

# Incentivizing Parametric Knowledge via Reinforcement Learning with Verifiable Rewards for Cross-Cultural Entity Translation

Anonymous ACL submission

## Abstract

Cross-cultural entity translation remains challenging for large language models (LLMs) as literal or phonetic renderings are usually yielded instead of culturally appropriate translations in context. However, relevant knowledge may already be encoded in model parameters during large-scale pre-training. To incentivize the effective use of parametric knowledge, we propose **EA-RLVR** (Entity-Anchored Reinforcement Learning with Verifiable Rewards), a training framework that optimizes cross-cultural entity translation without relying on external knowledge bases. EA-RLVR anchors supervision on a verifiable, entity-level reward signal and incorporates lightweight structural gates to stabilize optimization. This design steers the model toward learning a robust reasoning process rather than merely imitating reference translations. We evaluate EA-RLVR on XC-Translate and observe consistent improvements in both entity translation accuracy and out-of-domain generalization. Specifically, training on merely 7k samples boosts Qwen3-14B’s entity translation accuracy from 23.66% to 31.87% on a 50k test set comprising **entirely unseen entities**. The learned entity translation ability also transfers to general translation, yielding +1.35 XCOMET on WMT24pp, which scales to +1.59 with extended optimization. Extensive analyses of  $pass@k$  dynamics and reward formulations attribute these gains to superior sampling efficiency and a stable optimization landscape.

## 1 Introduction

At its core, machine translation aspires to make culturally situated texts accessible across languages. Despite substantial progress with multilingual large language models, current systems often fall short of this goal in settings where translation hinges on culturally grounded entities such as books, films, places, songs and idioms (Yao et al., 2024). In these cases, producing an accurate, culture-aligned

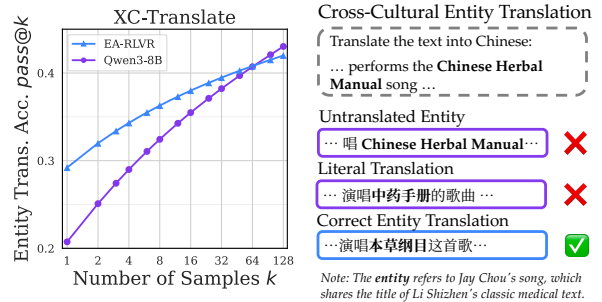


Figure 1: (Left) Entity translation accuracy (%)  $pass@k$  curves demonstrate the base model possesses latent knowledge (high accuracy at large  $k$ ) that EA-RLVR effectively activates at  $k = 1$ . (Right) An illustration of the challenge in cross-cultural entity translation.

translation requires identifying, in context, which real-world entity is being referred to and how it is conventionally named in the target culture (Moghe et al., 2025). Recent evaluations have shown that even frontier proprietary LLMs frequently default to literal or phonetic renderings that are grammatically well formed but semantically inappropriate in context, thereby altering or obscuring the intended meaning of the source text (Conia et al., 2024).

A widely adopted workaround for this limitation is to equip translation systems with external knowledge, e.g., through online retrieval, knowledge graphs, or curated databases (Conia et al., 2024; Khandelwal et al.). These approaches can improve accuracy when relevant information is successfully retrieved. However, they also introduce practical and structural constraints. The performance of such systems depends critically on how well the task aligns with underlying database (Agrawal et al., 2023), and in practice often requires task-specific retrievers that must be trained or tuned (Wang et al., 2025b). Moreover, it fundamentally shifts the bottleneck from contextual entity reasoning to the structure and coverage of the external knowledge source, making translation quality contingent on what can be retrieved.

On the other hand, as trained on corpora spanning trillions of tokens across diverse domains and languages, LLMs implicitly encode a wide range of entity correspondences, cultural references, and real-world usage conventions (Yang et al., 2025; Qwen et al., 2025). In principle, such knowledge should support cross-cultural translation. As illustrated in Figure 1 (Left), the correct cultural entities are often present in the base model’s probability distribution, evidenced by high accuracy when multiple sampling attempts ( $pass@128$ ). However, such knowledge remains effectively inaccessible during standard single-pass generation ( $pass@1$ ). Consequently, models frequently default to verbatim copying or literal renderings that obscure the intended meaning, such as retaining the source term or translating a song title literally as a medical manual (Figure 1, Right). These observations suggest that the core difficulty lies less in the availability of knowledge itself, but more in the absence of mechanisms that incentivize the model to surface that knowledge in a context-sensitive manner.

To incentivize LLMs to leverage their parametric knowledge effectively, we propose **EA-RLVR (Entity-Anchored RL with Verifiable Rewards)**, a framework for cross-cultural entity translation driven by fully rule-based, automatically verifiable reward. We cast cross-cultural entity translation as a sequence decision problem: given a source sentence, the model produces its own candidate translations, and a deterministic verifier evaluates whether the output expresses the correct target-culture entities. Rather than imitating reference translations, the model learns from verifiable rewards assigned to its own trajectories, reinforcing the reasoning that produces the correct entities. Concretely, EA-RLVR uses an **entity-matching reward** based on normalized substring matching between the predicted entity and the gold entity set. To stabilize optimization and reduce degenerate behaviors, we further introduce **structural gates** that modulate the reward according to lightweight output constraints (e.g., a prescribed reasoning format and translation length). This design avoids neural reward models that require additional computation and can be vulnerable to reward hacking in long-horizon RL, and it also addresses our empirical finding that neural metrics fails to provide supervision for culturally grounded entity choices. With these verifiable rewards and an efficient critic-free policy optimization recipe, EA-RLVR established a stable RLVR training framework for cross-cultural

entity translation.

We conduct extensive experiments to evaluate EA-RLVR, yielding three key insights into its efficacy and underlying mechanisms: **(1) EA-RLVR incentivizes parametric knowledge.** Training on only 7k examples generalizes to a 50k test set whose entities are entirely unseen during training, improving entity translation accuracy by +8.21%–9.06% across different model scales. **(2) The learned strategy transfers beyond entity evaluation.** On WMT24pp, our models achieve improvements of XCOMET by +1.25–1.35 points, even though XCOMET is never used as supervision. When scaling training to the full dataset and extending optimization to 1,000 steps, the gains increase to +1.59–1.68 points. **(3) In-depth analyses clarify the dynamics of learning.**  $pass@k$  evaluation, neural reward comparison, cross-lingual generalization, and examinations of reward-hacking behavior all point to the same pattern: our method improves sampling efficiency and induces stable, cross-cultural translation strategies, rather than encouraging memorization.

Our contributions are as follows: **(1)** We propose a novel framework for cross-cultural machine translation based on RLVR, showing that entity translation can be improved without access to external databases by directly incentivizing context-appropriate entity choices. **(2)** We empirically demonstrate that this approach improves entity translation accuracy and general translation quality across languages and model scales, including settings involving entirely unseen entities. **(3)** We provide several analyses that explain the mechanism behind these improvements.

## 2 Related Work

**Cross-Cultural and Entity-Centric Machine Translation** Prior work addresses the challenge of culturally grounded entities largely through two avenues: external knowledge integration and targeted data augmentation. Retrieval-based methods explicitly ground translation in external sources, utilizing multilingual knowledge graphs (e.g., KG-MT; Conia et al., 2024) or document stores (e.g., RAGtrans; Wang et al., 2025b) to resolve entity ambiguities. While these approaches mitigate hallucinations, they introduce a dependency on the availability and quality of auxiliary databases. On the training side, recent works enhance entity robustness by synthesizing code-switched or entity-

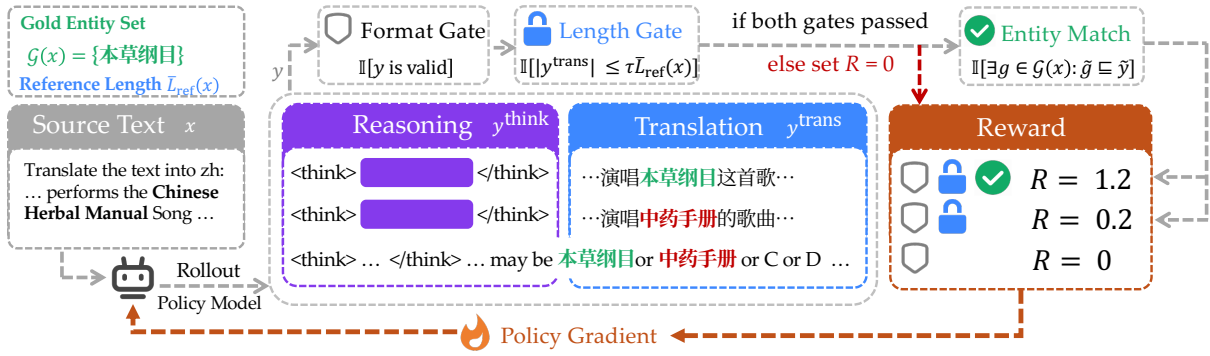


Figure 2: EA-RLVR framework.

replaced data for denoising pre-training (Hu et al., 2022; Liang et al., 2024), or by jointly optimizing translation with entity alignment tasks (Rikters and Miwa, 2024). Unlike these approaches, EA-RLVR does not require external retrieval at test time nor complex data synthesis pipelines. Instead, we cast entity translation as a reasoning problem, employing RLVR to activate and stabilize the parametric knowledge already present in the pre-trained model.

**RLVR and Reasoning in Translation** Recent post-training paradigms of LLMs leverage Reinforcement Learning with Verifiable Rewards (RLVR) to induce reasoning capabilities, a process theoretically understood as improving sampling efficiency to activate latent knowledge already present in the base model (Lambert et al., 2025; DeepSeek-AI et al., 2026; Yue et al., 2025). In machine translation, recent initiatives have actively explored integrating reasoning capabilities, for instance by employing multi-agent frameworks to synthesize long chain-of-thought trajectories for distillation (Wang et al., 2025a) or harnessing feedback from LLM judges and neural quality metrics to guide optimization (Feng et al., 2025a; Wang et al., 2025c; Feng et al., 2025b). Complementing these advances, EA-RLVR introduces a distinct paradigm centered on strict, rule-based verifiable rewards. We treat cultural entity translation as a precise reasoning task, employing deterministic rewards to directly surface parametric knowledge, thereby offering an alternative to distillation or neural-based objectives.

### 3 Method

We propose **EA-RLVR** (Entity-Anchored Reinforcement Learning with Verifiable Rewards), a framework designed to incentivize LLMs to accu-

rately ground cultural entities during translation without external knowledge. As illustrated in Figure 2, our approach treats cross-cultural translation as a sequential decision process optimized via reinforcement learning.

As shown in Figure 2, the framework consists of three core components: (1) A **reasoning-aware policy** that generates a thinking trajectory before the final translation, allowing the model to elicit latent knowledge; (2) A **verifiable reward mechanism** that anchors supervision on deterministic entity matching, safeguarded by structural gates to prevent reward hacking; and (3) A **critic-free optimization algorithm** that stabilizes training using sequence-level importance ratio. In the following sections, we detail the task formulation (§3.1), the reward design (§3.2), and the policy optimization objective (§3.3).

