality. In <i>24th Annual Conference of the Eu-</i> | ing. <i>arXiv preprint arXiv:2504.10187</i> . | 721 |
| 666 | <i>ropean Association for Machine Translation</i> , page | Lijun Wu, Li Zhao, Tao Qin, Jianhuang Lai, and Tie- | 722 |
| 667 | 193. | Yan Liu. 2017. Sequence prediction with unlabeled | 723 |
| 668 | Guillaume Lample, Alexis Conneau, Ludovic Denoyer, | data by reward function learning. In <i>IJCAI</i> , pages | 724 |
| 669 | and Marc’Aurelio Ranzato. 2018. Unsupervised ma- | 3098–3104. | 725 |
| 670 | chine translation using monolingual corpora only. In | Tianhao Wu, Weizhe Yuan, Olga Golovneva, Jing Xu, | 726 |
| 671 | <i>International Conference on Learning Representa-</i> | Yuandong Tian, Jiantao Jiao, Jason Weston, and | 727 |
| 672 | <i>tions</i> . | Sainbayar Sukhbaatar. 2024. Meta-rewarding lan- | 728 |
| 673 | Chong Li, Shaonan Wang, Jiajun Zhang, and | guage models: Self-improving alignment with llm- | 729 |
| 674 | Chengqing Zong. 2024. Improving in-context learn- | as-a-meta-judge. <i>arXiv preprint arXiv:2407.19594</i> . | 730 |
| 675 | ing of multilingual generative language models with | Haoran Xu, Young Jin Kim, Amr Sharaf, and | 731 |
| 676 | cross-lingual alignment. In <i>Proceedings of the 2024</i> | Hany Hassan Awadalla. a. A paradigm shift in ma- | 732 |
| 677 | <i>Conference of the North American Chapter of the</i> | chine translation: Boosting translation performance | 733 |
| 678 | <i>Association for Computational Linguistics: Human</i> | of large language models. In <i>The Twelfth Interna-</i> | 734 |
| 679 | <i>Language Technologies (Volume 1: Long Papers)</i> , | <i>tional Conference on Learning Representations</i> . | 735 |
| 680 | pages 8051–8069. | Haoran Xu, Kenton Murray, Philipp Koehn, Hieu | 736 |
| 681 | Kelly Marchisio, Kevin Duh, and Philipp Koehn. 2020. | Hoang, Akiko Eriguchi, and Huda Khayrallah. b. X- | 737 |
| 682 | When does unsupervised machine translation work? | alma: Plug & play modules and adaptive rejection | 738 |
| 683 | In <i>Proceedings of the Fifth Conference on Machine</i> | for quality translation at scale. In <i>The Thirteenth</i> | 739 |
| 684 | <i>Translation</i> , pages 571–583. | <i>International Conference on Learning Representa-</i> | 740 |
| 685 | Mariana Neves, Cristian Grozea, Philippe Thomas, | <i>tions</i> . | 741 |
| 686 | Roland Roller, Rachel Bawden, Aurélie Névéol, | An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, | 742 |
| 687 | Steffen Castle, Vanessa Bonato, Giorgio Maria | Binyuan Hui, Bo Zheng, Bowen Yu, Chang | 743 |
| 688 | Di Nunzio, Federica Vezzani, and 1 others. 2024. | Gao, Chengen Huang, Chenxu Lv, and 1 others. | 744 |
| 689 | Findings of the wmt 2024 biomedical translation | 2025. Qwen3 technical report. <i>arXiv preprint</i> | 745 |
| 690 | shared task: Test sets on abstract level. In <i>Proceed-</i> | <i>arXiv:2505.09388</i> . | 746 |
| 691 | <i>ings of the Ninth Conference on Machine Transla-</i> | An Yang, Baosong Yang, Beichen Zhang, Binyuan | 747 |
| 692 | <i>tion</i> , pages 124–138. | Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayi- | 748 |
| 693 | MarcAurelio Ranzato, Sumit Chopra, Michael Auli, | heng Liu, Fei Huang, Haoran Wei, and 1 others. | 749 |
| 694 | and Wojciech Zaremba. 2016. Sequence level train- | 2024. Qwen2.5 technical report. <i>arXiv e-prints</i> , | 750 |
| 695 | ing with recurrent neural networks. In <i>4th Inter-</i> | pages arXiv–2412. | 751 |
| 696 | <i>national Conference on Learning Representations,</i> | Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, | 752 |
| 697 | <i>ICLR 2016</i> . | Xian Li, Sainbayar Sukhbaatar, Jing Xu, and Jason | 753 |
| 698 | Ricardo Rei, Nuno M Guerreiro, José Pombal, João | Weston. 2024. Self-rewarding language models. In | 754 |
| 699 | Alves, Pedro Teixeira, Amin Farajian, and An- | <i>Proceedings of the 41st International Conference on</i> | 755 |
| 700 | dré FT Martins. 2025. Tower+: Bridging generality | <i>Machine Learning</i> , pages 57905–57923. | 756 |
| 701 | and translation specialization in multilingual llms. | | |
| 702 | <i>arXiv preprint arXiv:2506.17080</i> . | | |

