

# Remedy-R: Generative Reasoning for Machine Translation Evaluation without Error Annotations

Anonymous ACL submission

## Abstract

Over the years, scalar MT metrics have advanced rapidly on benchmarks. Yet they remain black boxes, offering little insight into their decisions and sometimes degrading under out-of-distribution inputs. We introduce Remedy-R, a reasoning-driven generative MT metric trained with reinforcement learning from pairwise translation preferences, without requiring error-span annotations or distillation from closed LLMs. Unlike scalar MT metrics that only outputs translation quality scores, Remedy-R produces step-by-step analyses of accuracy, fluency, and completeness, enabling more interpretable assessments. With only 60K pairwise training samples across two language pairs, Remedy-R remains competitive with top scalar metrics and GPT-4-based judges on WMT22–24 metric benchmarks, generalizes to other languages, and shows strong robustness on OOD stress tests. Moreover, Remedy-R generates self-reflective feedback that can be reused for translation refinement. We validate the faithfulness of such feedback with GPT-4 and show that a simple evaluate–revise pipeline leveraging Remedy-R’s analyses consistently improves translation quality across diverse models without any task-specific tuning.<sup>1</sup>

## 1 Introduction

Recent neural machine translation (MT) evaluation metrics like xCOMET (Guerreiro et al., 2024), MetricX (Juraska et al., 2024), and Remedy (Tan and Monz, 2025) achieve strong correlations with human preferences, and in some settings even report agreement that exceeds expert annotators (Proietti et al., 2025). Yet, they remain black boxes, producing a single scalar score with little insight into why a translation is good or bad.

This opacity matters because translation quality is multi-dimensional, involving criteria like accuracy, fluency, and completeness. A single number

does not reveal which dimension drives the judgment, and it also complicates robustness: without explicit decision-making process, metrics may exploit spurious cues learned during training and degrade under out-of-distribution (OOD) inputs such as source copy and input perturbations (Lo et al., 2023; Knowles et al., 2024; Moghe et al., 2025). As a result, existing MT metrics are powerful but hard to interpret and diagnose.

Recent work has attempted to enhance interpretability through translation error-span prediction, highlighting words or phrases that contribute to errors (Guerreiro et al., 2024; Treviso et al., 2024). Despite low overlap with human-annotated spans, such predictions are inherently local. In particular, span-based feedback struggles with omissions, discourse-level issues, and sentence-level phenomena not well captured by span annotations (e.g., translations in wrong or mixed languages): a translation may have few or no highlighted spans yet still be inadequate overall.

Another direction uses LLMs-as-judges to perform multi-dimensional assessment like MQM (Freitag et al., 2021). However, many LLMs-as-judges approaches still reduce the analysis to a final quality score (Kocmi and Federmann, 2023b; Fernandes et al., 2023; Qian et al., 2024), and existing works often rely on closed commercial LLMs (Lu et al., 2024; Kocmi and Federmann, 2023a; Zhang et al., 2025; Feng et al., 2025b), while open models struggle to match them (Lu et al., 2025). This raises a key question: how can we train an open, reasoning-based MT metric that produces informative analyses and aligns with human preferences, without expensive error annotations?

We introduce *Remedy-R*, a generative MT metric that produces step-by-step analyses along three dimensions (accuracy, fluency, completeness) and a final score. Inspired by recent pairwise preference framework in Remedy (Tan and Monz,

<sup>1</sup><https://anonymous.4open.science/r/Remedy-R-anonymous-5978/>

## Remedy-R ref-free Evaluation Example

**Source:** Even if I'm not in the income level that would pay more taxes, I can still buy fewer things that are more expensive so that they can be made by people making a living wage and I can be willing to wait a few days to get it so that some gig worker doesn't have to be worked to the bone...

**Translation:** Sogar wenn ich nicht in der Einkommensklasse bin, in der mehr Steuern gezahlt werden, kann ich immer noch weniger teure Dinge kaufen, sodass sie von Menschen hergestellt werden können, die maßgeblich dazu beitragen, und ich bin bereit, ein paar Tage zu warten, bevor der Arbeitslohn abläuft...

### Remedy-R's Evaluation:

Let's evaluate the translation step by step:

#### • Accuracy

- "sodass sie von Menschen hergestellt werden können, die maßgeblich dazu beitragen" is not clear and does not accurately reflect the original meaning. The original text states that the products should be made by people who are paid a living wage, but the translation does not convey this clearly.
- "bevor der Arbeitslohn abläuft" is incorrect and does not make sense in this context. The original text states "so that some gig worker doesn't have to be worked to the bone" which means the gig worker is not overworked. The translation should reflect this meaning.

#### • Fluency

- The sentence structure is somewhat complex and could be clearer.
- There are some awkward phrasings, such as "in der Einkommensklasse bin, in der mehr Steuern gezahlt werden," which could be simplified.
- The use of "maßgeblich dazu beitragen" is not clear and does not fit the context.

#### • Completeness

- The original text also states that the speaker is willing to wait a few days to get the product so that a gig worker doesn't have to be overworked. This part is also not clearly conveyed in the translation.

#### Score: 65.

2025), Remedy-R is trained with Reinforcement Learning via Verifiable Rewards (RLVR) from *pairwise human preferences*. Specifically, we define a deterministic reward that checks whether the model-implied preference (from parsed final scores) matches the human preference label, and optimize Remedy-R with PPO (Schulman et al., 2017) to maximize verifiable rewards. Unlike xTower, which depends on GPT-4 distillation and supervised learning with error spans (Treviso et al., 2024), Remedy-R trains directly with RLVR on MT preference data, without requiring error-span annotations, distillation, or SFT cold-start.

We assess whether Remedy-R's analyses are faithful and practically useful in two complementary ways. First, we validate explanation faithfulness with GPT-4o-mini (§5.1). Second, we introduce *Remedy-R Agent*, a simple training-free evaluate–revise loop that reuses Remedy-R's analyses as feedback to refine translations (§5). Despite never being trained for refinement, the agent consistently improves translation quality across diverse models (Qwen2.5 (Yang et al., 2024), ALMA-R (Xu et al., 2024b), GPT-4o-mini (Achiam et al., 2023), and Gemini-2.0-Flash) and generalizes to 11 language pairs beyond its two training language

pairs. Our contributions are:

- We propose *Remedy-R*, a generative MT evaluator that produces step-by-step analyses and a final score, trained directly via RLVR without error span annotation or distillation.
- Remedy-R achieves competitive performance on WMT22–24 meta-evaluation with only 60K training pairs from two language pairs, and exhibits strong OOD behavior.
- Remedy-R Agent, as a training-free framework, improves translations across diverse model families, and maintains cross-lingual generalization beyond its training languages.

## 2 Method

### 2.1 Task Formulation

We revisit MT evaluation as follows. Given a source sentence  $src$ , a translation  $mt$ , and an optional reference  $ref^*$  (with  $ref^* = \emptyset$  in reference-free settings), a metric  $M$  outputs a scalar quality score. Conventional learned metrics directly map  $(src, mt, ref^*)$  to a single number, which limits interpretability and reusability for downstream refinement. In Remedy-R, we formulate MT evaluation as *conditional text generation* guided by

an instruction that specifies the task and criteria, namely accuracy, fluency, and completeness. The model follows a reason then score protocol: it first writes a short analysis and then outputs a numeric score that we can parse for evaluation. Here, let the input be:

$$\mathbf{x} = \langle inst, src, mt, ref^* \rangle.$$

where the instruction *inst* specifies the evaluation task and criteria (e.g., accuracy, fluency, completeness). The model produces an output sequence  $\mathbf{y}$  in a *reason-then-score* format: a step-by-step reasoning analysis and a final numeric score in  $[0, 100]$ . We parameterize a conditional policy  $\pi_\theta(\mathbf{y} | \mathbf{x})$  and generate autoregressively:

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

Here  $T$  denotes the output length. The final score in the output is directly parsed for evaluation and will support the verifiable rewards used in §2.2.

## 2.2 Remedy-R Reward Design

We optimize Remedy-R with reinforcement learning using a simple, fully verifiable reward aligned with human judgments. The reward combines (i) a sparse *pairwise ranking* signal that matches the model-implied preference to the human preference label, and (ii) a reward shaping term that encourages predicted scores to be close to human scores, turning the sparse signal into a richer, continuous one. We maximize the expected reward with PPO; training details are in §2.3 (and Appendix A.1).

### 2.2.1 Pairwise ranking reward.

For each source sentence, we have two translations  $mt_A, mt_B$  with human scores  $g_A, g_B \in [0, 100]$ . We define the human preference label by comparing  $g_A$  and  $g_B$  and exclude ties ( $g_A = g_B$ ). We instruct the model to produce a reasoning COTs path and two final quality scores  $s_A, s_B \in [0, 100]$  in a fixed format (see the training template below). Instead of predicting a ranking label directly, we ask the model to evaluate A and B independently and to assign scores to each; this enables single segment evaluation at inference time without quadratic pairwise comparisons. We randomize the order of A/B during training to reduce position bias (Sproat et al., 2025).

The pairwise ranking reward then checks whether the model’s predicted ranking matches the

human ranking:

$$r_{\text{rank}} = \begin{cases} 1, & \text{if model and human rankings agree,} \\ 0, & \text{otherwise.} \end{cases} \quad (2)$$

This ranking reward is sparse and binary, but fully verifiable because it depends only on parsed rankings and human labels. If the score block cannot be parsed or falls outside the valid range, the reward is zero, and ties are treated as zero reward as well.

### 2.2.2 Reward shaping.

The pairwise ranking reward encourages correct relative preferences, but it does not calibrate score magnitudes: many score pairs can yield the same ranking and thus the same  $r_{\text{rank}}$ . For instance, if the model outputs  $(s_A, s_B) = (100, 99)$  and the human scores are  $(g_A, g_B) = (100, 0)$ , both cases receive  $r_{\text{rank}} = 1$ . This is undesirable because we ultimately want Remedy-R’s numeric scores to be meaningful and comparable across translations and systems, rather than only producing correct orderings. To provide a denser learning signal while remaining robust to noisy human scores, we add a calibration term that penalizes the deviation between predicted scores and human scores.

Specifically, let  $e_A = s_A - g_A$  and  $e_B = s_B - g_B$ . We adopt a Huber-style penalty:

$$\rho_c(e) = \begin{cases} \frac{1}{2} \frac{e^2}{c}, & |e| \leq c, \\ |e| - \frac{1}{2}c, & |e| > c, \end{cases} \quad (3)$$

where  $c$  defines a tolerance region (we set  $c=5$  for all experiments). The Huber form penalizes small errors smoothly while being less sensitive to large outliers, which is helpful given annotator noise and near-tie cases (Freitag et al., 2021; Tan and Monz, 2025). We normalize and average the penalties to obtain

$$\psi = \frac{1}{2} \left( \frac{\rho_c(e_A)}{c} + \frac{\rho_c(e_B)}{c} \right). \quad (4)$$

The final shaped reward is

$$r = r_{\text{rank}} \cdot (1 - \beta \psi), \quad (5)$$

where  $\beta$  controls shaping strength. We apply shaping only when the ranking is correct (through multiplication by  $r_{\text{rank}}$ ), preserving strict verifiability while improving score calibration. Appendix A.2 presents an ablation study on reward shaping. We also tested an auxiliary explanation-quality penalty using genRM, which provided only marginal gains. Accordingly, we adopt the ranking reward with Huber shaping for all main results.