#### 3.1 Task Formulation

Given a source sentence  $x$ , we aim to generate a target-language translation  $y^{\text{trans}}$  that correctly renders the culturally grounded entity mention(s) in context. We treat an autoregressive LLM as a stochastic policy  $\pi_\theta$  over output tokens, i.e.,

$$\pi_\theta(y | x) = \prod_{t=1}^{|y|} \pi_\theta(y_t | x, y_{<t}). \quad 233$$

Generation induces an episodic decision process in which the state at step  $t$  is  $(x, y_{<t})$ , the action is the next token  $y_t$ , and the episode terminates when an  $\langle \text{eos} \rangle$  token is produced.

**Reasoning and translation segments.** Following recent reasoning-based post-training (DeepSeek-AI et al., 2026), the model is encouraged to produce a reasoning trace enclosed by  $\langle \text{think} \rangle$  and  $\langle \text{/think} \rangle$  before emitting the final

translation. When the response format is valid, we decompose the output  $y$  as

$$y = \langle \text{think} \rangle y^{\text{think}} \langle \text{/think} \rangle y^{\text{trans}},$$

where  $y^{\text{think}}$  contains deliberation and  $y^{\text{trans}}$  is the final translation. Since the reasoning portion may contain exploratory candidate entities rather than the model’s final decision, all content-based evaluation will be applied exclusively to  $y^{\text{trans}}$ .

### 3.2 Stable and Verifiable Reward Design for Cross-Cultural Entity Translation

Our main design goal is to construct a reward that is (i) *verifiable* from dataset annotations without a learned reward model, (ii) directly *aligned* with cross-cultural entity correctness, and (iii) *stable* under policy optimization and robust to reward hacking.

**Normalized entity matching.** Each example is annotated with a comprehensive set of acceptable target entities  $\mathcal{G}(x)$ , derived from the Wiki-data alias field. This set captures legitimate variations, minimizing false negatives where the model predicts a valid entity surface form that differs from the primary reference. For brevity we denote  $\tilde{y} = \text{norm}(y^{\text{trans}})$  and  $\tilde{g} = \text{norm}(g)$ , where  $\text{norm}(\cdot)$  lowercases text and removes diacritics. We define a deterministic match function using normalized substring matching:

$$m(y, \mathcal{G}(x)) = \mathbb{I}[\exists g \in \mathcal{G}(x) : \tilde{g} \sqsubseteq \tilde{y}], \quad (1)$$

where  $\mathbb{I}[\cdot]$  is the indicator function and  $a \sqsubseteq b$  denotes that  $a$  is a substring of  $b$ . This captures whether the model produces an appropriate target-language realization of the entity.

**Structural gates: format and length.** To ensure that rewards reflect meaningful translations rather than degenerate behaviors, we introduce two hard gates on the output. First, the response must follow the required  $\langle \text{think} \rangle / \langle \text{/think} \rangle$  structure. Second, the translation segment must remain within a reasonable length relative to the references. Formally, let  $g_{\text{fmt}}(y) \in \{0, 1\}$  indicate whether the format is valid, and define

$$g_{\text{len}}(x, y) = \mathbb{I}[|y^{\text{trans}}| \leq \tau \cdot \bar{L}_{\text{ref}}(x)], \quad (2)$$

where  $\bar{L}_{\text{ref}}(x)$  is the average reference length and  $\tau > 0$  controls tolerance. Only responses that satisfy *both* gates are eligible to receive any reward.

We empirically demonstrate the necessity of these constraints in Appendix A, showing that removing them leads to catastrophic length explosion and reward hacking via keyword enumeration.

**Final reward.** Combining these components, the terminal reward is

$$R(x, y) = g_{\text{fmt}}(y)g_{\text{len}}(x, y) \left( \alpha + m(y, \mathcal{G}(x)) \right). \quad (3)$$

Intuitively, a response receives no reward if it fails either structural gate. If both gates are satisfied, it obtains a base reward  $\alpha$  for producing a well-formed answer, and an additional bonus when the target entity is correctly realized. We set  $\alpha = 0.2$  and  $\tau = 2$  in all experiments.

### 3.3 Policy Optimization

Following recent advances in RL post-training for large language models, we optimize  $\pi_\theta$  using a clipped policy-gradient objective from the PPO family (Schulman et al., 2017). To reduce training cost and improve stability, we adopt a critic-free variant and incorporate group-normalized advantages, sequence-level importance ratios, and asymmetric clipping, building on GRPO (Shao et al., 2024), GSPO (Zheng et al., 2025), and DAPO (Liu et al., 2025).

**Objective.** Given an input  $x$ , we sample  $G$  candidate responses  $\{y_i\}_{i=1}^G$  from the old policy  $\pi_{\theta_{\text{old}}}$ . The policy parameters are updated by maximizing the clipped surrogate:

$$\mathcal{J}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, \{y_i\} \sim \pi_{\theta_{\text{old}}}(\cdot|x)} \left[ \frac{1}{G} \sum_{i=1}^G \min \left( s_i(\theta) \hat{A}_i, \text{clip}(s_i(\theta), 1 - \varepsilon_{\text{low}}, 1 + \varepsilon_{\text{high}}) \hat{A}_i \right) \right], \quad (4)$$

where  $\mathcal{D}$  denotes the training dataset containing source sentences  $x$ , the  $\varepsilon_{\text{low}}$  and  $\varepsilon_{\text{high}}$  are clipping thresholds that bound the policy update to prevent training instability, and the  $s_i(\theta)$  is the sequence-level importance ratio defined in Eq. (6).

**Group-normalized advantages.** Rather than relying on a learned critic, we compute an advantage for each sampled response relative to other responses from the same prompt:

$$\hat{A}_i = \frac{R(x, y_i) - \text{mean}(\{R(x, y_j)\}_{j=1}^G)}{\text{std}(\{R(x, y_j)\}_{j=1}^G)}, \quad (5)$$

where  $\text{mean}(\cdot)$  and  $\text{std}(\cdot)$  denote the sample mean and standard deviation within the group. This normalization stabilizes training and makes the reward scale largely irrelevant.

**Sequence-level importance ratios.** We compute the importance ratio at the sequence level with length normalization:

$$s_i(\theta) = \left( \frac{\pi_\theta(y_i | x)}{\pi_{\theta_{\text{old}}}(y_i | x)} \right)^{\frac{1}{|y_i|}}$$

$$= \exp \left( \frac{1}{|y_i|} \sum_{t=1}^{|y_i|} \log \frac{\pi_\theta(y_{i,t} | x, y_{i,<t})}{\pi_{\theta_{\text{old}}}(y_{i,t} | x, y_{i,<t})} \right). \quad (6)$$

This formulation discourages overly aggressive updates on long sequences while still allowing meaningful policy shifts when rewards are consistently better.

## 4 Experiments

Our experiments aim to verify two hypotheses: (1) that RE-RLVR training can effectively elicit latent cultural knowledge solely from the model’s pre-trained parameters, and (2) that this entity-centric optimization does not compromise general translation quality. After outlining our setup in §4.1, we present empirical evidence supporting the activation of parametric knowledge in §4.2 and demonstrate positive transfer effects to general translation in §4.3.

### 4.1 Experimental Setup

**Datasets and Benchmarks.** To evaluate the activation of cultural knowledge, we utilized **XC-Translate** (Conia et al., 2024), a benchmark specializing in cross-cultural entity translation. We trained our models using 7,278 examples for training and the official test set of 49,606 examples. Each sample is annotated with a list of gold entity aliases derived from Wikidata, which serves as the reference set for our verifiable reward. Crucially, **the training and test sets share no overlapping entities**. The dataset covers ten language pairs (English  $\rightarrow$  X), detailed in Table 1. For general translation capability, we evaluated on **WMT24++** (Deutsch et al., 2025) across the corresponding languages, using the official test sets without any domain-specific fine-tuning. Further details on data composition are provided in Appendix E.

**Models and Baselines.** We employed Qwen3-8B and Qwen3-14B (Yang et al., 2025) as our backbone models, which are pre-trained on 36T tokens and possess native reasoning capabilities. Our proposed EA-RLVR was compared against two primary internal baselines: (1) the base model

and (2) Supervised Fine-Tuning (SFT) on the same 7k examples to control for data exposure. To further contextualize our performance, we included several strong external baselines: (i) frontier proprietary LLMs: GPT-5-mini; (ii) strong open-source model: Qwen3-235B-A22B, Marco-o1 (Zhao et al., 2024a), a multilingual reasoning model, and DeepTrans-7B (Wang et al., 2025a), a specialized reasoning-based translation model. More evaluation details are in Appendix E.

**Training and Implementation.** We implemented EA-RLVR using the ver1 framework (Sheng et al., 2024). All models were trained using the policy optimization algorithm described in §3.3. Full implementation details are provided in Appendix E.

**Evaluation Metrics.** We report three primary metrics: (1) **Entity Translation Accuracy:** Consistent with the normalized substring matching defined in Eq. 1, Entity Translation Accuracy measures the percentage of test samples where the generated translation  $y$  successfully includes the correct cultural entity (i.e.,  $m(y, \mathcal{G}(x)) = 1$ ). (2) **chrF:** We report the sentence-level Character F-score (chrF) (Popović, 2015) to assess the overall quality of the generated translations. (3) **XCOMET-XL** (Guerreiro et al., 2024): A state-of-the-art reference-based neural metric used to assess the general quality and fluency of the translations on WMT24++.

### 4.2 EA-RLVR incentivizes parametric knowledge

Table 1 reports the entity translation accuracy on the XC-Translate test set. The results provide empirical support for our hypothesis regarding parametric knowledge activation, revealing a fundamental divergence in effectiveness between imitation-based and reasoning-based optimization.

**Breaking the Imitation Ceiling via Reasoning.** Standard SFT yields negligible gains (e.g., +0.22% for Qwen3-14B), despite using the same 7k data as EA-RLVR. This outcome is predictable as our train-test entities are disjoint. Unlike style transfer, where learning generalizable patterns suffices, entity translation inherently biases towards memorization, limiting the effectiveness of SFT. EA-RLVR, however, reframes this task as a reasoning problem. Consequently, it achieves substantial improvements (+8.21%), **enabling the 14B model (31.87%) to**

Table 1: Entity translation accuracy (%) on XC-Translate across ten language directions (en  $\rightarrow$  X). Train and test sets share no overlapping entities. RLVR consistently outperforms SFT under the same data budget, while also maintaining strong character-level faithfulness.