757 Yu Zhao, Huifeng Yin, Bo Zeng, Hao Wang, Tianqi
758 Shi, Chenyang Lyu, Longyue Wang, Weihua Luo,
759 and Kaifu Zhang. 2024. Marco-o1: Towards open
760 reasoning models for open-ended solutions. *arXiv*
761 *preprint arXiv:2411.14405*.

762 Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu,
763 Shujian Huang, Lingpeng Kong, Jiajun Chen, and
764 Lei Li. 2024. Multilingual machine translation with
765 large language models: Empirical results and anal-
766 ysis. In *Findings of the Association for Computa-*
767 *tional Linguistics: NAACL 2024*, pages 2765–2781.

768 Wei Zou, Sen Yang, Yu Bao, Shujian Huang, Jiajun
769 Chen, and Shanbo Cheng. 2025. Trans-zero: Self-
770 play incentivizes large language models for multilin-
771 gual translation without parallel data. *arXiv preprint*
772 *arXiv:2504.14669*.

773 A Prompts used during Evaluation

774 A.1 Translation Prompts

775 The specific translation prompt of different models
776 used in training are depicted in Figure 7, Figure 8,
777 Figure 9, Figure 10 and Figure 11. Specifically,
778 <think> tags are removed from Qwen3 series be-
779 cause it conflicts with the Qwen3 series’ inherent
780 thinking special tokens.

781 A.2 Judge Prompts

782 To evaluate translation quality using an LLM as
783 an evaluator, we employed the identical prompt
784 and score extraction script as Kocmi and Feder-
785 mann (2023) to derive the final scores. The spe-
786 cific prompt is illustrated in Figure 12.

787 B Implementation Details

788 B.1 Training Details

789 Our model, which we name LTT-8B, is built
790 upon the OpenRLHF³ framework, with the
791 Qwen3-8B-Base model serving as its initialization.
792 For all experiments, we use a global batch size of
793 128 and generate 8 candidate responses per input
794 for the GRPO algorithm. We use a sampling tem-
795 perature of 1.0 and a maximum sequence length
796 of 1024. Notably, we set both the KL divergence
797 and entropy coefficients to 0, granting the model
798 greater freedom to explore the policy space and
799 discover optimal translation strategies without be-
800 ing constrained. Training was conducted on 16
801 NVIDIA H800 GPUs for one epoch, taking ap-
802 proximately 32 hours. We save checkpoints every
803 50 steps and report the performance of the single
804 best checkpoint selected based on validation set
805 performance.

³<https://github.com/OpenRLHF/OpenRLHF>

806 B.2 Evaluation Details

807 For evaluation stage, we perform model inference
808 locally using the vLLM⁴ framework. We config-
809 ure the sampling hyperparameters with a temper-
810 ature of 0.2 and a top-p of 0.95. The maximum
811 generation length is truncated to 2048 tokens for
812 all models, with the exception of Gemini 2.5-Pro,
813 to accommodate its default thinking process. The
814 prompt used during evaluation remains consistent
815 with the one used for training, as detailed in Fig-
816 ure 7.