Type	Methods	$\theta$	System-Level	Segment-Level $acc_{eq}^*$			Avg	
			Acc (3 LPs)	Avg	En-De	En-Ru	Zh-En	Corr
Scalar Metrics	COMET-22-DA	0.5B	82.8	54.5	58.2	49.5	55.7	68.7
	COMET-22 (ensemble)	5x0.5B	83.9	57.3	60.2	54.1	57.7	70.6
	MetricX-XXL	13B	85.0	58.8	61.1	54.6	60.6	71.9
	PaLM-2 BISON FT	>100B	88.0	57.3	61.0	51.5	59.5	72.7
	ReMedy	9B	91.2	58.9	61.0	60.4	55.4	75.1
LLM Judges	EAPrompt (Llama2)	70B	85.4	52.3	55.2	51.4	50.2	68.9
	EAPrompt (Mistral)	8x7B	84.0	50.9	53.8	50.6	48.2	67.5
	EAPrompt (GPT3.5-Turbo)	>100B	91.2	53.3	56.7	53.3	50.0	72.3
	GEMBA-MQM (Qwen)	72B	84.7	53.8	56.0	54.7	50.6	69.3
	MQM-APE (Qwen)	72B	85.8	54.5	56.4	55.7	51.4	70.2
	MQM-APE (Mistral)	8x22B	88.3	54.2	56.9	55.1	50.6	71.3
	GEMBA-DA (GPT4)	>100B	89.8	55.6	58.2	55.0	53.4	72.7
	PaLM	540B	90.1	50.8	55.4	48.6	48.5	70.5
	Remedy-R	7B	89.1	54.8	58.0	56.0	50.4	71.9
Remedy-R	14B	88.7	56.0	58.0	55.8	54.2	72.4	
Remedy-R	32B	91.6	55.2	57.8	55.7	52.2	73.4	

Table 1: Evaluation on WMT22 MQM set. Following official WMT22 settings, we report system-level Pairwise Accuracy (Acc) and segment-level pairwise accuracy with tie calibration ( $acc_{eq}^*$ ), using Perm-Both statistical significance test (Deutsch et al., 2021). orange and blue indicate the best performing scalar and LLM judge metrics. Our Remedy-R 7B outperforms all LLM judges with 70B parameters, and Remedy-R 32B surpasses GEMBA-DA (GPT4).

## 4 Automatic Translation Evaluation Results and Analyses

In this section, we evaluate Remedy-R on automatic translation evaluation tasks and address three questions: (a) how Remedy-R compares with top scalar MT metrics and LLM-as-judge (§4.1); (b) how performance changes under test-time-scaling with multiple evaluation passes (§4.2); and (c) whether the model behaves robustly on out-of-domain stress tests (§4.3).

### 4.1 Meta Evaluation on WMT benchmarks

As shown in Table 1, Remedy-R achieves leading performance on the WMT22 metric benchmark compared to both commercial LLM judges and open metrics. Specifically, Remedy-R 7B surpasses all open LLM-as-judge baselines with 70B parameters (e.g., EAPrompt and MQM-APE), while Remedy-R 32B exceeds GEMBA-DA (GPT-4). Among scalar metrics, Remedy-R also outperforms the regression-based PaLM-2 BISON and MetricX-XXL, showing that our model aligns with human ratings on human translation quality judgments.

In addition, we also provide additional results on WMT23 and WMT24 metric benchmarks in

Table 12 and Table 6. The results highlight that Remedy-R achieves on-par or superior performance to the current SOTA metrics. For example, on wmt23, Remedy-R-7B outperforms EAPrompt that is based on GPT4o-mini, and Remedy-R 14B and 32B performs on par with KIWI-XXL and MetricX-23. On wmt24, Remedy-R outperforms GEMBA-ESA (GPT4), MetricX-24-Hybrid, and XCOMET-XXL (see details in Table 6).

### 4.2 Test Time Scaling with multiple Evaluation Passes

Remedy-R’s generative reasoning nature enables the application of *Test-Time Scaling (TTS)*, where multiple evaluation passes are performed with different reasoning trajectories and their quality scores are aggregated. In this setting, we adopt a simple implementation that averages the quality scores from multiple independent evaluations.

We found performing more evaluation trajectories at test time consistently enhances performance across all model sizes (see Figure 3 in Appendix). Notably, Remedy-R-14B reaches an average correlation of 74.9, matching the strongest GEMBA-MQM performance. The steady improvement from 7B to 32B suggests that iterative reasoning stabi-

	ref?	empty mt	empty src+ref	src copy	wrong lang	mix lang	unrelated mt
COMET-22	✓	57.00%	58.81%	69.85%	67.84%	65.56%	45.23%
KIWI	✗	54.87%	67.72%	52.15%	82.64%	78.75%	41.95%
XCOMET	✓	73.79%	64.12%	82.04%	85.65%	71.77%	20.31%
MetricX-24-XXL	✓	-9.59	-5.85	-12.59	-3.06	-10.08	-24.15
MetricX-24-XXL	✗	-7.34	-5.85	-11.36	-2.51	-7.78	-24.25
GEMBA-ESA	✗	14.00%	13.5%	11.12%	14.32%	18.08%	1.27%
ReMedy-R-7B	✓	1.00%	7.07%	76.92%	43.69%	60.6%	1.5%
ReMedy-R-14B	✓	0.00%	0.00%	11.35%	14.6%	37.6%	0.6%
ReMedy-R-32B	✓	0.00%	0.00%	2.76%	8.30%	46.0%	1.3%

Table 2: Averaged quality scores of different metric models on MSLC24 OOD set. For all classes except *mix-lang*, a robust metric should output low scores; for *mix-lang*, the translation preserves the source meaning but contains code-switching, so its quality scores should be moderately high rather than near zero. MetricX scores are ranged from -25 to 0 (higher is better). We provide additional reference-free Remedy-R results in Table 7.

lizes evaluation outcomes and reduces stochastic variance, yielding more robust and reliable quality assessments.

Interestingly, we observe that TTS primarily improves *segment-level*  $acc_{eq}^*$  rather than system-level correlation (see Table 12). We hypothesize that this phenomenon is mostly due to the limitations of current meta-evaluation metrics. As noted by [Perrella et al. \(2024\)](#), tie-calibrated pairwise accuracy ( $acc_{eq}^*$ ) tends to favor metrics that output continuous rather than discrete scores. Averaging multiple predictions effectively smooths discrete outputs into continuous scores, improving agreement with tie-calibrated accuracy, which favors metrics with finer score granularity.

### 4.3 Analyses on Challenge sets

We evaluate Remedy-R under out-of-distribution (OOD) and adversarial conditions on MSLC24, and additionally construct an *unrelated translation* category by sampling 50 target-language sentences from Flores-200 ([Costa-Jussà et al., 2022](#)). For most perturbations (empty outputs, source copies, wrong-language, unrelated MT), a reliable metric should assign near-zero scores, while mixed-language outputs with correct meaning but code-switching should receive moderate scores.

Table 2 shows that scalar metrics (e.g., COMET-22, MetricX-24) can behave inconsistently under OOD inputs, sometimes assigning high scores to degenerate cases; for example, XCOMET scores 73.8%, 64.1%, and 82.0% on empty translation, empty source/reference, and source copy, respectively. In contrast, Remedy-R is robust and well-calibrated: it outputs near-zero scores for empty

translations, sharply penalizes source-copy and wrong-language cases, and assigns moderate scores to mixed-language inputs.

Unlike MetricX-24 and XCOMET, which use synthetic augmentation that includes empty or hallucinated translations, Remedy-R is trained without such data yet still exhibits coherent behavior across perturbations, indicating stronger generalization beyond standard WMT conditions.

## 5 ReMedy-R Agent

So far, our analyses have focused on evaluation at the final score level, i.e., how well Remedy-R’s predicted quality scores align with human ratings. However, such correlations do not necessarily imply that the generated analyses are faithful or practically useful.

We therefore evaluate Remedy-R’s analyses in two complementary ways: (i) we prompt GPT-4o-mini to score explanation faithfulness given only the source and the translation (§5.1); and (ii) we reuse the analyses as feedback in a simple training-free evaluate–revise loop (*Remedy-R Agent*) and measure translation improvements (§5.2.1–§5.2.2).

### 5.1 How faithful are Remedy-R’s evaluation explanations?

We first assess explanation faithfulness by prompting GPT-4o-mini to score how well Remedy-R’s reasoning is supported by the source sentence and the translation hypothesis only. We conduct this analysis on the WMT22 metric evaluation test set and sample 300 examples for each MQM language pair (en-de, en-ru, zh-en), resulting in 900 exam-

$M_{feedback}$	$M_{refinement}$	cs-uk	en-cs	en-de	en-es	en-hi	en-is	en-ja	en-ru	en-uk	en-zh	ja-zh	Avg
$M_{base} = \text{Qwen2.5-it-14B} \mid \text{Remedy-R} = 14\text{B} \mid \text{x-Tower} = 14\text{B} \text{ (w. XComet-XL} = 3.5\text{B)}$													
-	-	69.4	63.6	88.4	83.7	47.6	32.6	74.8	75.7	58.6	83.9	72.5	68.2
-	Base	71.0 $\uparrow_{0.6}$	67.2 $\uparrow_{3.6}$	89.4 $\uparrow_{1.0}$	85.0 $\uparrow_{1.3}$	51.5 $\uparrow_{3.9}$	34.4 $\uparrow_{1.9}$	77.9 $\uparrow_{3.1}$	77.9 $\uparrow_{2.2}$	63.2 $\uparrow_{4.6}$	84.5 $\uparrow_{0.6}$	72.7 $\uparrow_{0.2}$	70.4 $\uparrow_{2.2}$
x-Tower	x-Tower	77.1 $\uparrow_{7.7}$	62.5 $\downarrow_{1.1}$	91.4 $\uparrow_{3.0}$	86.5 $\uparrow_{2.8}$	45.0 $\downarrow_{2.6}$	34.1 $\uparrow_{1.6}$	65.6 $\downarrow_{9.2}$	80.2 $\uparrow_{4.5}$	70.4 $\uparrow_{11.8}$	82.8 $\downarrow_{1.1}$	68.3 $\downarrow_{4.2}$	69.4 $\uparrow_{1.2}$
Remedy-R	Remedy-R	74.1 $\uparrow_{4.7}$	68.3 $\uparrow_{4.7}$	89.8 $\uparrow_{1.4}$	84.9 $\uparrow_{1.2}$	52.6 $\uparrow_{5.0}$	35.8 $\uparrow_{3.2}$	77.8 $\uparrow_{2.9}$	77.4 $\uparrow_{1.7}$	64.3 $\uparrow_{5.7}$	83.8 $\downarrow_{0.1}$	72.4 $\downarrow_{0.1}$	71.0 $\uparrow_{2.8}$
x-Tower	Base	73.0 $\uparrow_{3.6}$	68.5 $\uparrow_{4.9}$	90.2 $\uparrow_{1.8}$	85.3 $\uparrow_{1.7}$	51.3 $\uparrow_{3.7}$	34.6 $\uparrow_{2.0}$	78.6 $\uparrow_{3.8}$	78.5 $\uparrow_{2.8}$	65.1 $\uparrow_{6.5}$	84.3 $\uparrow_{0.4}$	71.9 $\uparrow_{0.6}$	71.0 $\uparrow_{2.8}$
Remedy-R	Base	74.8 $\uparrow_{5.4}$	68.2 $\uparrow_{4.6}$	90.3 $\uparrow_{1.9}$	85.6 $\uparrow_{1.9}$	53.1 $\uparrow_{5.5}$	36.3 $\uparrow_{3.7}$	79.3 $\uparrow_{4.5}$	79.2 $\uparrow_{3.5}$	64.5 $\uparrow_{5.9}$	84.1 $\uparrow_{0.2}$	72.4 $\downarrow_{0.1}$	71.6 $\uparrow_{3.4}$

Table 3: Agent MT experiments on WMT24 using Qwen2.5-14B models as the initial  $M_{base}$  translator (gray background). We report XCOMET-XXL in this table and provide results of more model variants in Appendix A.7

ples in total. Table 4 reports the average faithfulness scores. Overall, Remedy-R explanations receive consistently high faithfulness scores, and faithfulness improves with model scale.