Model	Entity Translation Accuracy on XC-Translate ( $en \rightarrow X$ )										Avg.
	ar	de	es	fr	it	ja	ko	th	tr	zh	Acc/charF
<i>Baselines</i>											
GPT-5-mini	35.03	36.03	42.71	35.85	35.37	37.27	30.20	15.96	40.44	29.85	33.87/62.30
Qwen3-235B-A22B	27.93	32.06	41.79	34.29	33.07	29.07	30.22	15.32	34.32	32.38	31.05/61.99
Marco-o1	11.88	19.72	26.47	22.09	21.26	11.77	8.42	3.45	17.71	16.23	15.90/49.68
DeepTrans-7B	11.41	21.58	30.10	23.49	22.44	13.21	8.44	2.90	17.37	15.96	16.69/47.88
<i>Ours</i>											
Qwen3-8B	14.91	23.64	31.04	25.97	25.01	17.13	11.12	5.63	23.94	23.50	20.19/56.07
+ SFT	15.28	23.28	30.76	25.32	25.44	16.86	11.12	5.86	23.56	23.08	20.06/56.20
+ EA-RLVR	<b>25.23</b>	<b>31.01</b>	<b>44.44</b>	<b>34.89</b>	<b>35.35</b>	<b>24.02</b>	<b>25.70</b>	<b>14.74</b>	<b>27.79</b>	<b>29.31</b>	<b>29.25/59.86</b>
Qwen3-14B	20.48	26.86	34.88	28.66	28.70	18.81	17.67	7.98	25.96	26.57	23.66/58.22
+ SFT	20.21	26.97	35.13	28.75	30.05	18.32	17.65	8.01	25.93	27.75	23.88/59.43
+ EA-RLVR	<b>28.17</b>	<b>33.71</b>	<b>44.70</b>	<b>35.10</b>	<b>37.62</b>	<b>27.64</b>	<b>31.29</b>	<b>17.09</b>	<b>29.96</b>	<b>33.42</b>	<b>31.87/62.27</b>

424 **outperform the much larger Qwen3-235B-A22B**  
425 **baseline (31.05 %)**. RLVR never presents the gold  
426 entity to the model during training, therefore its  
427 performance gains cannot stem from memorizing  
428 entity mappings. By optimizing self-generated tra-  
429 jectories against verifiable rewards, the model is  
430 incentivized to internalize a better strategy rather  
431 than fitting specific entity mappings. As detailed  
432 in our qualitative analysis (Appendix G), this be-  
433 havior manifests as explicitly grounding cultural  
434 contexts within the reasoning trace before generat-  
435 ing the translation. This learned capability further  
436 exhibits robustness across typologically diverse lan-  
437 guage groups (Appendix C), confirming that the  
438 performance gains stem from the activation of lat-  
439 ent parametric knowledge.

### 440 4.3 Transfer to General Machine Translation

441 A potential concern with specialized reinforcement  
442 learning is the risk of “alignment tax,” where opti-  
443 mizing for a narrow objective (entity correctness)  
444 degrades general capabilities. We investigate this  
445 on the WMT24++ benchmark (Table 2) and ob-  
446 serve the opposite effect.

447 **Reasoning Improves General Quality.** Despite  
448 being trained solely on the XC-Translate dataset  
449 (which contains only 7k samples), EA-RLVR mod-  
450 els consistently improve general translation qual-  
451 ity across all evaluated languages. Qwen3-14B +  
452 EA-RLVR achieves an average XCOMET score of  
453 78.41, a +1.35 point improvement over the base  
454 model. This suggests that the reasoning strategies  
455 learned for entity translation—such as attending  
456 more carefully to source semantics and deliberat-

ing before generating—are transferable. The model  
becomes less prone to hallucination and more faith-  
ful to the source text, benefiting general translation  
tasks.

461 **Scalability and Robustness.** Table 2 also com-  
462 pares models trained on the standard 7k set versus  
463 the full dataset. While SFT performance stagnates  
464 or even slightly degrades when scaling data (likely  
465 due to overfitting on the specific formatting of the  
466 entity dataset), EA-RLVR continues to improve.  
467 The “Full Data” setting yields further gains, push-  
468 ing the average XCOMET score to 78.65 for the  
469 14B model. This indicates that our outcome-based  
470 reward formulation provides a stable optimization  
471 landscape that scales effectively with data, unlike  
472 likelihood-based SFT which may suffer from dis-  
473 tribution shift.

## 474 5 Analysis

### 475 5.1 Improved Sampling Efficiency: Unlocking 476 Dormant Knowledge

477 To analyze the mechanism behind the performance  
478 gains, we adopt the  $pass@k$  evaluation framework  
479 recently utilized to study the boundaries of RLVR  
480 in reasoning tasks (Yue et al., 2025). Formally,  
481  $pass@k$  estimates the probability that at least one  
482 correct translation exists within  $k$  independent  
483 samples generated for a given input. Following (Chen  
484 et al., 2021), we calculate the unbiased estimator  
485 (see Appendix F for details). By observing how this  
486 probability scales with  $k$ , we can distinguish be-  
487 tween *knowledge injection* (learning new informa-  
488 tion) and *knowledge activation* (surfacing existing  
489 information). Figure 3 compares the entity transla-

Table 2: XCOMET-XL score on WMT24++ across ten language directions. All models are trained only on the XC-Translate cross-cultural entity dataset, without using WMT data or general MT supervision. Despite this, EA-RLVR consistently improves performance on general machine translation, and further benefits appear when scaling to the full XC-Translate dataset.

Model	XCOMET score on WMT24++ ( $en \rightarrow X$ )										Avg.
	ar	de	es	fr	it	ja	ko	th	tr	zh	
<i>Baselines</i>											
GPT-5-mini	72.86	91.71	86.93	84.30	86.34	82.15	82.56	80.29	79.49	77.76	82.44
Qwen3-235B-A22B	69.72	90.42	85.53	81.93	84.59	78.54	79.76	77.32	73.92	75.88	79.76
Marco-o1	53.50	84.09	80.23	75.38	74.84	67.04	67.04	67.04	48.15	70.47	68.78
DeepTrans-7B	56.51	84.69	81.45	76.31	74.86	70.24	62.94	65.80	48.78	71.53	69.31
<i>Ours</i>											
Qwen3-8B	61.88	88.06	82.36	78.33	81.09	72.86	71.97	72.43	62.48	73.11	74.46
+ SFT	62.73	88.29	82.91	78.52	80.23	72.99	71.73	72.20	62.47	73.67	74.57
+ EA-RLVR	<b>64.65</b>	<b>88.85</b>	<b>83.29</b>	<b>79.48</b>	<b>81.16</b>	<b>73.60</b>	<b>74.04</b>	<b>72.92</b>	<b>64.91</b>	<b>74.20</b>	<b>75.71</b>
+ SFT (full data)	62.51	87.78	82.25	78.79	80.48	72.32	71.37	71.57	61.83	73.31	74.22
+ EA-RLVR (full data)	<b>65.29</b>	<b>89.00</b>	<b>83.48</b>	<b>79.82</b>	<b>82.17</b>	<b>74.52</b>	<b>73.62</b>	<b>73.88</b>	<b>65.58</b>	<b>74.06</b>	<b>76.14</b>
Qwen3-14B	66.37	89.29	83.79	80.15	81.91	76.02	75.73	74.87	67.40	75.10	77.06
+ SFT	66.20	89.12	84.10	80.09	82.37	76.38	75.35	75.23	67.70	75.15	77.17
+ EA-RLVR	<b>68.30</b>	<b>89.79</b>	<b>84.28</b>	<b>80.78</b>	<b>83.56</b>	<b>77.61</b>	<b>77.49</b>	<b>76.30</b>	<b>70.16</b>	<b>75.84</b>	<b>78.41</b>
+ SFT (full data)	65.27	89.11	84.07	79.87	82.47	76.26	75.50	74.92	67.82	75.18	77.05
+ EA-RLVR (full data)	<b>68.00</b>	<b>90.11</b>	<b>85.06</b>	<b>81.27</b>	<b>84.10</b>	<b>78.14</b>	<b>77.57</b>	<b>75.97</b>	<b>70.25</b>	<b>75.99</b>	<b>78.65</b>

tion accuracy of Qwen3-8B (Base) and EA-RLVR across varying sample sizes  $k \in [1, 128]$ .

We observe a distinct convergence pattern across most languages. At  $k = 1$ , EA-RLVR holds a substantial lead over the base model, confirming that our policy optimization successfully concentrates probability mass on the correct entity translations. However, as  $k$  increases, the base model’s accuracy rises steeply, often converging with or appearing to surpass the RLVR model at  $k = 128$ . This phenomenon has a critical implication: **the base model inherently possesses the necessary cultural knowledge** to translate these entities correctly (evidenced by high performance at large  $k$ ), but it fails to rank these correct translations as the most probable candidates during standard decoding. EA-RLVR functions as a steering mechanism that activates this dormant knowledge, transforming low-probability correct candidates into high-probability deterministic outputs, rather than memorizing new mappings from external supervision. We further isolate the contribution of the explicit reasoning phase in Appendix B, finding that the “thinking” workspace is essential for absorbing the complexity of cultural alignment without sacrificing general fluency.

**Sampling Efficiency and Determinism.** The slope of the curves in Figure 3 further elucidates the shift in model behavior. The base model exhibits a high-entropy distribution over entities, requiring

extensive sampling ( $k \gg 1$ ) to uncover the correct answer. In contrast, the RLVR curves are notably flatter, indicating a more deterministic policy where the model is confident in its reasoning path. While high determinism theoretically reduces diversity (explaining the slight underperformance at  $k = 128$  in some high-resource languages where the base model’s broad search space is advantageous), it is the desired behavior for a translation system: users expect the correct cultural translation in a single attempt ( $k = 1$ ), not after filtering through a hundred generations. EA-RLVR effectively optimizes for this *sampling efficiency*.