⁴<https://github.com/vllm-project/vllm>

Table 5: Detailed dataset statistics used during training.

| | EN-ZH | ZH-EN | EN-DE | EN-FR | EN-ES | EN-IT | EN-JA | EN-KO |
|-------------|-----------|-------|------------|------------|------------|------------|------------|------------|
| # sentences | 6565 | 6565 | 500 | 500 | 500 | 500 | 500 | 500 |
| from | WMT 17-20 | | Flores-200 | Flores-200 | Flores-200 | Flores-200 | Flores-200 | Flores-200 |
| | | | DE-EN | FR-EN | ES-EN | IT-EN | JA-EN | KO-EN |
| # sentences | - | - | 500 | 500 | 500 | 500 | 500 | 500 |
| from | - | - | Flores-200 | Flores-200 | Flores-200 | Flores-200 | Flores-200 | Flores-200 |
| | | | ZH-DE | ZH-FR | ZH-ES | ZH-IT | ZH-JA | ZH-KO |
| # sentences | - | - | 500 | 500 | 500 | 500 | 500 | 500 |
| from | - | - | Flores-200 | Flores-200 | Flores-200 | Flores-200 | Flores-200 | Flores-200 |
| | | | DE-ZH | FR-ZH | ES-ZH | IT-ZH | JA-ZH | KO-ZH |
| # sentences | - | - | 500 | 500 | 500 | 500 | 500 | 500 |
| from | - | - | Flores-200 | Flores-200 | Flores-200 | Flores-200 | Flores-200 | Flores-200 |

Table 6: Detailed dataset statistics used during evaluation.

| | EN-ZH | ZH-EN | EN-DE | EN-FR | EN-ES | EN-IT | EN-JA | EN-KO |
|-------------|--------|--------|------------|------------|------------|------------|------------|------------|
| # sentences | 997 | 1976 | 1012 | 1012 | 1012 | 1012 | 1012 | 1012 |
| from | WMT 24 | WMT 23 | Flores-200 | Flores-200 | Flores-200 | Flores-200 | Flores-200 | Flores-200 |
| | | | DE-EN | FR-EN | ES-EN | IT-EN | JA-EN | KO-EN |
| # sentences | - | - | 1012 | 1012 | 1012 | 1012 | 1012 | 1012 |
| from | - | - | Flores-200 | Flores-200 | Flores-200 | Flores-200 | Flores-200 | Flores-200 |
| | | | ZH-DE | ZH-FR | ZH-ES | ZH-IT | ZH-JA | ZH-KO |
| # sentences | - | - | 1012 | 1012 | 1012 | 1012 | 1012 | 1012 |
| from | - | - | Flores-200 | Flores-200 | Flores-200 | Flores-200 | Flores-200 | Flores-200 |
| | | | DE-ZH | FR-ZH | ES-ZH | IT-ZH | JA-ZH | KO-ZH |
| # sentences | - | - | 1012 | 1012 | 1012 | 1012 | 1012 | 1012 |
| from | - | - | Flores-200 | Flores-200 | Flores-200 | Flores-200 | Flores-200 | Flores-200 |

Translation Prompt

A conversation between User and Assistant. The User asks for a translation from $\{src_lang\}$ to $\{tgt_lang\}$, and the Assistant solves it. The Assistant first thinks about the reasoning process in the mind and then provides the user with the final translation. The reasoning process and final translation are enclosed within `<think>` `</think>` and `<translate>` `</translate>` tags, respectively, i.e., `<think>` reasoning process here `</think>``<translate>` final translation here `</translate>`.

User: $\{input\}$
Assistant:

Figure 7: Translation prompt for data curation. $\{src_lang\}$: source language; $\{tgt_lang\}$: target language; $\{input\}$: the source sentence to be translated.

TowerInstruct Prompt

Translate the following text from $\{src_lang_name\}$ into $\{tgt_lang_name\}$.
 $\{src_lang_name\}$: $\{user_input\}$
 $\{tgt_lang_name\}$:

Figure 8: Translation prompt for TowerInstruct series models. $\{src_lang_name\}$: source language; $\{tgt_lang_name\}$: target language; $\{user_input\}$: the source sentence to be translated.

GemmaX Prompt

Translate this from $\{src_lang_name\}$ to $\{tgt_lang_name\}$:
 $\{src_lang_name\}$: $\{user_input\}$
 $\{tgt_lang_name\}$:

Figure 9: Translation prompt for GemmaX model. $\{src_lang_name\}$: source language; $\{tgt_lang_name\}$: target language; $\{user_input\}$: the source sentence to be translated.