Remedy-R	en-de	en-ru	zh-en	Avg
7B	79.10	76.57	74.95	76.87
14B	80.02	76.75	78.05	78.27
32B	81.18	79.05	78.38	79.54

Table 4: GPT-4o-mini faithfulness scores for Remedy-R explanations on WMT22 MQM test samples (300 per language pair). GPT-4o-mini is given the source, translation, and explanation. Higher score means more faithful explanations.

## 5.2 Reusing Remedy-R Analyses for Training-Free Refinement

In the evaluate–revise setting, three models are involved:

- $M_{base}$ : a translator that produces an initial translation for a source sentence;
- $M_{feedback}$ : an evaluator that inspects the translation and generates an analysis as feedback;
- $M_{refinement}$ : a refiner that revises the translation conditioned on  $(src, mt, feedback)$ .

Instantiating these roles yields a training-free loop (*Remedy-R Agent*) that reuses Remedy-R’s analyses for translation refinement. We evaluate the loop under controlled configurations and on strong open (Table 3 and 10) and closed (Figure 1) MT systems .

### 5.2.1 Controlled agent experiments with Qwen translators

We first test whether Remedy-R’s analyses provide useful guidance beyond self-refinement by

fixing  $M_{refinement} = M_{base}$  and comparing: (1) self-refinement ( $M_{feedback} = -$ ), where the translator revises its own output without external feedback; and (2) feedback-refinement ( $M_{feedback} = \text{Remedy-R}$ ), where the translator revises using Remedy-R’s analysis as additional context. Across model scales, feedback-refinement yields consistent improvements (e.g., +3.4 xCOMET score for 14B Qwen).

We further compare Remedy-R feedback to x-Tower (Treviso et al., 2024), a GPT-4–distilled evaluator that relies on external span supervision from xCOMET-XL (Guerreiro et al., 2024). At similar scales, Remedy-R achieves larger and more stable gains (14B Qwen: +3.4 with Remedy-R feedback vs. +1.2 for x-Tower as refiner in Table 3).

Finally, we test a unified setting where  $M_{feedback} = M_{refinement} = \text{Remedy-R}$ , i.e., Remedy-R generates an analysis and then refines based on its own feedback. This unified setting performs comparably to using Remedy-R for feedback and the base translator for refinement (Table 3), indicating that Remedy-R’s analyses can be effectively reused for refinement even without any post-editing training.

### 5.2.2 Applying Remedy-R Agent to strong closed LLMs

Finally, we test whether Remedy-R Agent remains effective when applied to strong closed LLM translators. Following prior work on multi-pass translation refinement (Briakou et al., 2024; Wu et al., 2025), we evaluate on the WMT24++ benchmark (Kocmi et al., 2024) and set  $M_{base}$  to GPT-4o-mini and Gemini-2.0-Flash. We use Remedy-R-32B as  $M_{feedback}$  and apply iterative training-free evaluate–revise steps. For comparison, we include each model’s internal self-refinement baselines, including step-by-step (Briakou et al., 2024) and translate-again (Wu et al., 2025).

As shown in Figure 1, Remedy-R Agent consistently improves translation quality on both

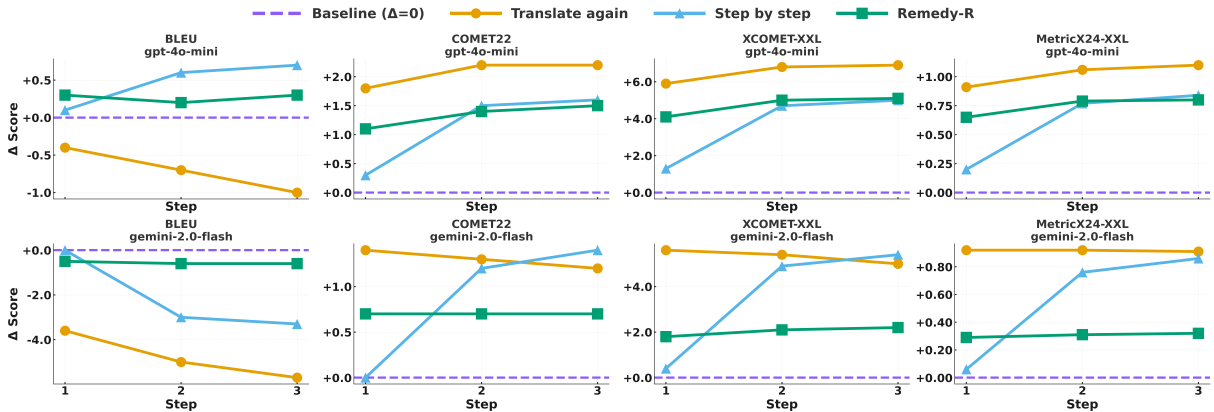


Figure 1: Refinement performance comparison on the initial translations from GPT-4o-mini and Gemini-2.0-Flash using paragraph-level WMT24++ benchmark. We use ref-based XCOMET-XXL to measure the translation quality. Here Translate-again and Step-by-Step utilizes self-refinement (using GPT-4o-mini and gemini-2.0-flash themselves for refinement), while we use the 32B version for Remedy-R. More detailed results are in Appendix (Table 11).

closed LLMs. On GPT-4o-mini, Remedy-R Agent achieves gains comparable to step-by-step refinement, while on Gemini-2.0-Flash it attains a substantial fraction of the gains from self-refinement. Notably, Remedy-R Agent uses only a 32B open model for feedback, yet remains effective on translations produced by significantly larger closed LLMs, highlighting the generality of Remedy-R’s feedback across model families.

In addition, we observe consistent improvements when applying Remedy-R Agent to ALMA-R (Xu et al., 2024a), see Appendix A.7.

## 6 Related Work

### 6.1 Automatic Translation Evaluation

Early MT metrics such as BLEU (Papineni et al., 2002) and ChrF (Popović, 2015) rely on surface matching and correlate weakly with human judgments (Freitag et al., 2022). Learned scalar metrics (Rei et al., 2020; Juraska et al., 2023; Tan and Monz, 2025) achieve higher correlations with human, but remain opaque and can degrade under OOD conditions (Lo et al., 2023; Knowles et al., 2024). For interpretability, xCOMET (Guerreiro et al., 2024) predicts MQM-based error spans, and xTower (Treviso et al., 2024) distills GPT-4 explanations conditioned on such spans.

### 6.2 Reasoning for Machine Translation

Reasoning has been explored for both MT evaluation and generation. For evaluation, recent LLM-as-judge systems use multi-stage or multi-agent designs to structure assessments (Feng et al., 2025b;

Zhang et al., 2025). For generation, RL-based approaches optimize translation quality using learned rewards, e.g., MT-R1-Zero (Feng et al., 2025a) and Hunyuan-MT (Zheng et al., 2025).

### 6.3 Translation Agents

Agentic translation methods iteratively improve translations via self-refinement or decomposed sub-tasks, e.g., translate-again (Wu et al., 2025) and step-by-step (Briakou et al., 2024). Closest to our setting, xTower (Treviso et al., 2024) can be used to support post-editing via span-based explanations. Our goal is not to propose a new translation pipeline; we use a simple evaluate–revise loop as a lightweight diagnostic to test whether Remedy-R’s feedback can guide refinement.

## 7 Conclusions

A longstanding challenge in MT evaluation is limited explainability, which undermines metric reliability and robustness. We introduce Remedy-R, a generative reasoning-based MT metric trained with RLVR from pairwise human preferences using a verifiable reward, without error-span annotations. Remedy-R produces multi-dimensional analyses (accuracy, fluency, completeness) followed by a final score, achieves competitive performance on WMT22–24 benchmarks, and remains robust on OOD tests. Beyond score-level correlation, we validate the faithfulness and practical value of its quality analyses via Remedy-R Agent: a training-free evaluate–revise loop that consistently improves translations across models and language pairs, including strong open and commercial systems.