## 5.2 The Fluency Trap: Why Neural Rewards Fail Cultural Entities

A natural alternative to our rule-based framework is to optimize state-of-the-art neural quality metrics directly. To investigate this, we trained a **Comet-RL** baseline on Qwen3-8B using the same RL setup but replacing our normalized entity matching reward with sentence-level comet scores, specifically the wmt22-comet-da (Rei et al., 2022). Figure 4 presents a striking divergence in optimization outcomes, revealing what we term the “*Fluency Trap*”:

**High Fluency, Low Grounding.** As shown in the right panel, the Comet-RL model achieves substantial gains in general translation quality, boosting the XCOMET score on WMT24++ from 74.46

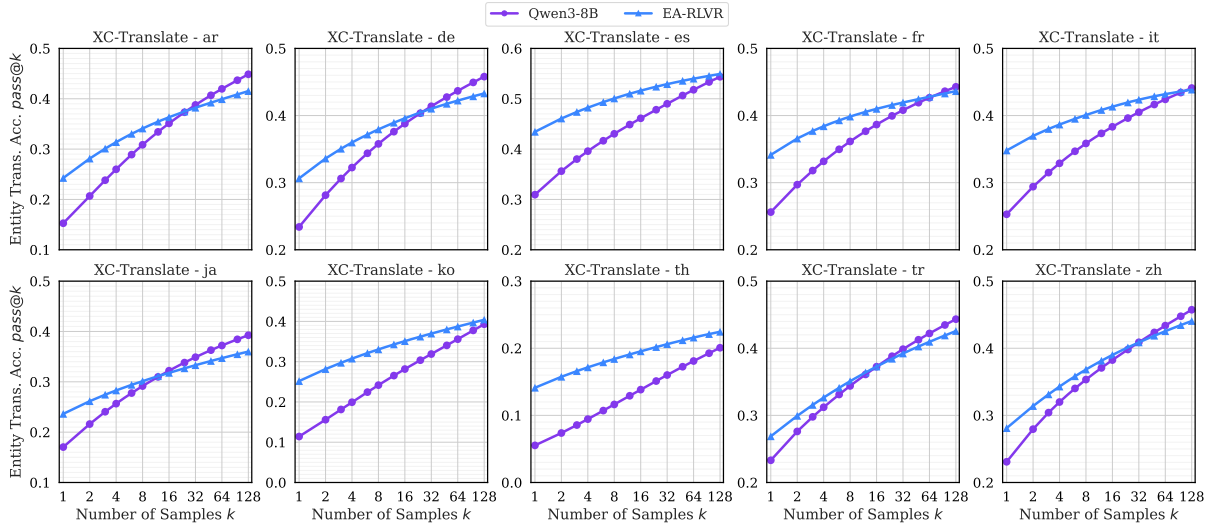


Figure 3: Entity Translation Accuracy  $pass@k$  curves across ten languages. The Base model (purple) shows poor performance at  $k = 1$  but improves rapidly as  $k$  increases, indicating latent knowledge is present but buried. EA-RLVR (blue) significantly boosts  $pass@1$  accuracy, effectively surfacing this parametric knowledge. The convergence at high  $k$  confirms that improvements stem from better utilization of existing knowledge rather than learning new facts.

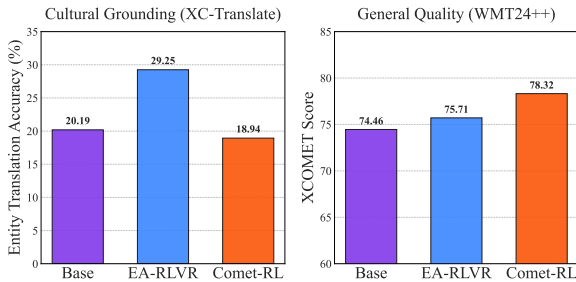


Figure 4: **Verifiable vs. Neural Rewards.** While Comet-RL maximizes general fluency (Right) at the cost of cultural accuracy (Left), falling into a “fluency trap”, EA-RLVR achieves robust improvements across both cultural grounding and general translation quality.

to 78.32. This confirms that RL effectively optimized the reward signal. However, the left panel reveals a critical failure: on the entity-dense XC-Translate benchmark, Comet-RL’s entity accuracy actually *degrades* from 20.19% (Base) to 18.94%. The neural metric, lacking fine-grained resolution for specific entities, fails to penalize these fluent but culturally incorrect entity errors, which echos prior observations (Rei et al., 2023).

**Verifiable Rewards Ensure Alignment.** In contrast, EA-RLVR escapes this trap. By anchoring supervision on verifiable outcomes, it forces the model to prioritize semantic correctness. This yields a massive improvement in entity accuracy (+9.06% absolute) while still conferring ro-

bust gains in general translation quality (+1.25 XCOMET). This result fundamentally justifies our approach: while holistic neural metrics drive fluency, verifiable constraints are indispensable for aligning models with entity translation tasks.

## 6 Conclusion

In this work, we investigate the underutilized potential of parametric knowledge in cross-cultural translation: while LLMs possess extensive latent cultural knowledge ( $pass@128$  performance), they struggle to utilize it during standard decoding ( $pass@1$ ). Motivated by this observation, we propose **EA-RLVR**, a framework that transforms cross-cultural translation from a memorizing task into a reasoning-intensive process. Our extensive experiments reveals that incentivizing verifiable correctness is superior to optimizing neural quality metrics, which often lead models into a “fluency trap”, generating smooth but culturally inaccurate entities. By anchoring supervision on entity matching and allowing the model to reasoning, EA-RLVR enables a 14B model to outperform a 235B baseline significantly on unseen entities. Ultimately, our findings suggest a potential paradigm shift for knowledge-intensive translation: moving beyond mere imitation of references (SFT) or reliance on external retrieval, toward internalizing self-evolving strategies that effectively unlock the model’s inherent potential.

## 7 Limitations

**Gap to Latent Potential.** While EA-RLVR significantly outperforms SFT, a notable disparity remains between the optimized policy’s greedy performance ( $pass@1$ ) and the model’s theoretical upper bound estimated by rejection sampling ( $pass@128$ , as shown in Figure 1). This gap highlights a shared limitation across current RLVR methodologies: standard policy gradient algorithms often struggle to fully explore and converge to the global optimum within a limited sample budget. Future research could bridge this gap by developing more sample-efficient optimization algorithms or by scaling the number of rollout trajectories ( $G$ ). Although computationally intensive, expanding the exploration horizon offers a promising avenue for approaching the model’s intrinsic capability ceiling.

**Knowledge Boundaries.** Our framework is designed for *knowledge elicitation*, not *knowledge injection*. EA-RLVR optimizes the retrieval of long-tail cultural concepts that exist within the model’s pre-training data but are suppressed during standard decoding. Consequently, it cannot generate correct translations for entities entirely absent from the pre-training corpus. However, we observed that our method synergizes effectively with retrieval-augmented systems rather than acting as a simple alternative, providing a combined benefit that we evaluate in Appendix D.

**Reward Rigidity vs. Flexibility.** We prioritize optimization stability via rigid substring matching to prevent the reward hacking often observed with neural metrics. Although we mitigate the risk of false negatives by employing a comprehensive gold set of aliases (derived from Wikidata) rather than a single reference, this strict verification process may still occasionally penalize valid but unlisted stylistic variations. Developing rewards that balance verifiable strictness with semantic flexibility remains an open challenge for the field.

## References

Sweta Agrawal, Chunting Zhou, Mike Lewis, Luke Zettlemoyer, and Marjan Ghazvininejad. 2023. [In-context examples selection for machine translation](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 8857–8873, Toronto, Canada. Association for Computational Linguistics.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, and 39 others. 2021. [Evaluating large language models trained on code](#). *Preprint*, arXiv:2107.03374.

Simone Conia, Daniel Lee, Min Li, Umar Farooq Minhas, Saloni Potdar, and Yunyao Li. 2024. [Towards cross-cultural machine translation with retrieval-augmented generation from multilingual knowledge graphs](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 16343–16360.

Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2022. [Flashattention: Fast and memory-efficient exact attention with io-awareness](#). *Preprint*, arXiv:2205.14135.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang, Shirong Ma, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, and 181 others. 2026. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#). *Preprint*, arXiv:2501.12948.

Daniel Deutsch, Eleftheria Briakou, Isaac Rayburn Caswell, Mara Finkelstein, Rebecca Galor, Juraj Juraska, Geza Kovacs, Alison Lui, Ricardo Rei, Jason Riesa, and 1 others. 2025. [Wmt24++: Expanding the language coverage of wmt24 to 55 languages & dialects](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 12257–12284.

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

Zhaopeng Feng, Jiahao Ren, Jiayuan Su, Jiamei Zheng, Hongwei Wang, and Zuozhu Liu. 2025b. [Mt-rewardtree: A comprehensive framework for advancing llm-based machine translation via reward modeling](#). *Preprint*, arXiv:2503.12123.

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

Junjie Hu, Hiroaki Hayashi, Kyunghyun Cho, and Graham Neubig. 2022. [DEEP: DENOISING ENTITY PRE-TRAINING FOR NEURAL MACHINE TRANSLATION](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1753–1766, Dublin, Ireland. Association for Computational Linguistics.

698	Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2021. <a href="#">Unsupervised dense information retrieval with contrastive learning</a> .	
699		
700		
701		
702	Urvashi Khandelwal, Angela Fan, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. Nearest neighbor machine translation. In <i>International Conference on Learning Representations</i> .	
703		
704		
705		
706	Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, and 4 others. 2025. <a href="#">Tulu 3: Pushing frontiers in open language model post-training</a> . <i>Preprint</i> , arXiv:2411.15124.	
707		
708		
709		
710		
711		
712		
713		
714		
715	Tian Liang, Xing Wang, Mingming Yang, Yujiu Yang, Shuming Shi, and Zhaopeng Tu. 2024. <a href="#">Addressing entity translation problem via translation difficulty and context diversity</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 11628–11638, Bangkok, Thailand. Association for Computational Linguistics.	
716		
717		
718		
719		
720		
721		
722	Jiacai Liu, Chaojie Wang, Chris Yuhao Liu, Liang Zeng, Rui Yan, Yiwen Sun, and Yang Liu. 2025. <a href="#">DAPO: Improving multi-step reasoning abilities of large language models with direct advantage-based policy optimization</a> . In <i>The Thirty-ninth Annual Conference on Neural Information Processing Systems</i> .	
723		
724		
725		
726		
727		
728	Nikita Moghe, Arnisa Fazla, Chantal Amrhein, Tom Kocmi, Mark Steedman, Alexandra Birch, Rico Senrich, and Liane Guillou. 2025. <a href="#">Machine translation meta evaluation through translation accuracy challenge sets</a> . <i>Computational Linguistics</i> , 51(1):73–137.	
729		
730		
731		
732		
733	Maja Popović. 2015. <a href="#">chrF: character n-gram F-score for automatic MT evaluation</a> . In <i>Proceedings of the Tenth Workshop on Statistical Machine Translation</i> , pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.	
734		
735		
736		
737		
738	Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, and 25 others. 2025. <a href="#">Qwen2.5 technical report</a> . <i>Preprint</i> , arXiv:2412.15115.	
739		
740		
741		
742		
743		
744		
745	Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022. <a href="#">COMET-22: Unbabel-IST 2022 submission for the metrics shared task</a> . In <i>Proceedings of the Seventh Conference on Machine Translation (WMT)</i> , pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.	
746		
747		
748		
749		
750		
751		
752		
753	Ricardo Rei, Nuno M. Guerreiro, Marcos Treviso, Luisa Coheur, Alon Lavie, and André Martins. 2023. <a href="#">The</a>	
754		
	<a href="#">inside story: Towards better understanding of machine translation neural evaluation metrics</a> . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 1089–1105, Toronto, Canada. Association for Computational Linguistics.	755
		756
		757
		758
		759
		760
	Matiss Rikters and Makoto Miwa. 2024. <a href="#">Entity-aware multi-task training helps rare word machine translation</a> . In <i>Proceedings of the 17th International Natural Language Generation Conference</i> , pages 47–54, Tokyo, Japan. Association for Computational Linguistics.	761
		762
		763
		764
		765
		766
	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. <a href="#">Proximal policy optimization algorithms</a> . <i>Preprint</i> , arXiv:1707.06347.	767
		768
		769
		770
	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. <a href="#">Deepseekmath: Pushing the limits of mathematical reasoning in open language models</a> . <i>Preprint</i> , arXiv:2402.03300.	771
		772
		773
		774
		775
		776
	Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. 2024. <a href="#">Hybridflow: A flexible and efficient rlhf framework</a> . <i>arXiv preprint arXiv:2409.19256</i> .	777
		778
		779
		780
		781
	Jiaan Wang, Fandong Meng, Yunlong Liang, and Jie Zhou. 2025a. <a href="#">Drt: Deep reasoning translation via long chain-of-thought</a> . <i>Preprint</i> , arXiv:2412.17498.	782
		783
		784
	Jiaan Wang, Fandong Meng, Yingxue Zhang, and Jie Zhou. 2025b. <a href="#">Retrieval-augmented machine translation with unstructured knowledge</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2025</i> , pages 5858–5871, Suzhou, China. Association for Computational Linguistics.	785
		786
		787
		788
		789
		790
	Jiaan Wang, Fandong Meng, and Jie Zhou. 2025c. <a href="#">Deeptrans: Deep reasoning translation via reinforcement learning</a> . <i>Preprint</i> , arXiv:2504.10187.	791
		792
		793
	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. <a href="#">Qwen3 technical report</a> . <i>arXiv preprint arXiv:2505.09388</i> .	794
		795
		796
		797
		798
	Binwei Yao, Ming Jiang, Tara Bobinac, Diyi Yang, and Junjie Hu. 2024. <a href="#">Benchmarking machine translation with cultural awareness</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 13078–13096, Miami, Florida, USA. Association for Computational Linguistics.	799
		800
		801
		802
		803
		804
	Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Yang Yue, Shiji Song, and Gao Huang. 2025. <a href="#">Does reinforcement learning really incentivize reasoning capacity in llms beyond the base model?</a> <i>Preprint</i> , arXiv:2504.13837.	805
		806
		807
		808
		809