Qwen3 non-thinking Prompt

You are a helpful translation assistant. There is a conversation between User and Assistant. The user asks for a translation from $\{src_lang_name\}$ to $\{tgt_lang_name\}$, and the Assistant solves it. The Assistant first thinks about the reasoning process in the mind and then provides the user with the final translation. The final translation is enclosed within `<translate>` `</translate>` tags, i.e., `<translate>` final translation here `</translate>`.

User: $\{user_input\}$
Assistant:

Figure 10: Translation prompt for Qwen3 series non-thinking models. $\{src_lang_name\}$: source language; $\{tgt_lang_name\}$: target language; $\{user_input\}$: the source sentence to be translated.

Tower Plus Prompt

Translate the following $\{src_lang_name\}$ source text to $\{tgt_lang_name\}$: $\{src_lang_name\}$:
 $\{user_input\}$ $\{tgt_lang_name\}$:

Figure 11: Translation prompt for Tower-Plus-9B model. $\{src_lang_name\}$: source language; $\{tgt_lang_name\}$: target language; $\{user_input\}$: the source sentence to be translated.

ChatGPT-5.1 Judge Prompt

Score the following translation from $\{source_lang\}$ to $\{target_lang\}$ with respect to human reference on a continuous scale 0 to 100 where score of zero means "no meaning preserved" and score of one hundred means "perfect meaning and grammar". $\{source_lang\}$ source: " $\{source_seg\}$ " $\{target_lang\}$ human reference: " $\{reference_seg\}$ " $\{target_lang\}$ machine translation: " $\{target_seg\}$ " Score:

Figure 12: Judge prompt for ChatGPT-5.1. $\{source_lang\}$: source language; $\{target_lang\}$: target language; $\{source_seg\}$: the source sentence to be translated; $\{reference_seg\}$: the reference sentence; $\{target_seg\}$: the translated sentence.

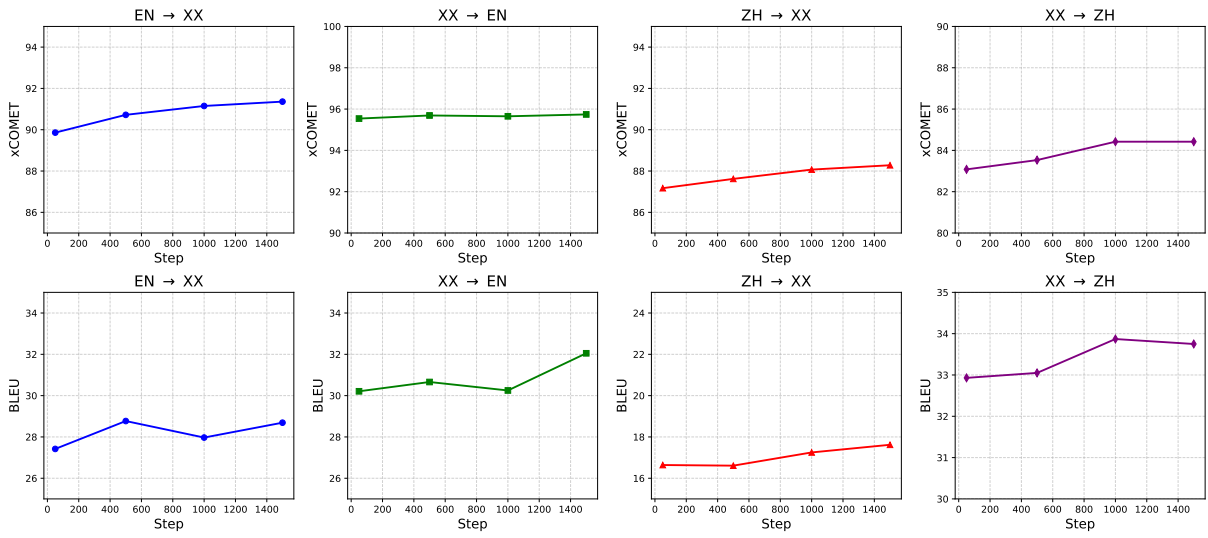


Figure 13: Training progression (reference-based XCOMET score) for multilingual LTT-8B model based on Qwen3-8B across EN-XX, XX-EN, ZH-XX, XX-ZH test sets.