559	<b>Limitations</b>		
560	In this paper, we do not conduct human evaluation		
561	to measure the faithfulness of Remedy-R’s		
562	evaluation explanations due to limited budgets.		
563	However, we conduct this experiments with GPT-		
564	4O-mini with 900 random samples in the official		
565	WMT22 metric test set. As a generative evaluator,		
566	Remedy-R produces natural-language reasoning		
567	explanations that may not always perfectly reflect		
568	the underlying scoring process. We therefore include		
569	faithfulness checks and qualitative analyses;		
570	nonetheless, fully guaranteeing reasoning faithfulness		
571	remains an open challenge for all reasoning		
572	LLMs.		
573	<b>Broader Impact</b>		
574	We acknowledge there might be several ethical considerations		
575	in MT evaluation research such as gender bias.		
576	We prioritize high-quality open-sourced data and model		
577	in this research. We acknowledge that automatic metrics		
578	can reflect societal biases that is in training data,		
579	such as human biases.		
580	<b>References</b>		
581	Josh Achiam, Steven Adler, Sandhini Agarwal, Lama		
582	Ahmad, Ilge Akkaya, Florencia Leoni Aleman,		
583	Diogo Almeida, Janko Altenschmidt, Sam Altman,		
584	Shyamal Anadkat, et al. 2023. Gpt-4 technical report.		
585	<i>arXiv preprint arXiv:2303.08774</i> .		
586	Eleftheria Briakou, Jiaming Luo, Colin Cherry, and		
587	Markus Freitag. 2024. Translating step-by-step: De-		
588	composing the translation process for improved translation		
589	quality of long-form texts. In <i>Proceedings of the Ninth</i>		
590	<i>Conference on Machine Translation</i> , pages 1301–1317.		
591			
592	Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,		
593	Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul		
594	Barham, Hyung Won Chung, Charles Sutton, Sebastian		
595	Gehrmann, et al. 2023. Palm: Scaling language		
596	modeling with pathways. <i>Journal of Machine Learning</i>		
597	<i>Research</i> , 24(240):1–113.		
598	Marta R Costa-Jussà, James Cross, Onur Çelebi, Maha		
599	Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe		
600	Kalbassi, Janice Lam, Daniel Licht, Jean Maillard,		
601	et al. 2022. No language left behind: Scaling		
602	human-centered machine translation. <i>arXiv preprint</i>		
603	<i>arXiv:2207.04672</i> .		
604	Daniel Deutsch, Rotem Dror, and Dan Roth. 2021. A		
605	statistical analysis of summarization evaluation metrics		
606	using resampling methods. <i>Transactions of the</i>		
607	<i>Association for Computational Linguistics</i> , 9:1132–		
608	1146.		
	Daniel Deutsch, George Foster, and Markus Freitag.	609	
	2023. Ties matter: Meta-evaluating modern metrics	610	
	with pairwise accuracy and tie calibration. In <i>Pro-</i>	611	
	<i>ceedings of the 2023 Conference on Empirical Meth-</i>	612	
	<i>ods in Natural Language Processing</i> , pages 12914–	613	
	12929.	614	
	Zhaopeng Feng, Shaosheng Cao, Jiahao Ren, Jiayuan	615	
	Su, Ruizhe Chen, Yan Zhang, Zhe Xu, Yao Hu,	616	
	Jian Wu, and Zuozhu Liu. 2025a. Mt-r1-zero:	617	
	Advancing llm-based machine translation via r1-	618	
	zero-like reinforcement learning. <i>arXiv preprint</i>	619	
	<i>arXiv:2504.10160</i> .	620	
	Zhaopeng Feng, Jiayuan Su, Jiamei Zheng, Jiahao Ren,	621	
	Yan Zhang, Jian Wu, Hongwei Wang, and Zuozhu	622	
	Liu. 2025b. M-MAD: Multidimensional multi-agent	623	
	debate for advanced machine translation evaluation.	624	
	In <i>Proceedings of the 63rd Annual Meeting of the</i>	625	
	<i>Association for Computational Linguistics (Volume</i>	626	
	<i>1: Long Papers)</i> , pages 7084–7107.	627	
	Patrick Fernandes, Daniel Deutsch, Mara Finkelstein,	628	
	Parker Riley, André FT Martins, Graham Neubig,	629	
	Ankush Garg, Jonathan H Clark, Markus Freitag,	630	
	and Orhan Firat. 2023. The devil is in the errors:	631	
	Leveraging large language models for fine-grained	632	
	machine translation evaluation. In <i>Proceedings of the</i>	633	
	<i>Eighth Conference on Machine Translation</i> , pages	634	
	1066–1083.	635	
	Markus Freitag, George Foster, David Grangier, Viresh	636	
	Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021.	637	
	Experts, errors, and context: A large-scale study of	638	
	human evaluation for machine translation. <i>Transac-</i>	639	
	<i>tions of the Association for Computational Linguistics</i> ,	640	
	9:1460–1474.	641	
	Markus Freitag, Nitika Mathur, Daniel Deutsch, Chi-	642	
	Kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian	643	
	Thompson, Frederic Blain, Tom Kocmi, Jiayi Wang,	644	
	David Ifeoluwa Adelani, Marianna Buchicchio,	645	
	Chrysoula Zerva, and Alon Lavie. 2024. <a href="#">Are LLMs</a>	646	
	<a href="#">breaking MT metrics? results of the WMT24 metrics</a>	647	
	<a href="#">shared task</a> . In <i>Proceedings of the Ninth Confer-</i>	648	
	<i>ence on Machine Translation</i> , pages 47–81, Miami,	649	
	Florida, USA. Association for Computational Lin-	650	
	guistics.	651	
	Markus Freitag, Nitika Mathur, Chi-kiu Lo, Elefthe-	652	
	rios Avramidis, Ricardo Rei, Brian Thompson, Tom	653	
	Kocmi, Frederic Blain, Daniel Deutsch, Craig Stew-	654	
	art, Chrysoula Zerva, Sheila Castilho, Alon Lavie,	655	
	and George Foster. 2023. <a href="#">Results of WMT23 metrics</a>	656	
	<a href="#">shared task: Metrics might be guilty but references</a>	657	
	<a href="#">are not innocent</a> . In <i>Proceedings of the Eighth Con-</i>	658	
	<i>ference on Machine Translation</i> , pages 578–628, Sin-	659	
	gapore. Association for Computational Linguistics.	660	
	Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo,	661	
	Craig Stewart, Eleftherios Avramidis, Tom Kocmi,	662	
	George Foster, Alon Lavie, and André F. T. Martins.	663	
	2022. <a href="#">Results of WMT22 metrics shared task: Stop</a>	664	
	<a href="#">using BLEU – neural metrics are better and more</a>	665	
	<a href="#">robust</a> . In <i>Proceedings of the Seventh Conference</i>	666	

667		on Machine Translation (WMT), pages 46–68, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.		723
668				724
669				725
670	Nuno M Guerreiro, Ricardo Rei, Daan Van Stigt, Luísa Coheur, Pierre Colombo, and André FT Martins.			726
671				727
672		2024. xcomet: Transparent machine translation evaluation through fine-grained error detection. <i>Transactions of the Association for Computational Linguistics</i> , 12:979–995.		728
673				729
674				730
675				731
676	Jingcheng Hu, Yinmin Zhang, Qi Han, Daxin Jiang, Xiangyu Zhang, and Heung-Yeung Shum.			732
677		2025. Openreasoner-zero: An open source approach to scaling up reinforcement learning on the base model. <i>arXiv preprint arXiv:2503.24290</i> .		733
678				734
679				735
680				736
681	Juraj Juraska, Daniel Deutsch, Mara Finkelstein, and Markus Freitag.			737
682		2024. Metricx-24: The google submission to the wmt 2024 metrics shared task. In <i>Proceedings of the Ninth Conference on Machine Translation</i> , pages 492–504.		738
683				739
684				740
685				741
686	Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag.			742
687		2023. Metricx-23: The google submission to the wmt 2023 metrics shared task. In <i>Proceedings of the Eighth Conference on Machine Translation</i> , pages 756–767.		743
688				744
689				745
690				746
691				747
692	Rebecca Knowles, Samuel Larkin, and Chi-Kiu Lo.			748
693		2024. Mslc24: Further challenges for metrics on a wide landscape of translation quality. In <i>Proceedings of the Ninth Conference on Machine Translation</i> , pages 475–491.		749
694				750
695				751
696				752
697	Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, et al.			753
698		2024. Findings of the wmt24 general machine translation shared task: The llm era is here but mt is not solved yet. In <i>Proceedings of the Ninth Conference on Machine Translation</i> , pages 1–46.		754
699				755
700				756
701				757
702				758
703				759
704				760
705	Tom Kocmi and Christian Federmann.			761
706		2023a. Gembamqm: Detecting translation quality error spans with gpt-4. In <i>Proceedings of the Eighth Conference on Machine Translation</i> , pages 768–775.		762
707				763
708				764
709	Tom Kocmi and Christian Federmann.			765
710		2023b. Large language models are state-of-the-art evaluators of translation quality. In <i>Proceedings of the 24th Annual Conference of the European Association for Machine Translation</i> , pages 193–203.		766
711				767
712				768
713				769
714	Tom Kocmi, Christian Federmann, Roman Grundkiewicz, Marcin Junczys-Dowmunt, Hitokazu Matsushita, and Arul Menezes.			770
715		2021. To ship or not to ship: An extensive evaluation of automatic metrics for machine translation. In <i>Proceedings of the Sixth Conference on Machine Translation</i> , pages 478–494.		771
716				772
717				773
718				774
719				775
720	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica.			776
721		2023. Efficient		777
722		memory management for large language model serving with pagedattention. In <i>Proceedings of the 29th Symposium on Operating Systems Principles</i> , pages 611–626.		778
				779
				780
				781
				782
				783
				784
				785
				786
				787
				788
				789
				790
				791
				792
				793
				794
				795
				796
				797
				798
				799
				800
				801
				802
				803
				804
				805
				806
				807
				808
				809
				810
				811
				812
				813
				814
				815
				816
				817
				818
				819
				820
				821
				822
				823
				824
				825
				826
				827
				828
				829
				830
				831
				832
				833
				834
				835
				836
				837
				838
				839
				840
				841
				842
				843
				844
				845
				846
				847
				848
				849
				850
				851
				852
				853
				854
				855
				856
				857
				858
				859
				860
				861
				862
				863
				864
				865
				866
				867
				868
				869
				870
				871
				872
				873
				874
				875
				876
				877
				878
				879
				880
				881
				882
				883
				884
				885
				886
				887
				888
				889
				890
				891
				892
				893
				894
				895
				896
				897
				898
				899
				900