810 Yu Zhao, Huifeng Yin, Bo Zeng, Hao Wang, Tianqi  
811 Shi, Chenyang Lyu, Longyue Wang, Weihua Luo,  
812 and Kaifu Zhang. 2024a. [Marco-o1: Towards open](#)  
813 [reasoning models for open-ended solutions](#). *Preprint*,  
814 arXiv:2411.14405.

815 Yuze Zhao, Jintao Huang, Jinghan Hu, Xingjun Wang,  
816 Yunlin Mao, Daoze Zhang, Zeyinzi Jiang, Zhikai Wu,  
817 Baole Ai, Ang Wang, Wenmeng Zhou, and Yingda  
818 Chen. 2024b. [Swift:a scalable lightweight infrastruc-](#)  
819 [ture for fine-tuning](#). *Preprint*, arXiv:2408.05517.

820 Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui  
821 Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong  
822 Liu, Rui Men, An Yang, Jingren Zhou, and Jun-  
823 yang Lin. 2025. [Group sequence policy optimization](#).  
824 *Preprint*, arXiv:2507.18071.

## 825 A Impact of Structural Gates

826 To validate the necessity of the structural con-  
827 straints introduced in Section 3.2, we conduct an  
828 ablation study using the full XC-Translate dataset  
829 with extended optimization (1,000 steps). We iso-  
830 late the contributions of the format and length gates  
831 by comparing three configurations:

- 832 • **Soft Format:** A relaxed baseline where any  
833 output containing a valid `<think>` block is  
834 eligible for rewards, with no length penalty  
835 applied.
- 836 • **Format Gate Only:** The strict format verifi-  
837 cation ( $g_{\text{fmt}}$ ) is applied, but the relative length  
838 constraint ( $g_{\text{len}}$ ) is removed.
- 839 • **EA-RLVR:** Our proposed framework, which  
840 enforces both strict format verification and the  
841 relative length constraint ( $g_{\text{len}}$ ).

842 **Training Dynamics and Stability.** The training  
843 dynamics, visualized in Figure 5, demonstrate that  
844 structural constraints are critical for optimization  
845 stability. The most prominent difference lies in the  
846 *Response Length* (Center). In the absence of the  
847 length constraint ( $g_{\text{len}}$ ), both ablated variants suf-  
848 fer from a catastrophic “length explosion,” where  
849 generation length increases uncontrollably. This  
850 instability is mirrored in the *Actor Entropy* (Right),  
851 where the ablated models exhibit erratic spikes, in-  
852 dicating that the policy fails to converge to a stable  
853 reasoning strategy. In contrast, EA-RLVR main-  
854 tains a consistent length and stable entropy profile  
855 throughout training.

**Reward Hacking via Enumeration.** Qualitative  
analysis reveals that the length explosion is a symp-  
tom of reward hacking. As illustrated in the case  
study (Figure 6), without the length penalty, the  
optimization landscape encourages a degenerate  
solution: the model learns to “brute-force” the ver-  
ification condition ( $m(y, \mathcal{G}(x))$ ) by enumerating  
synonymous entities or repeating candidates. This  
strategy maximizes the recall of the gold entity at  
the expense of precision and structural integrity. By  
treating the translation task as a keyword-stuffing  
exercise, the model achieves technically high re-  
wards but produces unusable translations.

**The Deception of High Rewards.** The *Average*  
*Reward* curves (Left) present a counter-intuitive  
trend: the weaker constraints yield higher raw re-  
ward values. The *Soft Format* setting achieves  
the highest reward trajectory despite exhibiting  
the earliest collapse in generation quality. This  
phenomenon is a classic manifestation of **Good-**  
**hart’s Law:** when the unconstrained entity-match  
metric becomes the sole target, the model exploits  
its loopholes (e.g., infinite generation) rather than  
improving the intended task utility. EA-RLVR’s  
lower reward curve reflects a constrained, harder-  
to-optimize landscape that successfully steers the  
model away from these degenerate local optima  
and toward concise, correct translations.

## 884 B Impact of Reasoning: Does Thinking 885 Matter?

886 A central premise of EA-RLVR is that a dedicated  
887 reasoning phase (“thinking”) allows the model to  
888 navigate the optimization landscape more effec-  
889 tively than immediate generation. To isolate the  
890 impact of this reasoning process, we conduct a  
891 controlled comparison using the Qwen3-4B-2507  
892 family.

**Experimental Control.** We compare two  
specific checkpoints: the standard instruction-  
tuned model, Qwen3-4B-Instruct-2507,  
and the reasoning-enhanced model,  
Qwen3-4B-Thinking-2507. We apply the  
EA-RLVR framework to both models with a  
crucial adaptation for the Instruct baseline: since  
Qwen3-4B-Instruct-2507 does not generate  
reasoning traces (i.e., no `<think>` tokens), we  
remove the format verification gate ( $g_{\text{fmt}}$ ) from  
the reward function. However, to ensure a fair  
comparison of supervision signals, we retain

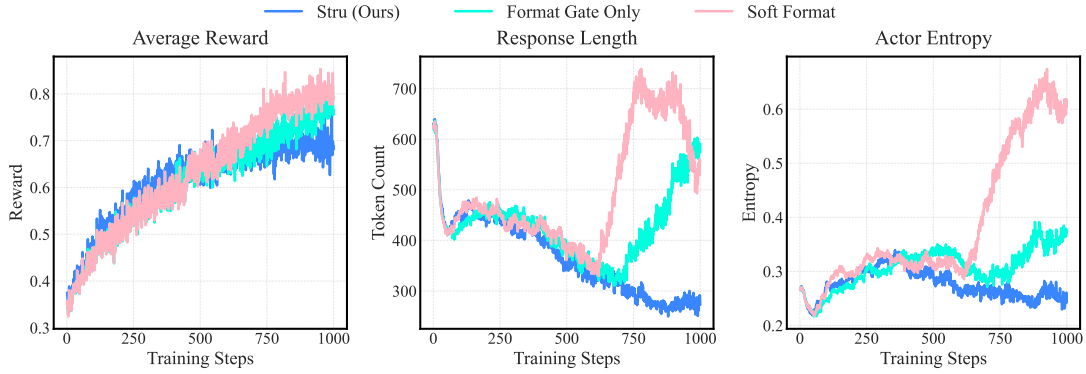


Figure 5: **Ablation Dynamics.** Training curves for Average Reward, Response Length, and Actor Entropy. Without the full structural gates (EA-RLVR, Blue), the model suffers from reward hacking, characterized by an explosion in response length (Center) and unstable entropy (Right), despite achieving higher raw rewards (Left).

Table 3: Cross-lingual generalization on XC-Translate with Qwen3-8B. The model is trained on one group of languages and evaluated on both. **Bold text** indicates zero-shot cross-lingual transfer (e.g., Train A  $\rightarrow$  Test B), while **gray text** indicates in-domain performance. Improvements over the Base model on unseen language groups demonstrate the acquisition of a transferable reasoning strategy.

Train split (Qwen3-8B)	Group A					Group B					Group A Avg.	Group B Avg.
	ar	ja	ko	th	zh	de	es	fr	it	tr		
Base	14.91	17.13	11.12	5.63	23.50	23.64	31.04	25.97	25.01	23.94	14.46	25.92
Group A	48.27	43.48	47.70	25.76	53.24	<b>29.70</b>	<b>40.61</b>	<b>33.78</b>	<b>34.95</b>	<b>26.89</b>	43.69	<b>33.19</b>
Group B	<b>24.81</b>	<b>21.77</b>	<b>18.32</b>	<b>8.62</b>	<b>30.20</b>	51.28	62.50	49.66	51.61	50.19	<b>20.74</b>	53.05

both the length constraint ( $g_{len}$ ) and the entity matching reward ( $m(y, \mathcal{G}(x))$ ). This setup allows us to strictly evaluate whether the presence of a thinking process facilitates better alignment with the verifiable reward.

### Thinking Unlocks Higher Entity Accuracy.

Figure 7 (Center) illustrates the progression of entity translation accuracy on the XC-Translate test set. The *Non-Think Model* (Cyan) plateaus early at approximately 21% accuracy. In contrast, the *Think Model* (Blue) achieves a significantly higher peak of  $\sim 26\%$ , despite starting from a lower baseline. This confirms that the reasoning trace provides the necessary computational workspace to resolve complex cultural entity mappings that are inaccessible to a single-pass decoding policy.

**Dynamics of Length and Stability.** The training dynamics reveal a crucial interaction between reasoning length and reward. As shown in Figure 7 (Left), the Think Model initially receives near-zero rewards. This is an artifact of the base model’s instability: the untrained Qwen3-4B-Thinking-2507 frequently generates extremely long chain-of-thought traces that exceed our training context window of 4096 tokens, causing the samples to be truncated and penalized.

However, EA-RLVR rapidly corrects this behavior. The *Response Length* curve (Right) shows a dramatic reduction in token count within the first 50 steps. The model learns to be concise, condensing its reasoning into an efficient path that fits the constraints while maximizing the entity-match reward. This demonstrates that EA-RLVR serves not only as a task optimizer but also as a length regularizer for reasoning models.

### Thinking Mitigates the “Alignment Tax”.

We further evaluate the impact of these strategies on general translation quality using WMT24++. Table 4 presents the XCOMET scores. A striking divergence is observed:

- **Standard Model (Instruct):** Applying EA-RLVR to Qwen3-4B-Instruct-2507 leads to a slight regression in general quality (Avg. 89.11  $\rightarrow$  88.46). Without a reasoning buffer, the model is forced to overload its generation weights to satisfy the strict entity constraints, leading to a “fluency tax” where general translation quality is sacrificed.
- **Reasoning Model (Thinking):** Conversely, Qwen3-4B-Thinking-2507 improves with



Figure 6: **Reward Hacking Case Study.** In the absence of a length constraint ( $g_{len}$ ), the model exploits the unconstrained reward by endlessly enumerating possible entity translations to ensure verification success.

EA-RLVR (Avg. 90.36  $\rightarrow$  90.67). The reasoning trace absorbs the complexity of the entity task, allowing the final translation generation to remain fluent and robust.

**Note on Evaluation Subset.** It is important to note that the baseline Qwen3-4B-Thinking-2507 is highly unstable for translation tasks, often generating endless thought loops exceeding 32k tokens. Consequently, it fails to produce any translation for a large portion of the WMT24++ test set. To ensure a scientifically valid comparison, the results in Table 4 are calculated on the **common intersection** of sentences where the baseline model successfully produced an output. Table 5 details the data statistics. The valid subset size ranges from 68 to 260 samples per language (out of 960), highlighting the severity of the length explosion issue in the baseline model and the necessity of this filtering step.

## C Cross-Lingual Generalization

A core hypothesis of this work is that EA-RLVR does not merely fit language-specific translation

patterns, but rather induces a generalizable *reasoning strategy*—specifically, the mechanism of identifying entity-rich contexts and querying parametric knowledge. To decouple this reasoning capability from language-specific supervision, we conduct a zero-shot cross-lingual transfer experiment.

**Experimental Setup.** We partition the ten target languages into two typologically distinct clusters: **Group A (Asian & Semitic)**, comprising Arabic, Japanese, Korean, Thai, and Chinese, which largely utilize non-Latin scripts; and **Group B (European & Turkic)**, comprising German, Spanish, French, Italian, and Turkish. We train the Qwen3-8B model exclusively on one group and evaluate its performance on the disjoint group. Crucially, in the transfer setting (e.g., Train A  $\rightarrow$  Test B), the model receives *no supervision or reward signal* for the target languages during the RLVR stage.

**Zero-Shot Transfer Results.** Table 3 presents the results. We observe robust positive transfer in both directions, validating the hypothesis that the learned optimization is language-agnostic.

- **A  $\rightarrow$  B Transfer:** The model trained solely on Group A achieves an average entity accuracy of 33.19% on Group B, surpassing the base model (25.92%) by a significant margin of **+7.27%**.
- **B  $\rightarrow$  A Transfer:** Conversely, training on Group B improves performance on Group A languages to 20.74%, a gain of **+6.28%** over the base baseline.

**Discussion: Learning a Meta-Reasoning Skill.** The fact that RLVR improves entity translation accuracy on unseen languages—with distinct scripts and grammatical structures—suggests that the policy has acquired a meta-reasoning skill. The model learns *when* to pause generation to attend to specific entities and *how* to retrieve the corresponding cultural concepts from its pre-trained weights. This internal “lookup and verification” process is independent of the surface form of the target language, allowing the reasoning patterns tailored for one linguistic family to effectively facilitate translation in another.

## D Synergy with Retrieval-Augmented Generation

A central question in modern translation systems is the interplay between optimizing internal pa-

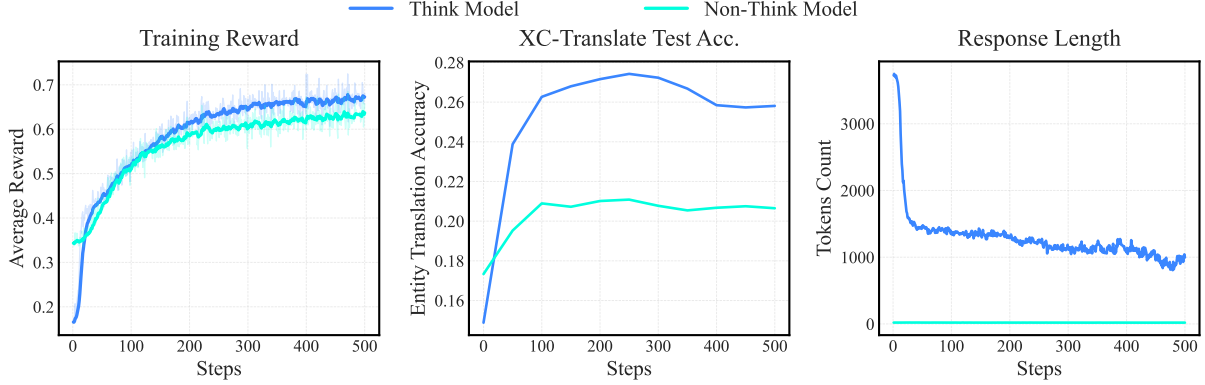


Figure 7: **Training Dynamics of Thinking vs. Non-Thinking Models.** (Left) The Think Model (Blue) initially suffers from low rewards due to context length overflows ( $> 4096$  tokens) but eventually surpasses the Non-Think Model (Cyan). (Center) The reasoning capability unlocks a significantly higher ceiling for entity translation accuracy on the XC-Translate test set. (Right) EA-RLVR acts as a strong regularizer for reasoning, rapidly curbing the “infinite thought” tendency of the base model to a stable, efficient length.

Table 4: Impact of reasoning on general translation quality (XCOMET on WMT24++). Due to the instability of the baseline Qwen3-4B-Thinking model (which frequently generates infinite reasoning traces exceeding 32k tokens), this evaluation is restricted to the **common subset** of sentences where the baseline model successfully produced a valid output. On this subset, EA-RLVR improves the reasoning model’s quality, whereas it degrades the standard model, suggesting that a thinking workspace is necessary to absorb the complexity of entity constraints without sacrificing fluency.

Model	XCOMET score on WMT24++ Subset ( $en \rightarrow X$ )										Avg.
	ar	de	es	fr	it	ja	ko	th	tr	zh	
<i>Standard Instruction Backbone</i>											
Qwen3-4B-Instruct	87.09	95.14	92.01	89.75	90.78	89.70	86.69	90.37	86.25	83.32	89.11
+ EA-RLVR	87.41	95.17	91.12	88.21	89.79	88.37	84.65	90.03	86.94	82.90	88.46
<i>Reasoning Backbone</i>											
Qwen3-4B-Thinking	88.84	95.74	92.48	90.20	92.21	90.66	89.39	<b>91.04</b>	<b>90.38</b>	82.65	90.36
+ EA-RLVR	<b>89.03</b>	<b>96.17</b>	<b>93.25</b>	<b>91.16</b>	<b>93.09</b>	<b>90.72</b>	<b>88.66</b>	90.79	89.69	<b>84.19</b>	<b>90.67</b>

rameters (via RLVR) and utilizing external non-parametric knowledge (via RAG). To determine whether our method complements retrieval-based approaches, we conduct an ablation study using a standard RAG pipeline.

**Experimental Setup.** We employ **mContriever** (Izacard et al., 2021), a widely used multilingual dense retriever, without any task-specific fine-tuning. We index the aliases of entities present in the XC-Translate test set (extracted from Wikidata via the provided QIDs). For each source sentence, we retrieve the top- $k$  most relevant entity aliases and prepend them to the system prompt as context. For the combined setting (**EA-RLVR + RAG**), we utilize the Qwen3-8B model trained via EA-RLVR and provide it with the same retrieved context during inference.

**Parametric Optimization Outperforms Naive Retrieval.** As shown in Table 6, the standard RAG baseline improves the base model’s entity accuracy from 20.19% to 23.14% and slightly boosts general quality (chrF 57.43), validating the effectiveness of our retrieval setup. However, **EA-RLVR alone significantly outperforms the RAG baseline** (29.25%). This result highlights a critical insight: for cross-cultural translation, the bottleneck is often not the *availability* of knowledge (which RAG provides), but the model’s ability to *align* that knowledge with the translation context. EA-RLVR addresses this alignment directly via optimization, proving more effective than passively injecting context.

**Additive Gains.** Crucially, the two approaches are synergistic. The **EA-RLVR + RAG** configuration achieves the highest overall performance

Table 5: Statistics of the evaluation subset for WMT24++. Due to the endless thinking issue, the Qwen3-4B-Thinking baseline yields valid outputs for only a fraction of the test set. We report results on the **Common** intersection to ensure fair comparison.

Language Pair	Total (Test Set)	Common Subset
en-ar	960	68
en-de	960	143
en-es	960	260
en-fr	960	193
en-it	960	211
en-ja	960	144
en-ko	960	135
en-th	960	124
en-tr	960	99
en-zh	960	241

(30.49% Acc, 60.05 chrF). This demonstrates that the reasoning patterns learned by EA-RLVR are robust; the model does not “overfit” to its internal weights but retains the flexibility to incorporate external evidence. By transforming the model into an active reasoner, EA-RLVR enables it to utilize retrieved context to resolve tail cases that neither parametric knowledge nor retrieval could solve alone.

## E Implementation Details

**Data Construction.** We conduct experiments on the XC-Translate benchmark across the ten language directions listed in Table 7. As the benchmark does not provide a dedicated training set, we repurpose the official validation set as our training split. Table 8 summarizes the statistics of our data split. Crucially, to strictly evaluate the model’s ability to generalize rather than memorize, we ensure that **the training and test sets share no overlapping entities**. The final setup comprises 7,278 examples for training and 49,606 examples for testing.

**Training Implementation.** We implement EA-RLVR using the ver1 library (Sheng et al., 2024), a framework designed for efficient RLHF post-training. Unless otherwise specified, all RL experiments across different model scales (8B, 14B) and variants (Instruct, Thinking) use the unified set of hyperparameters reported in Table 9.

**Compute and Environment.** All models were trained on  $32 \times$  NVIDIA H100 80GB GPUs. The 7k samples training for the 8B model takes approximately 12 hours, while the 14B model takes 24 hours. And the scaled training using full XC-

Translate for the 8B model takes approximately 24 hours, while the 14B model takes 48 hours. We use FlashAttention for efficient computation (Dao et al., 2022).

**SFT Baseline.** For the SFT baseline, we fine-tune the base models on the same 7k training examples for 2 epochs, supervised by reference translation. We use a learning rate of  $1e-6$  with a cosine decay schedule and a global batch size of 64. We SFT our model using ms-swift framework (Zhao et al., 2024b).