779	<i>Empirical Methods in Natural Language Processing</i> ,	Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan,	834
780	pages 3660–3674.	Lingfeng Shen, Benjamin Van Durme, Kenton Mur-	835
781	Ricardo Rei, José GC De Souza, Duarte Alves,	ray, and Young Jin Kim. 2024b. Contrastive pref-	836
782	Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova,	erence optimization: Pushing the boundaries of llm	837
783	Alon Lavie, Luisa Coheur, and André FT Martins.	performance in machine translation. In <i>International</i>	838
784	2022. Comet-22: Unbabel-ist 2022 submission for	<i>Conference on Machine Learning</i> , pages 55204–	839
785	the metrics shared task. In <i>Proceedings of the Sev-</i>	55224. PMLR.	840
786	<i>enth Conference on Machine Translation (WMT)</i> ,	Qwen An Yang, Baosong Yang, Beichen Zhang,	841
787	pages 578–585.	Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan	842
788	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon	Li, Dayiheng Liu, Fei Huang, Guanting Dong, Hao-	843
789	Lavie. 2020. <b>COMET: A neural framework for MT</b>	ran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei	844
790	<b>evaluation</b> . In <i>Proceedings of the 2020 Conference</i>	Zhang, Jianxin Yang, Jiaxin Yang, Jingren Zhou, Jun-	845
791	<i>on Empirical Methods in Natural Language Pro-</i>	yang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin	846
792	<i>cessing (EMNLP)</i> , Online. Association for Computa-	Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin	847
793	tional Linguistics.	Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia,	848
794	John Schulman, Filip Wolski, Prafulla Dhariwal,	Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang	849
795	Alec Radford, and Oleg Klimov. 2017. Proxi-	Su, Yi-Chao Zhang, Yuyang Wan, Yuqi Liu, Zeyu	850
796	mal policy optimization algorithms. <i>arXiv preprint</i>	Cui, Zhenru Zhang, Zihan Qiu, Shanghaoran Quan,	851
797	<i>arXiv:1707.06347</i> .	and Zekun Wang. 2024. <b>Qwen2.5 technical report</b> .	852
798	Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin	<i>ArXiv</i> , abs/2412.15115.	853
799	Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin	Shijie Zhang, Renhao Li, Songsheng Wang, Philipp	854
800	Lin, and Chuan Wu. 2024. Hybridflow: A flexible	Koehn, Min Yang, and Derek F Wong. 2025. Hi-	855
801	and efficient rlhf framework. <i>arXiv preprint</i>	mate: A hierarchical multi-agent framework for	856
802	<i>arXiv:2409.19256</i> .	machine translation evaluation. <i>arXiv preprint</i>	857
803	Richard Sproat, Tianyu Zhao, and Llion Jones. 2025.	<i>arXiv:2505.16281</i> .	858
804	Transevalnia: Reasoning-based evaluation and rank-	Mao Zheng, Zheng Li, Bingxin Qu, Mingyang	859
805	ing of translations. <i>arXiv preprint</i>	Song, Yang Du, Mingrui Sun, and Di Wang.	860
806	<i>arXiv:2507.12724</i> .	2025. Hunyuan- <i>mt</i> technical report. <i>arXiv preprint</i>	861
807	Shaomu Tan and Christof Monz. 2025. <b>ReMedy: Learn-</b>	<i>arXiv:2509.05209</i> .	862
808	<b>ing machine translation evaluation from human pref-</b>	<b>A Appendix</b>	863
809	<b>erences with reward modeling</b> . In <i>Proceedings of</i>	<b>A.1 RLVR Training Details</b>	864
810	<i>the 2025 Conference on Empirical Methods in Natu-</i>	We train Remedy-R using reinforcement learning	865
811	<i>ral Language Processing</i> , pages 4370–4387, Suzhou,	with the verifiable reward in §2.2. The model acts	866
812	China. Association for Computational Linguistics.	as a policy that, given an input, generates reasoning	867
813	Brian Thompson, Nitika Mathur, Daniel Deutsch, and	steps and final scores, and then receives a scalar re-	868
814	Huda Khayrallah. 2024. Improving statistical sig-	ward that reflects alignment with human judgments.	869
815	nificance in human evaluation of automatic metrics	This objective is to update the model so that high-	870
816	via soft pairwise accuracy. In <i>Proceedings of the</i>	reward behaviors become more likely over time.	871
817	<i>Ninth Conference on Machine Translation</i> , pages	Formally, the model defines a conditional policy	872
818	1222–1234.	$\pi_\theta(\mathbf{y} \mid \mathbf{x})$ over output sequences. After generating	873
819	Marcos Treviso, Nuno Guerreiro, Sweta Agrawal, Ri-	a response, we compute a scalar reward and maxi-	874
820	cardo Rei, José Pombal, Tânia Vaz, Helena Wu, Beat-	mize the expected return with gradients estimated	875
821	riz Silva, Daan Stigt, and André FT Martins. 2024.	via the policy gradient theorem:	876
822	xtower: A multilingual llm for explaining and correct-	$\mathcal{J}(\theta) = \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \pi_\theta} [r(\mathbf{y}, \mathbf{x})],$	
823	ing translation errors. In <i>Findings of the Association</i>	$\nabla_\theta \mathcal{J}(\theta) = \mathbb{E} \left[ \sum_{t=1}^T \nabla_\theta \log \pi_\theta(y_t \mid y_{<t}, \mathbf{x}) A_t \right].$	
824	<i>for Computational Linguistics: EMNLP 2024</i> , pages	(6)	877
825	15222–15239.	where $A_t$ denotes the token-level advantage.	878
826	Di Wu, Seth Aycock, and Christof Monz. 2025. Please	This estimator increases the likelihood of high-	879
827	translate again: Two simple experiments on whether	reward generations and decreases that of low-	880
828	human-like reasoning helps translation. <i>arXiv</i>	reward ones, aligning the model’s behavior with	881
829	<i>preprint arXiv:2506.04521</i> .		
830	Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Has-		
831	san Awadalla. 2024a. A paradigm shift in machine		
832	translation: Boosting translation performance of		
833	large language models. In <i>The Twelfth International</i>		
	<i>Conference on Learning Representations</i> .		

## Remedy-R Training Prompt Template

You are an expert machine translation evaluator. You need to assess the quality of two translations of the same source text. Your task is to evaluate the translation quality and provide scores from 0–100, where higher scores indicate better quality. You are also given a reference (not always perfect) to help you assess the quality.

### Evaluation Criteria:

- Accuracy: Whether the meaning expressed in the translation is correct and faithful to the source. Penalize mistranslation, unsupported additions/hallucinations, terminology errors, and untranslated text.
- Fluency: How natural and grammatically correct the translation reads in the target language. Consider grammar, agreement, word order, punctuation, spelling, register.
- Completeness: Is all information from the source conveyed without omissions?

**Instructions:** Think step by step about the quality of each translation and write your analysis first, then provide your final scores. Evaluate each translation independently rather than by comparison.

**Output Format:** Thinking through your evaluation first, then output the scores in exactly this format (do not give scores first):

#### A: [score] B: [score]

Now evaluate this \$SRC\_LANG-\$TGT\_LANG translation:

—

Source: \$SOURCE

Reference: \$REFERENCE

Translation A: \$TRANSLATION\_A

Translation B: \$TRANSLATION\_B

882 the reward signal. We optimize this objective using  
883 Proximal Policy Optimization (PPO) (Schulman  
884 et al., 2017), which stabilizes updates by clipping  
885 the ratio between the new and the old policy, while  
886 regularizing against a frozen reference model  $\pi_{\theta_{\text{ref}}}$   
887 to prevent reward over-optimization. During train-  
888 ing, each rollout prompts the model with a transla-  
889 tion pair, generates an evaluation output, parses the  
890 scores, computes the reward, estimates token level  
891 advantages, and finally applies PPO updates. We  
892 then optimize the following PPO objective:

$$\begin{aligned} \mathcal{L}_{\text{PPO}}(\theta) &= \mathbb{E}_t[\ell_t(\theta)] - \beta_{\text{KL}} \text{KL}(\pi_\theta \parallel \pi_{\text{ref}}), \\ \ell_t(\theta) &= \min\left(r_t(\theta)A_t, \text{clip}_\epsilon(r_t(\theta))A_t\right). \end{aligned} \quad (7)$$

894 Here,  $r_t(\theta) = \frac{\pi_\theta(y_t|y_{<t}, \mathbf{x})}{\pi_{\theta_{\text{old}}}(y_t|y_{<t}, \mathbf{x})}$  is the likelihood  
895 ratio,  $\epsilon$  is the clipping threshold, and  $\beta_{\text{KL}}$  controls  
896 the strength of the KL penalty that regularizes the  
897 updated policy toward a frozen reference policy  
898  $\pi_{\theta_{\text{ref}}}$  (the corresponding pretrained base model of  
899 the same size). Rewards are terminal and defined  
900 on the final parsed score. For advantage estimation,  
901 we use Generalized Advantage Estimation (GAE)  
902 with  $\lambda = 1$ , which trades off bias and variance by  
903 exponentially weighting multi step returns:

$$A_t = \sum_{l=0}^{\infty} (\gamma\lambda)^l \delta_{t+l}, \quad (8)$$

$$\delta_t = r_t + \gamma V_\phi(s_{t+1}) - V_\phi(s_t).$$

905 where  $V_\phi$  is the value function, and  $\gamma \in (0, 1]$

906 and  $\lambda \in [0, 1]$  are standard discounting hyperpa-  
907 rameters. Following recent work (Hu et al., 2025),  
908 we set  $\lambda = 1$ , which simplifies the estimator to a  
909 discounted Monte Carlo return and improves sta-  
910 bility in practice.

## A.2 Ablation study on WMT22 Metric Benchmark

911 We conduct an ablation study on the WMT22 met-  
912 ric benchmark to isolate the effect of our reward  
913 design during RLVR training. Unless otherwise  
914 specified, all ablations in this subsection are per-  
915 formed with **Remedy-R 7B**. We compare two re-  
916 ward settings: (1) a pure pairwise ranking reward  
917 that only verifies whether the model-implied pref-  
918 erence matches the human preference label; and (2)  
919 adding Huber reward shaping to encourage better  
920 calibration and reduce overly discrete score behav-  
921 iors.

922 As shown in Table 5, adding Huber reward shap-  
923 ing consistently improves accuracy at both the sys-  
924 tem and segment levels. In particular, Huber shap-  
925 ing improves system-level accuracy from 87.6% to  
926 89.1% and segment-level accuracy from 52.7% to  
927 54.8%, yielding a +1.7% absolute gain in average  
928 accuracy. These results suggest that Huber shap-  
929 ing provides a more informative learning signal  
930 beyond binary preference verification and leads to  
931 better-aligned quality judgments.

932 We further explore incorporating an explanation  
933 quality penalty into the reward. Instead of training  
934 a separate reward model, we reuse the same base  
935  
936

Model	Reward Setting	Sys	Seg	Avg
7B	Pairwise Ranking Reward	87.6%	52.7%	70.2%
	+ Huber Reward Shaping	89.1%	54.8%	71.9%
14B	Pairwise + Huber	88.7%	56.0%	72.4%
	Pairwise + Huber + genRM penalty	89.8%	56.3%	73.1%

Table 5: Ablation study on WMT22 comparing reward designs. We report accuracy at the system and segment levels, and their average.

instruct model (the RL initialization) as a generative rationale judge (genRM). Given the source, translation, and the model-generated explanation, genRM produces a scalar score intended to reflect the rationale’s faithfulness and relevance. We subtract this score (scaled by a fixed coefficient) from the RL reward, penalizing low-quality explanations while keeping the pairwise preference verification term unchanged.

We provide the reward dynamics during training in Figure 2. As shown in Table 5, adding this penalty results in only a marginal change in accuracy in our Remedy-R 14B setting. Overall, this suggests that the main gains come from the pairwise reward with Huber shaping, while explanation-based regularization provides at most a small additional benefit under our current design.

### A.3 Test Time Scaling with multiple Evaluation Passes

Remedy-R’s generative reasoning nature enables the application of **Test-Time Scaling (TTS)**, where multiple evaluation passes are performed with different reasoning trajectories and their quality scores are aggregated. In this setting, we adopt a simple implementation that averages the quality scores from multiple independent evaluations.

As shown in Figure 3, performing more evaluation trajectories at test time consistently enhances performance across all model sizes. Notably, Remedy-R-14B reaches an average correlation of 74.9, matching the strongest GEMBA-MQM performance. The steady improvement from 7B to 32B suggests that iterative reasoning stabilizes evaluation outcomes and reduces stochastic variance, yielding more robust and reliable quality assessments.

Interestingly, we observe that TTS primarily improves *segment-level*  $acc_{eq}^*$  rather than system-level correlation (see Table 12). We hypothesize that this phenomenon is mostly due to the limitations of current meta-evaluation metrics. As noted by Per-

rella et al. (2024), tie-calibrated pairwise accuracy ( $acc_{eq}^*$ ) tends to favor metrics that output continuous rather than discrete scores. Averaging multiple predictions effectively smooths discrete outputs into continuous scores, improving agreement with tie-calibrated accuracy, which favors metrics with finer score granularity.

### A.4 WMT23 Metric Benchmark

We report results on the WMT23 MQM metric shared task in Table 12. We compare Remedy-R against established reference-based metrics (e.g., XCOMET-XXL, MetricX-23), QE variants, and recent LLM-based judges. For Remedy-R, we additionally evaluate test-time scaling (TTS) by sampling multiple reasoning outputs and aggregating their predicted scores ( $TTS=k$ ). We report both system-level accuracy (**Acc**) and segment-level  $acc_{eq}^*$ , together with the average of the two.

As shown in Table 12, increasing TTS generally improves  $acc_{eq}^*$  more consistently than **Acc**, leading to steady gains in the overall average. This trend is most pronounced for smaller models (7B), where TTS closes a substantial portion of the gap to stronger baselines, and remains beneficial for 14B and 32B as well. These results support our claim that sampling-based aggregation helps mitigate discrete score behaviors and yields more stable segment-level judgments.