**EA-RLVR Configuration.** Our EA-RLVR optimization follows the critic-free policy gradient approach described in §3.3. We initialize the actor network with the Qwen3 weights post-trained by their official team, which endow the model basic reasoning capability. During the rollout phase, we sample  $G = 16$  responses for each prompt to compute the group-normalized advantages.

Table 9 lists the detailed hyperparameters used for the RLVR stage.

**Prompt Format.** We use the standard chat template of the Qwen3 family. For the input, we wrap the source sentence with the instruction: “Translate the following sentence into [Lang]...’”.

**Evaluation Configuration.** We treat the reasoning process (i.e., “thinking”) as an intrinsic capability of the models rather than a separate module. Consequently, we enable the thinking mode by default for all applicable models (e.g., Qwen3, Marcoo1 and GPT5-mini) and all settings including SFT and RAG. To ensure a fair comparison, we allocate a unified maximum generation budget of 4,096 tokens for all experiments. Regarding decoding strategies, we follow the best practices established by DeepSeek-R1 (DeepSeek-AI et al., 2026; Yang et al., 2025), setting the sampling temperature to 0.6 and top- $p$  to 0.95. This specific configuration is critical, as lower temperatures (e.g., greedy decoding) tend to induce severe repetition loops and infinite generation behaviors in reasoning-heavy models. Finally, to ensure statistical reliability, all reported results for open-weight models are averaged over three independent runs ( $pass@1$ ).

## F Unbiased $pass@k$ Estimator

While  $pass@k$  is defined as the probability of generating at least one correct sample in  $k$  attempts,

Table 6: **Synergy between Parametric and Non-Parametric Knowledge.** Comparison of EA-RLVR against a standard Retrieval-Augmented Generation (RAG) baseline on Qwen3-8B. Adding RAG to the base model yields moderate gains, validating the retrieval setup. However, EA-RLVR acting as a standalone method provides a substantially larger improvement. Crucially, the combined setting (**EA-RLVR + RAG**) achieves the highest accuracy and faithfulness (chrF), indicating that the reasoning capabilities induced by EA-RLVR effectively complement external knowledge retrieval.

Model	Entity Translation Accuracy on XC-Translate ( $en \rightarrow X$ )										Avg.
	ar	de	es	fr	it	ja	ko	th	tr	zh	Acc/ChrF
Qwen3-8B	14.91	23.64	31.04	25.97	25.01	17.13	11.12	5.63	23.94	23.50	20.19/56.07
+ RAG	17.13	32.85	36.62	30.58	28.82	16.35	11.79	6.18	27.81	23.25	23.14/57.43
+ EA-RLVR	25.23	31.01	44.44	34.89	35.35	24.02	25.70	14.74	27.79	29.31	29.25/59.86
+ EA-RLVR + RAG	<b>26.17</b>	<b>35.50</b>	<b>47.81</b>	<b>36.63</b>	<b>37.43</b>	<b>23.36</b>	<b>25.52</b>	<b>14.42</b>	<b>29.33</b>	<b>28.75</b>	<b>30.49/60.05</b>

Table 7: Languages used in our experiments, together with their ISO 639-1 codes, language-region locales, and English names.

ISO 639-1	Locale	English Name
ar	ar_SA	Arabic (Saudi Arabia)
de	de_DE	German (Germany)
es	es_MX	Spanish (Mexico)
fr	fr_FR	French (France)
it	it_IT	Italian (Italy)
ja	ja_JP	Japanese (Japan)
ko	ko_KR	Korean (Korea)
th	th_TH	Thai (Thailand)
tr	tr_TR	Turkish (Turkey)
zh	zh_TW	Chinese (Taiwan, Traditional)

Table 8: Statistics of the dataset used in our experiments. We utilize the official validation set of XC-Translate as our training split. The training and test sets are strictly disjoint in terms of entity coverage.

Language Pair	Train	Test	Total
English $\rightarrow$ Arabic	722	4,546	5,268
English $\rightarrow$ Chinese	722	5,181	5,903
English $\rightarrow$ French	724	5,464	6,188
English $\rightarrow$ German	731	5,875	6,606
English $\rightarrow$ Italian	730	5,097	5,827
English $\rightarrow$ Japanese	723	5,107	5,830
English $\rightarrow$ Korean	745	5,081	5,826
English $\rightarrow$ Spanish	739	5,337	6,076
English $\rightarrow$ Thai	710	3,446	4,156
English $\rightarrow$ Turkish	732	4,472	5,204
<b>Total</b>	<b>7,278</b>	<b>49,606</b>	<b>56,884</b>

1141 directly computing this probability typically re-  
1142 quires a very large number of samples to reduce  
1143 variance. To evaluate  $pass@k$  efficiently, we fol-  
1144 low the method proposed by Chen et al. (2021),  
1145 instead of just sampling  $k$  times, we generate a  
1146 larger number of samples  $n$  (where  $n \geq k$ ) for  
1147 each input and count the number of correct sam-  
1148 ples  $c$ . The unbiased estimator for  $pass@k$  is then

Table 9: Hyperparameters for EA-RLVR training.

Hyperparameter	Value
<i>Optimization</i>	
Optimizer	AdamW
Peak Learning Rate (Actor)	1e-6
Learning Rate Scheduler	Cosine
Warmup Ratio	0.05
Weight Decay	0.1
Train Batch Size	512
PPO Mini Batch Size	128
Total Training Steps	500
<i>PPO / Policy Gradient</i>	
Group Size ( $G$ )	16
Clip Ratio ( $\epsilon_{low}, \epsilon_{high}$ )	3e-4, 4e-4
Advantage Estimator	Group Normalization
<i>Generation / Rollout</i>	
Sampling Temperature	1.0
Top- $p$	1.0
Max Sequence Length	4096
<i>Reward Function</i>	
Format Reward ( $\alpha$ )	0.2
Length Tolerance ( $\tau$ )	2.0

calculated as:

$$pass@k := \mathbb{E}_{\text{problems}} \left[ 1 - \frac{\binom{n-c}{k}}{\binom{n}{k}} \right] \quad (7)$$

1149  
1150  
1151 where  $\binom{n}{k}$  denotes the number of combinations  
1152 of choosing  $k$  items from a set of  $n$ . Mathemati-  
1153 cally, this formula calculates the probability that a  
1154 randomly chosen subset of size  $k$  contains at least  
1155 one correct answer, derived from the complement  
1156 of the probability that all  $k$  chosen samples are  
1157 incorrect (i.e., chosen from the  $n - c$  incorrect sam-  
1158 ples).

1159 In our experiments, we set the total sample  
1160 budget  $n = 128$  and evaluate  $pass@k$  for  $k \in$   
1161  $\{1, \dots, 128\}$ . If  $n - c < k$ , the estimator returns  
1162 1.0, as it is impossible to choose  $k$  incorrect sam-  
1163 ples.

## G Case Study

**Qualitative Analysis** We present a qualitative analysis of four representative cases to illuminate the mechanism by which EA-RLVR improves entity translation. By examining the generated reasoning traces (denoted as LLM), we identify a distinct shift in cognitive patterns: while the baseline model (Qwen3-8B) relies on *literal semantic composition*, EA-RLVR exhibits *entity-aware deliberation* and *domain-specific retrieval*.

**Canonicalization of Historical Terminology (Case 1).** In Case 1, the user asks for the translation of “Great Ming Code”. The baseline Qwen3 adopts a compositional approach, translating “Great Ming” ( $\rightarrow$  明朝) and “Code” ( $\rightarrow$  法典) separately, resulting in the descriptive but non-standard phrase “明朝法典” (Ming Dynasty Code). In contrast, EA-RLVR’s reasoning trace explicitly triggers a hypothesis check: “*Great Ming Code is likely referring to the 大明律... I should confirm the standard translation.*”. By treating the phrase as a rigid proper noun rather than a translatable sentence fragment, EA-RLVR successfully retrieves the historiographically correct term “大明律”.

**Analogical Reasoning for Cultural Conventions (Case 2).** Case 2 illustrates how EA-RLVR leverages parametric knowledge for style transfer. For the film title “Once Upon a Time in Venezuela”, Qwen3 defaults to a dictionary translation of the idiom “Once Upon a Time” ( $\rightarrow$  很久很久以前), missing the cinematic context. EA-RLVR, however, employs **analogical reasoning**. The thinking trace reveals a crucial intermediate step: it recalls a prototype entity, “*Once Upon a Time in America is translated as 美往事*” and applies this naming convention to the target entity, synthesizing the culturally attuned title “委瑞拉往事” (Venezuela Chronicles/Past). This demonstrates the model’s ability to map new entities to existing cultural schemas.

**Domain-Specific Disambiguation (Case 3 & 4).** Polysemy poses a major challenge in entity translation. In Case 3, Qwen3 fails to resolve the term “Mahavira Hall” within a Buddhist context. Confused by the association of “Mahavira” with Jainism, it resorts to a phonetic transliteration “哈拉”. EA-RLVR correctly identifies the domain constraints: “*In Chinese Buddhist terminology... Mahavira is often transliterated as 大雄*”. It successfully navigates the semantic shift of “Mahavira” (Great Hero) in Buddhism to produce the correct

temple hall name “大雄殿”. Similarly, in Case 4 (French), EA-RLVR shows a willingness to localize media titles (“Emi magique”) rather than leaving them in English, reflecting a deeper adherence to target-language publication norms.

**Summary of Learned Reasoning Pattern.** Across these cases, a consistent meta-strategy emerges. The standard SFT/Base model tends to optimize for  $P(y|x)$  at the token level, favoring high-frequency phrases (fluency) over factual precision. EA-RLVR, driven by the entity-anchored reward, learns a “**Pause-Retrieve-Verify**” loop. It (1) detects potential cultural entities, (2) suspends immediate translation to search its parametric memory for domain equivalents, and (3) verifies the candidate against the target cultural context. This shift from *translating meaning* to *matching entities* is the core driver of the observed performance gains.