### A.5 WMT24 Metric Benchmark

Methods	Rank	Avg corr	Sys SPA	Seg $acc_{eq}^*$
ReMedy <sub>9B-24</sub>	1	<b>72.9</b>	85.9	<b>60.0</b>
Remedy-R-14B	2	72.6	<b>87.9</b>	57.2
MetricX-24-Hybrid	3	72.1	85.6	58.5
XCOMET-XXL	4	71.9	86.1	57.6
MetricX-24-Hybrid-QE	5	71.4	84.9	58.0
GEMBA-ESA (GPT4)	6	71.1	84.6	57.6
Remedy-R-32B	7	70.6	85.3	56.2
Remedy-R-7B	8	70.6	85.8	55.4
XCOMET-XXL-QE	9	69.5	83.3	55.7
Skywork-RM	10	69.0	83.2	54.7

Table 6: Evaluation on WMT24 MQM set. We report the official accuracy percentage (SPA and  $acc_{eq}^*$ ).

Table 6 reports results on the WMT24 MQM metric benchmark. Following the official shared

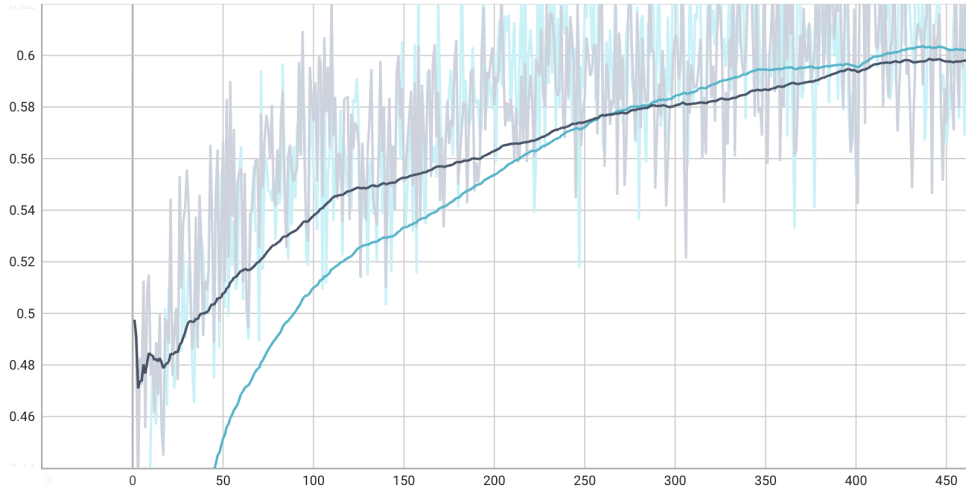


Figure 2: Reward curves during training for Pairwise + Huber (dark-blue) and Pairwise + Huber + genRM penalty (light-blue) settings.

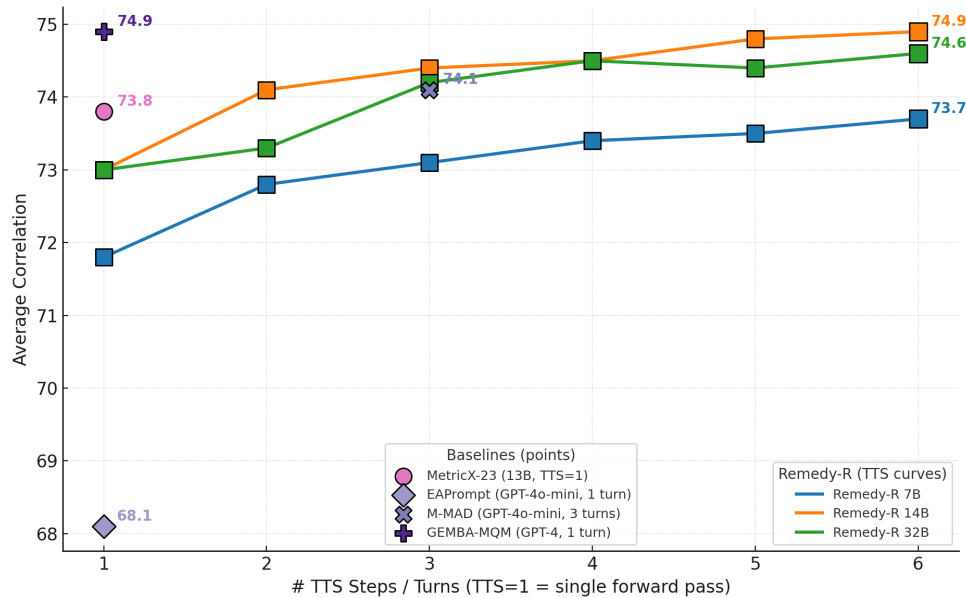


Figure 3: Average correlation across WMT23 MQM benchmarks under different numbers of Test-Time Scaling (TTS) evaluation passes. Each configuration aggregates multiple independent evaluations by averaging their final quality scores. TTS consistently improves correlation as the number of evaluation passes increases. Full results are shown in Table 12 in Appendix.

task protocol, we report system-level accuracy (SPA) and segment-level  $acc_{eq}^*$ , as well as their average correlation score (Avg corr) and the corresponding rank.

Remedy-R achieves competitive performance on WMT24 across model scales. In particular, Remedy-R-14B ranks second overall and attains the best system-level accuracy (SPA). While Remedy-R-32B and Remedy-R-7B show slightly lower segment-level  $acc_{eq}^*$  than top baselines, their overall average remains comparable to strong reference-based metrics such as MetricX-

24-Hybrid and XCOMET-XXL.

## A.6 GPT-4o-mini Faithfulness Prompt

We use the following prompt to assess the faithfulness of Remedy-R’s evaluation explanations. GPT-4 is given the source sentence, the translation hypothesis, and the explanation only (no reference translation), and returns a JSON object containing a faithfulness\_score and a brief reason.

System:

	ref?	empty mt	empty src+ref	src copy	wrong lang	mix lang	unrelated mt
COMET-22	✓	57.00%	58.81%	69.85%	67.84%	65.56%	45.23%
KIWI	✗	54.87%	67.72%	52.15%	82.64%	78.75%	41.95%
XCOMET	✓	73.79%	64.12%	82.04%	85.65%	71.77%	20.31%
MetricX-24-XXL	✓	-9.59	-5.85	-12.59	-3.06	-10.08	-24.15
MetricX-24-XXL	✗	-7.34	-5.85	-11.36	-2.51	-7.78	-24.25
GEMBA-ESA	✗	14.00%	13.5%	11.12%	14.32%	18.08%	1.27%
ReMedy-R-7B	✓	1.00%	7.07%	76.92%	43.69%	60.6%	1.5%
ReMedy-R-7B	✗	0.83%	5.40%	90.38%	33.15%	65.4%	2.0%
ReMedy-R-14B	✓	0.00%	0.00%	11.35%	14.6%	37.6%	0.6%
ReMedy-R-14B	✗	0.00%	0.00%	28.07%	12.1%	35.5%	1.0%
ReMedy-R-32B	✓	0.00%	0.00%	2.76%	8.30%	46.0%	1.3%
ReMedy-R-32B	✗	0.00%	0.00%	1.30%	8.0%	43.8%	3.5%

Table 7: Averaged quality scores of different metric models on MSLC24 OOD set. For all classes except *mix-lang*, a robust metric should output low scores; for *mix-lang*, the translation preserves the source meaning but contains code-switching, so its quality scores should be moderately high rather than near zero. MetricX scores are ranged from -25 to 0.

You are a strict verifier. Your job is to score the FAITHFULNESS of an evaluation explanation. You must return ONLY a single valid JSON object and nothing else.

User:

You are given: 1) `src_sent`: the source sentence; 2) `target_sent`: the translation hypothesis; 3) `explanation`: an evaluation text that comments on the translation quality

Task: Provide a score (0-100) indicating how FAITHFUL the explanation is to `src_sent` and `target_sent`.

Definition of faithfulness:

- Every key claim in the explanation must be supported by what is actually present in `src_sent` and/or `target_sent`.
- If the explanation invents content, mentions errors that are not evidenced, misquotes words, or contradicts `src/target`, score lower.

CRITICAL:

- You are NOT evaluating translation quality.
- A translation can be very bad, but an explanation can still be highly faithful if it correctly describes that badness.

Return ONLY JSON with:

```
{"faithfulness_score": int, // 0-100
 "brief_reason": string // <= 40 words,
 cite the biggest supported or unsupported
 claim}
```

Input:

```
src_lang: <src_lang>; tgt_lang: <tgt_lang>;
src_sent: <src>; target_sent: <tgt>; explanation:
<explanation>
```

## A.7 Remedy-R Agent experiments

Table 10 reports additional Remedy-R Agent results on WMT24 using strong open-source ALMA-

R translators as  $M_{base}$ . We evaluate two base settings: ALMA-R-7B with Remedy-R-7B, and ALMA-R-13B with Remedy-R-14B. For each language direction, we report multiple automatic metrics, including SacreBLEU, XCOMET-XXL, MetricX-24-XXL, and Remedy-R’s own scores. We include both the initial translations (row “- / -”) and the refined outputs produced by the Remedy-R Agent (row “Remedy-R / Remedy-R”), where Remedy-R generates the explanation and performs the refinement based on its own feedback.

Overall, Remedy-R Agent consistently improves translation quality across language pairs, with particularly large gains on en-zh. These results complement the main agent experiments and further demonstrate that Remedy-R’s reasoning can be reused to drive refinement even when the initial translations come from a strong external MT system.

## A.8 Remedy-R examples

We further provide several examples of the Remedy-R-14B model, including Chinese-English, English-Japanese, English-German, English-Czech, English-Spanish.

1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057

1031  
1032  
1033

$M_{feedback}$	$M_{refinement}$	cs-uk	en-cs	en-de	en-es	en-hi	en-is	en-ja	en-ru	en-uk	en-zh	ja-zh	Avg
$M_{base} = Qwen2.5-it-7B \mid Remedy-R = 7B \mid x-Tower = 14B$ (w. $XComet-XL = 3.5B$ )													
-	-	62.9	53.1	86.0	81.9	37.9	29.8	69.7	71.2	51.5	82.8	69.9	63.4
-	Base	64.3 $\uparrow_{1.4}$	55.0 $\uparrow_{1.9}$	86.4 $\uparrow_{0.4}$	82.1 $\uparrow_{0.2}$	39.1 $\uparrow_{1.2}$	30.7 $\uparrow_{0.9}$	69.6 $\downarrow_{0.1}$	72.1 $\uparrow_{0.9}$	53.8 $\uparrow_{2.3}$	82.7 $\downarrow_{0.1}$	70.4 $\uparrow_{0.5}$	64.2 $\uparrow_{0.8}$
x-Tower	x-Tower	75.6 $\uparrow_{12.6}$	57.0 $\uparrow_{3.9}$	90.7 $\uparrow_{4.7}$	85.8 $\uparrow_{3.9}$	40.3 $\uparrow_{2.4}$	33.6 $\uparrow_{3.7}$	62.3 $\downarrow_{7.4}$	79.3 $\uparrow_{8.0}$	69.2 $\uparrow_{17.7}$	82.2 $\downarrow_{0.6}$	66.5 $\downarrow_{3.4}$	67.5 $\uparrow_{4.2}$
Remedy-R	Remedy-R	66.0 $\uparrow_{3.1}$	56.0 $\uparrow_{2.9}$	87.2 $\uparrow_{1.1}$	83.2 $\uparrow_{1.3}$	40.5 $\uparrow_{2.6}$	30.9 $\uparrow_{1.1}$	71.8 $\uparrow_{2.1}$	73.4 $\uparrow_{2.1}$	54.6 $\uparrow_{3.1}$	83.2 $\uparrow_{0.5}$	69.7 $\downarrow_{0.2}$	65.1 $\uparrow_{1.8}$
x-Tower	Base	66.3 $\uparrow_{3.4}$	55.8 $\uparrow_{2.6}$	87.3 $\uparrow_{1.3}$	83.0 $\uparrow_{1.1}$	39.4 $\uparrow_{1.5}$	30.3 $\uparrow_{0.5}$	70.4 $\uparrow_{0.7}$	74.2 $\uparrow_{2.9}$	56.0 $\uparrow_{4.4}$	83.0 $\uparrow_{0.2}$	69.7 $\downarrow_{0.2}$	65.0 $\uparrow_{1.7}$
Remedy-R	Base	65.8 $\uparrow_{2.9}$	55.9 $\uparrow_{2.8}$	86.9 $\uparrow_{0.9}$	82.9 $\uparrow_{1.0}$	40.4 $\uparrow_{2.5}$	30.5 $\uparrow_{0.7}$	71.5 $\uparrow_{1.8}$	73.3 $\uparrow_{2.1}$	54.4 $\uparrow_{2.9}$	83.0 $\uparrow_{0.3}$	69.3 $\downarrow_{0.6}$	64.9 $\uparrow_{1.6}$
$M_{base} = Qwen2.5-it-14B \mid Remedy-R = 14B \mid x-Tower = 14B$ (w. $XComet-XL = 3.5B$ )													
-	-	69.4	63.6	88.4	83.7	47.6	32.6	74.8	75.7	58.6	83.9	72.5	68.2
-	Base	71.0 $\uparrow_{0.6}$	67.2 $\uparrow_{3.6}$	89.4 $\uparrow_{1.0}$	85.0 $\uparrow_{1.3}$	51.5 $\uparrow_{3.9}$	34.4 $\uparrow_{1.9}$	77.9 $\uparrow_{3.1}$	77.9 $\uparrow_{2.2}$	63.2 $\uparrow_{4.6}$	84.5 $\uparrow_{0.6}$	72.7 $\uparrow_{0.2}$	70.4 $\uparrow_{2.2}$
x-Tower	x-Tower	77.1 $\uparrow_{7.7}$	62.5 $\downarrow_{1.1}$	91.4 $\uparrow_{3.0}$	86.5 $\uparrow_{2.8}$	45.0 $\downarrow_{2.6}$	34.1 $\uparrow_{1.6}$	65.6 $\downarrow_{9.2}$	80.2 $\uparrow_{4.5}$	70.4 $\uparrow_{11.8}$	82.8 $\downarrow_{1.1}$	68.3 $\downarrow_{4.2}$	69.4 $\uparrow_{1.2}$
Remedy-R	Remedy-R	74.1 $\uparrow_{4.7}$	68.3 $\uparrow_{4.7}$	89.8 $\uparrow_{1.4}$	84.9 $\uparrow_{1.2}$	52.6 $\uparrow_{5.0}$	35.8 $\uparrow_{3.2}$	77.8 $\uparrow_{2.9}$	77.4 $\uparrow_{1.7}$	64.3 $\uparrow_{5.7}$	83.8 $\downarrow_{0.1}$	72.4 $\downarrow_{0.1}$	71.0 $\uparrow_{2.8}$
x-Tower	Base	73.0 $\uparrow_{3.6}$	68.5 $\uparrow_{4.9}$	90.2 $\uparrow_{1.8}$	85.3 $\uparrow_{1.7}$	51.3 $\uparrow_{3.7}$	34.6 $\uparrow_{2.0}$	78.6 $\uparrow_{3.8}$	78.5 $\uparrow_{2.8}$	65.1 $\uparrow_{6.5}$	84.3 $\uparrow_{0.4}$	71.9 $\downarrow_{0.6}$	71.0 $\uparrow_{2.8}$
Remedy-R	Base	74.8 $\uparrow_{5.4}$	68.2 $\uparrow_{4.6}$	90.3 $\uparrow_{1.9}$	85.6 $\uparrow_{1.9}$	53.1 $\uparrow_{5.4}$	36.3 $\uparrow_{3.7}$	79.3 $\uparrow_{4.5}$	79.2 $\uparrow_{3.5}$	64.5 $\uparrow_{5.9}$	84.1 $\uparrow_{0.2}$	72.4 $\downarrow_{0.1}$	71.6 $\uparrow_{3.4}$
$M_{base} = Qwen2.5-it-32B \mid Remedy-R = 32B \mid x-Tower = 14B$ (w. $XComet-XL = 3.5B$ )													
-	-	74.0	68.2	89.6	84.6	53.4	35.1	78.0	76.9	64.0	83.5	72.6	70.9
-	Base	75.0 $\uparrow_{1.0}$	69.5 $\uparrow_{1.3}$	90.0 $\uparrow_{0.4}$	85.5 $\uparrow_{0.9}$	55.4 $\uparrow_{2.0}$	36.5 $\uparrow_{1.4}$	78.7 $\uparrow_{0.7}$	77.2 $\uparrow_{0.3}$	66.7 $\uparrow_{2.7}$	83.8 $\uparrow_{0.3}$	73.6 $\uparrow_{1.0}$	72.0 $\uparrow_{1.1}$
x-Tower	x-Tower	77.5 $\uparrow_{3.5}$	63.5 $\downarrow_{4.7}$	91.7 $\uparrow_{2.1}$	86.1 $\uparrow_{1.5}$	45.5 $\downarrow_{7.9}$	34.2 $\downarrow_{0.9}$	66.7 $\downarrow_{11.2}$	80.2 $\uparrow_{3.2}$	71.6 $\uparrow_{7.6}$	82.5 $\downarrow_{1.0}$	68.2 $\downarrow_{4.5}$	69.8 $\downarrow_{1.1}$
Remedy-R	Remedy-R	76.6 $\uparrow_{2.6}$	71.5 $\uparrow_{3.3}$	91.1 $\uparrow_{1.5}$	85.6 $\uparrow_{1.0}$	57.1 $\uparrow_{3.7}$	37.4 $\uparrow_{2.2}$	80.4 $\uparrow_{2.5}$	78.5 $\uparrow_{1.6}$	67.8 $\uparrow_{3.8}$	84.0 $\uparrow_{0.6}$	73.3 $\uparrow_{0.7}$	73.0 $\uparrow_{2.1}$
x-Tower	Base	75.6 $\uparrow_{1.6}$	70.4 $\uparrow_{2.2}$	91.2 $\uparrow_{1.6}$	86.1 $\uparrow_{1.5}$	56.1 $\uparrow_{2.7}$	36.5 $\uparrow_{1.4}$	80.5 $\uparrow_{2.6}$	80.3 $\uparrow_{3.3}$	69.1 $\uparrow_{5.1}$	84.8 $\uparrow_{1.3}$	72.7 $\uparrow_{0.1}$	73.0 $\uparrow_{2.1}$
Remedy-R	Base	76.9 $\uparrow_{2.9}$	71.5 $\uparrow_{3.3}$	91.0 $\uparrow_{1.4}$	85.9 $\uparrow_{1.3}$	56.8 $\uparrow_{3.4}$	37.2 $\uparrow_{2.1}$	81.1 $\uparrow_{3.2}$	79.4 $\uparrow_{2.5}$	67.8 $\uparrow_{3.8}$	83.9 $\uparrow_{0.5}$	73.4 $\uparrow_{0.8}$	73.2 $\uparrow_{2.3}$

Table 8: Agent MT experiments on WMT24 benchmark using Qwen2.5 series models as the initial  $M_{base}$  translators (gray background). We report XCOMET-XXL (ref-based) in this table.

$M_{feedback}$	$M_{refinement}$	cs-uk	en-cs	en-de	en-es	en-hi	en-is	en-ja	en-ru	en-uk	en-zh	ja-zh	Avg
$M_{base} = Qwen2.5-it-7B \mid Remedy-R = 7B \mid x-Tower = 14B$ (w. $XComet-XL = 3.5B$ )													
-	-	-8.5	-11.1	-3.6	-4.1	-9.5	-21.2	-6.0	-5.8	-10.2	-2.8	-4.0	-7.9
xTower	xTower	-5.9 $\uparrow_{2.6}$	-9.6 $\uparrow_{1.4}$	-2.2 $\uparrow_{1.4}$	-3.2 $\uparrow_{1.0}$	-10.3 $\downarrow_{0.8}$	-19.6 $\uparrow_{1.6}$	-5.9 $\uparrow_{0.1}$	-4.2 $\uparrow_{1.7}$	-6.2 $\uparrow_{3.9}$	-2.9 $\downarrow_{0.1}$	-4.6 $\downarrow_{0.6}$	-6.8 $\uparrow_{1.1}$
Remedy-R	Remedy-R	-7.8 $\uparrow_{0.7}$	-10.1 $\uparrow_{1.0}$	-3.2 $\uparrow_{0.4}$	-3.8 $\uparrow_{0.4}$	-8.8 $\uparrow_{0.7}$	-21.1 $\uparrow_{0.2}$	-5.5 $\uparrow_{0.5}$	-5.2 $\uparrow_{0.6}$	-9.3 $\uparrow_{0.8}$	-2.7 $\uparrow_{0.1}$	-3.9 $\uparrow_{0.1}$	-7.4 $\uparrow_{0.5}$
$M_{base} = Qwen2.5-it-14B \mid Remedy-R = 14B \mid x-Tower = 14B$ (w. $XComet-XL = 3.5B$ )													
-	-	-7.0	-8.3	-2.8	-3.5	-7.1	-20.1	-5.4	-4.6	-8.3	-2.6	-3.7	-6.7
xTower	xTower	-5.6 $\uparrow_{1.4}$	-8.5 $\downarrow_{0.2}$	-2.1 $\uparrow_{0.8}$	-3.0 $\uparrow_{0.5}$	-9.1 $\downarrow_{2.0}$	-19.4 $\uparrow_{0.7}$	-5.6 $\downarrow_{0.2}$	-3.9 $\uparrow_{0.7}$	-5.9 $\uparrow_{2.4}$	-2.8 $\downarrow_{0.2}$	-4.4 $\downarrow_{0.8}$	-6.4 $\uparrow_{0.3}$
Remedy-R	Remedy-R	-6.0 $\uparrow_{1.0}$	-7.3 $\uparrow_{1.0}$	-2.3 $\uparrow_{0.5}$	-3.2 $\uparrow_{0.3}$	-6.3 $\uparrow_{0.8}$	-19.1 $\uparrow_{1.0}$	-4.8 $\uparrow_{0.5}$	-4.3 $\uparrow_{0.3}$	-7.0 $\uparrow_{1.2}$	-2.5 $\uparrow_{0.1}$	-3.6 $\uparrow_{0.1}$	-6.0 $\uparrow_{0.6}$
$M_{base} = Qwen2.5-it-32B \mid Remedy-R = 32B \mid x-Tower = 14B$ (w. $XComet-XL = 3.5B$ )													
-	-	-6.2	-7.3	-2.5	-3.4	-6.1	-19.2	-5.0	-4.6	-7.2	-2.6	-4.0	-6.2
xTower	xTower	-5.4 $\uparrow_{0.7}$	-8.3 $\downarrow_{0.9}$	-2.1 $\uparrow_{0.5}$	-3.0 $\uparrow_{0.4}$	-8.6 $\downarrow_{2.5}$	-18.8 $\uparrow_{0.4}$	-5.5 $\downarrow_{0.6}$	-3.9 $\uparrow_{0.7}$	-5.7 $\uparrow_{1.5}$	-2.8 $\downarrow_{0.1}$	-4.5 $\downarrow_{0.5}$	-6.2 $\downarrow_{0.0}$
Remedy-R	Remedy-R	-5.5 $\uparrow_{0.7}$	-6.6 $\uparrow_{0.7}$	-2.1 $\uparrow_{0.4}$	-3.0 $\uparrow_{0.3}$	-5.6 $\uparrow_{0.5}$	-18.3 $\uparrow_{0.9}$	-4.6 $\uparrow_{0.3}$	-4.0 $\uparrow_{0.6}$	-6.4 $\uparrow_{0.8}$	-2.5 $\uparrow_{0.1}$	-3.7 $\uparrow_{0.4}$	-5.7 $\uparrow_{0.5}$