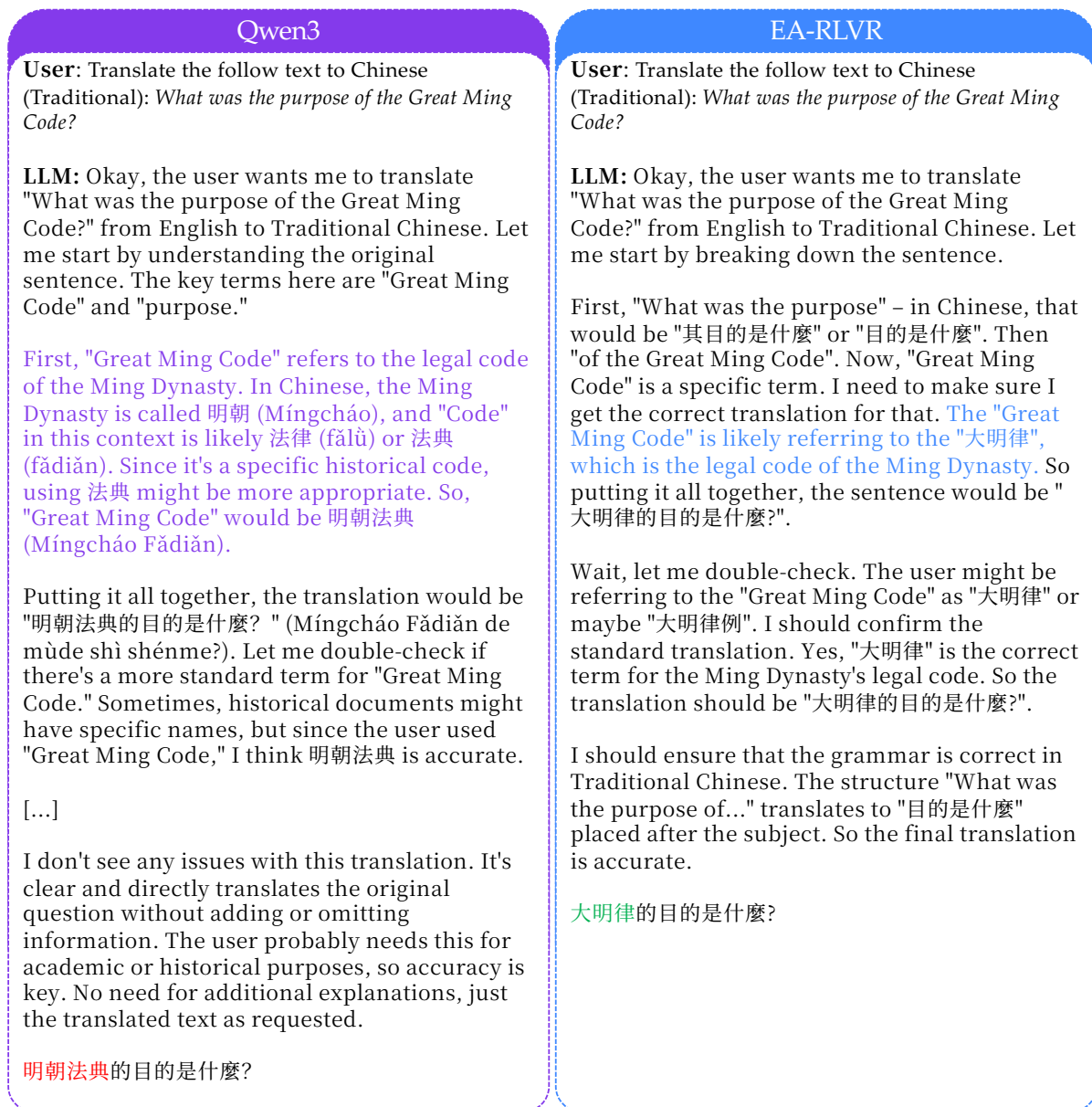


Figure 8: Case study 1: Canonicalization of Historical Terminology.

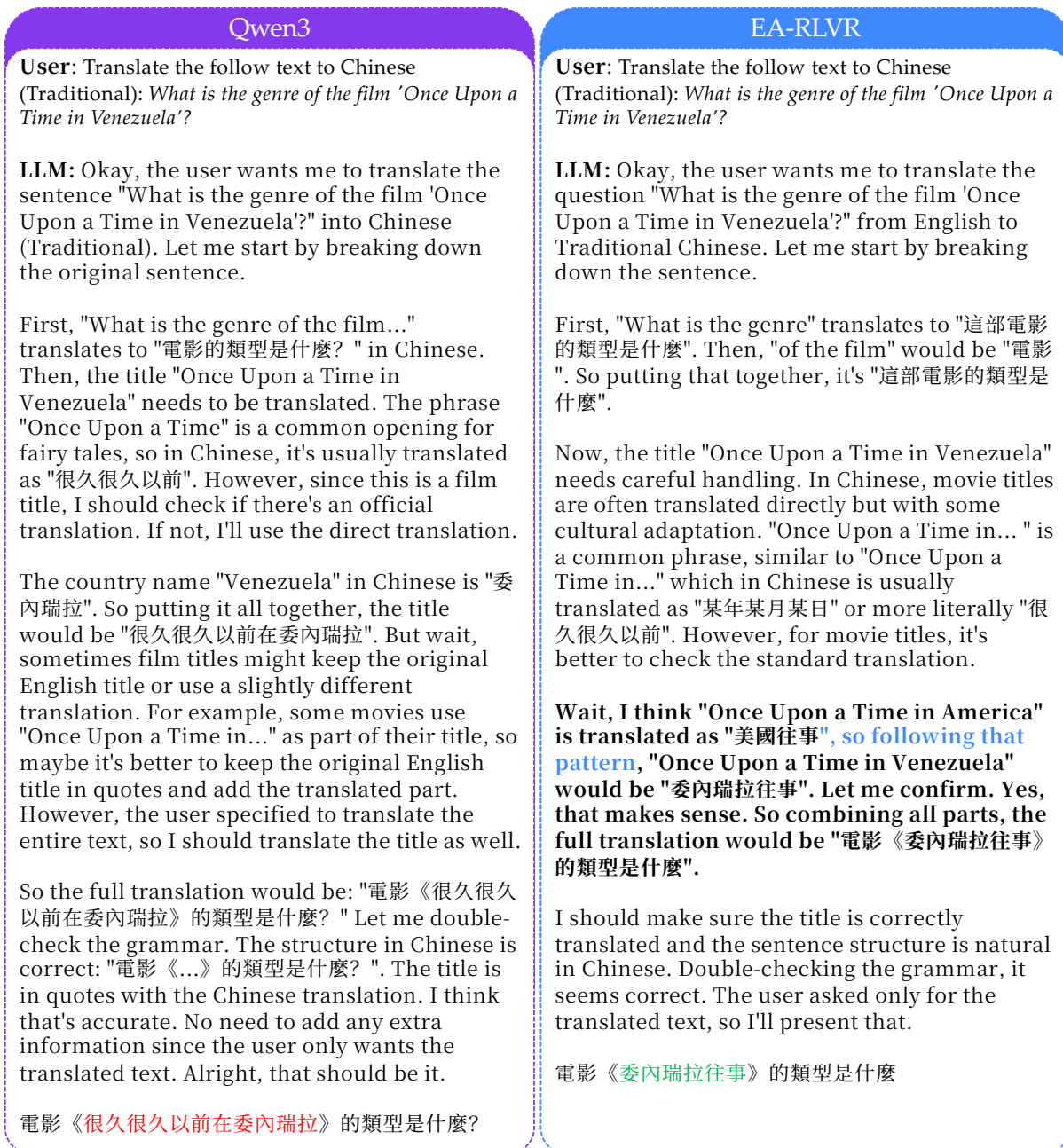


Figure 9: Case study 2: Analogical Reasoning for Cultural Conventions.

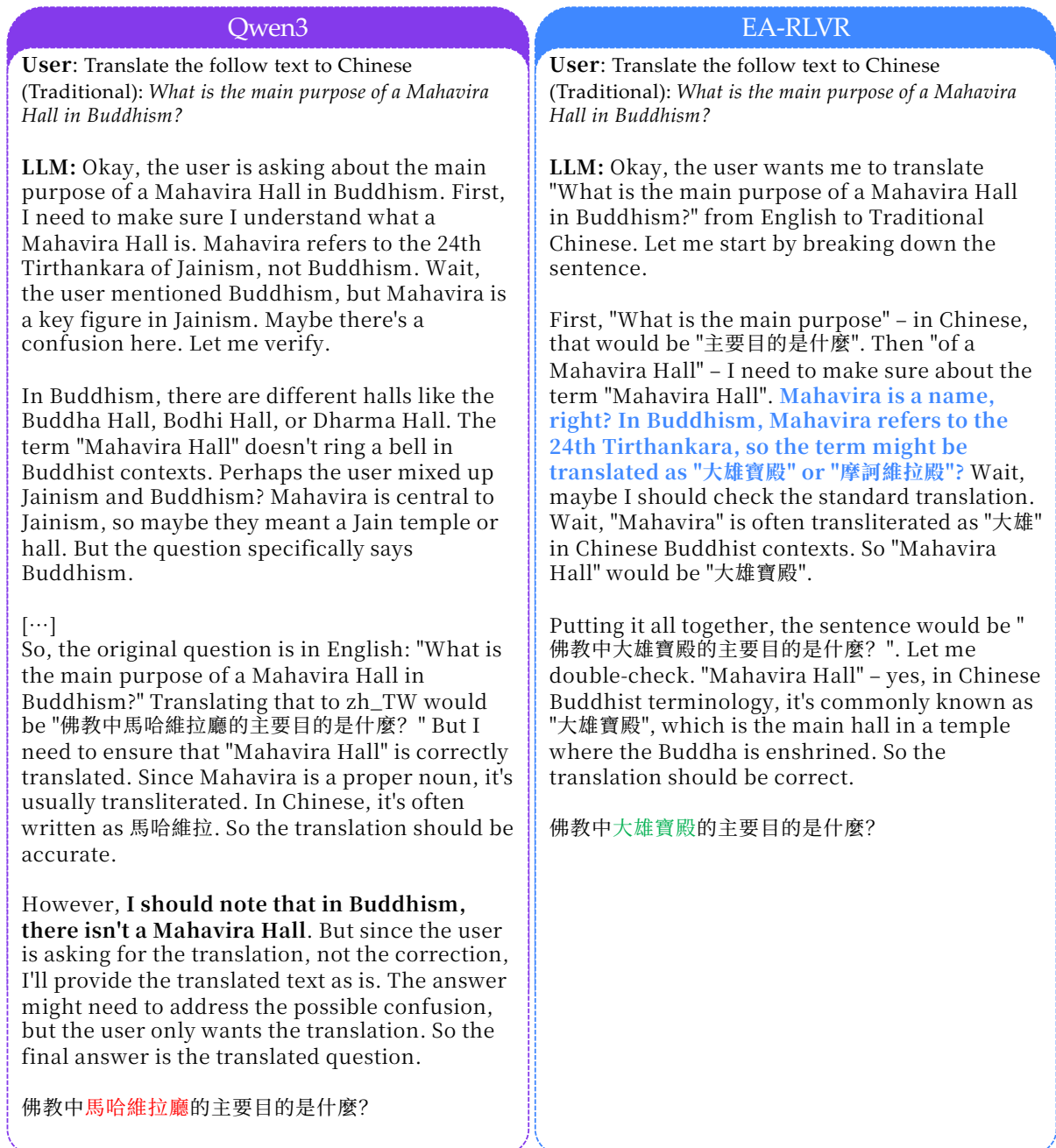


Figure 10: Case study 3: Domain-Specific Disambiguation .



Figure 11: Case study 4: Domain-Specific Disambiguation .