Table 9: Agent MT experiments on WMT24 benchmark using Qwen2.5 series models as the initial  $M_{base}$  translators (gray background). We report MetricX24-XXL-Hybrid (ref-based) in this table.

$M_{feedback}$	$M_{refinement}$	en-cs	en-de	en-ru	en-zh	Avg	$M_{feedback}$	$M_{refinement}$	en-cs	en-de	en-ru	en-zh	Avg
$Base = ALMA-R-7B \mid Remedy-R = 7B$													
SacreBLEU							XCOMET-XXL						
-	-	16.6	22.2	14.1	23.9	19.2	-	-	71.9	89.5	77.3	75.7	78.6
Remedy-R	Remedy-R	17.5 $\uparrow_{0.9}$	22.2 $\uparrow_0$	15.8 $\uparrow_{1.7}$	27.4 $\uparrow_{3.5}$	20.7 $\uparrow_{1.5}$	Remedy-R	Remedy-R	71.4 $\downarrow_{0.5}$	90.4 $\uparrow_{0.9}$	78.5 $\uparrow_{1.3}$	79.9 $\uparrow_{4.2}$	80.0 $\uparrow_{1.4}$
MetricX-24-XXL							Remedy-R						
-	-	-6.2	-2.5	-4.5	-3.5	-4.2	-	-	92.3	93.4	92.5	86.5	91.2
Remedy-R	Remedy-R	-6.3 $\uparrow_{0.1}$	-2.4 $\uparrow_{0.2}$	-4.2 $\uparrow_{0.3}$	-3.0 $\uparrow_{0.5}$	-4.0 $\uparrow_{0.2}$	Remedy-R	Remedy-R	94.2 $\uparrow_{1.9}$	95.4 $\uparrow_{2.0}$	94.6 $\uparrow_{2.1}$	92.0 $\uparrow_{5.5}$	94.1 $\uparrow_{2.9}$
$Base = ALMA-R-13B \mid Remedy-R = 14B$													
SacreBLEU							XCOMET-XXL						
-	-	18.2	22.8	15.9	26.2	20.8	-	-	76.5	91.0	80.4	79.6	81.9
Remedy-R	Remedy-R	20.1 $\uparrow_{1.9}$	23.2 $\uparrow_{0.4}$	17.3 $\uparrow_{1.4}$	34.8 $\uparrow_{8.6}$	23.8 $\uparrow_{3.1}$	Remedy-R	Remedy-R	77.3 $\uparrow_{0.8}$	91.4 $\uparrow_{0.3}$	80.7 $\uparrow_{0.2}$	84.0 $\uparrow_{4.5}$	83.3 $\uparrow_{1.4}$
MetricX-24-XXL							Remedy-R						
-	-	-5.2	-2.1	-3.8	-3.1	-3.6	-	-	90.5	91.4	90.5	88.7	90.3
Remedy-R	Remedy-R	-5.3 $\downarrow_{0.1}$	-2.0 $\uparrow_{0.1}$	-3.8 $\uparrow_{0.1}$	-2.5 $\uparrow_{0.5}$	-3.4 $\uparrow_{0.2}$	Remedy-R	Remedy-R	92.1 $\uparrow_{1.6}$	93.1 $\uparrow_{1.7}$	91.9 $\uparrow_{1.4}$	93.2 $\uparrow_{4.5}$	92.6 $\uparrow_{2.3}$

Table 10: Agent MT experiments on WMT24 using ALMA-R models as  $M_{base}$ . We report initial translations (“- / -”) and Remedy-R Agent refinements (“Remedy-R / Remedy-R”) under multiple metrics. Left: SacreBLEU. Right: XCOMET-XXL. We additionally report MetricX-24-XXL and Remedy-R scores in the lower blocks.

Method	Step	en-cs	en-de	en-fr	en-ja	en-ru	en-zh	Avg
<b>gpt-4o-mini as base translator model</b>								
gpt-4o-mini	base	54.5	78.6	61.7	58.7	58.0	57.3	61.5
translate_again	1	61.4 $\uparrow 6.9$	81.9 $\uparrow 3.4$	68.4 $\uparrow 6.6$	64.3 $\uparrow 5.6$	64.6 $\uparrow 6.6$	63.5 $\uparrow 6.3$	67.4 $\uparrow 5.9$
translate_again	2	62.5 $\uparrow 7.9$	82.0 $\uparrow 3.4$	69.8 $\uparrow 8.1$	65.6 $\uparrow 7.0$	65.1 $\uparrow 7.1$	64.6 $\uparrow 7.3$	68.3 $\uparrow 6.8$
translate_again	3	62.9 $\uparrow 8.4$	82.5 $\uparrow 3.9$	69.7 $\uparrow 7.9$	65.5 $\uparrow 6.8$	65.6 $\uparrow 7.6$	64.5 $\uparrow 7.2$	68.4 $\uparrow 7.0$
step_by_step	1	56.6 $\uparrow 2.0$	78.7 $\uparrow 0.1$	64.2 $\uparrow 2.4$	59.1 $\uparrow 0.5$	59.8 $\uparrow 1.8$	58.5 $\uparrow 1.2$	62.8 $\uparrow 1.3$
step_by_step	2	60.7 $\uparrow 6.2$	80.8 $\uparrow 2.3$	67.9 $\uparrow 6.2$	62.5 $\uparrow 3.8$	63.0 $\uparrow 5.0$	62.5 $\uparrow 5.2$	66.2 $\uparrow 4.8$
step_by_step	3	61.2 $\uparrow 6.7$	81.2 $\uparrow 2.6$	68.2 $\uparrow 6.5$	62.5 $\uparrow 3.8$	63.6 $\uparrow 5.6$	62.5 $\uparrow 5.2$	66.5 $\uparrow 5.1$
Remedy-R	1	57.6 $\uparrow 3.1$	80.3 $\uparrow 1.7$	66.3 $\uparrow 4.6$	63.0 $\uparrow 4.4$	63.2 $\uparrow 5.2$	63.4 $\uparrow 6.2$	65.6 $\uparrow 4.2$
Remedy-R	2	58.0 $\uparrow 3.5$	81.3 $\uparrow 2.7$	67.2 $\uparrow 5.5$	64.5 $\uparrow 5.8$	64.1 $\uparrow 6.1$	64.0 $\uparrow 6.7$	66.5 $\uparrow 5.1$
Remedy-R	3	58.4 $\uparrow 3.8$	81.3 $\uparrow 2.7$	67.3 $\uparrow 5.6$	64.2 $\uparrow 5.5$	64.3 $\uparrow 6.3$	64.2 $\uparrow 7.0$	66.6 $\uparrow 5.2$
<b>gemini-2.0-flash as base translator model</b>								
gemini-2.0-flash	base	63.0	80.3	64.5	67.6	65.0	62.3	67.1
translate_again	1	68.7 $\uparrow 5.6$	84.4 $\uparrow 4.1$	70.8 $\uparrow 6.4$	71.2 $\uparrow 3.6$	71.1 $\uparrow 6.1$	70.0 $\uparrow 7.7$	72.7 $\uparrow 5.6$
translate_again	2	68.5 $\uparrow 5.5$	84.1 $\uparrow 3.8$	69.6 $\uparrow 5.2$	70.9 $\uparrow 3.3$	71.2 $\uparrow 6.2$	70.8 $\uparrow 8.4$	72.5 $\uparrow 5.4$
translate_again	3	68.7 $\uparrow 5.7$	84.1 $\uparrow 3.8$	68.6 $\uparrow 4.1$	70.2 $\uparrow 2.6$	70.7 $\uparrow 5.7$	70.3 $\uparrow 7.9$	72.1 $\uparrow 5.0$
step_by_step	1	64.5 $\uparrow 1.4$	81.0 $\uparrow 0.7$	65.6 $\uparrow 1.2$	66.3 $\downarrow 1.3$	65.6 $\uparrow 0.6$	62.0 $\downarrow 0.3$	67.5 $\uparrow 0.4$
step_by_step	2	68.9 $\uparrow 5.8$	84.2 $\uparrow 3.8$	69.7 $\uparrow 5.3$	70.0 $\uparrow 2.4$	70.8 $\uparrow 5.8$	68.4 $\uparrow 6.1$	72.0 $\uparrow 4.9$
step_by_step	3	69.3 $\uparrow 6.2$	84.5 $\uparrow 4.1$	70.3 $\uparrow 5.8$	70.4 $\uparrow 2.8$	71.0 $\uparrow 6.0$	69.5 $\uparrow 7.2$	72.5 $\uparrow 5.4$
Remedy-R	1	62.5 $\downarrow 0.5$	82.0 $\uparrow 1.6$	67.8 $\uparrow 3.3$	69.6 $\uparrow 2.0$	66.5 $\uparrow 1.5$	65.1 $\uparrow 2.7$	68.9 $\uparrow 1.8$
Remedy-R	2	62.8 $\downarrow 0.2$	82.3 $\uparrow 2.0$	67.6 $\uparrow 3.1$	69.9 $\uparrow 2.3$	66.6 $\uparrow 1.6$	65.9 $\uparrow 3.5$	69.2 $\uparrow 2.0$
Remedy-R	3	62.8 $\downarrow 0.3$	81.9 $\uparrow 1.5$	67.9 $\uparrow 3.4$	69.9 $\uparrow 2.3$	66.7 $\uparrow 1.6$	65.8 $\uparrow 3.5$	69.2 $\uparrow 2.0$

Table 11: Refinement performance comparison on the initial translations from GPT-4o-mini and Gemini-2.0-Flash using paragraph-level WMT24++ benchmark. We report ref-based XCOMET-XXL to measure the translation quality. Step here means refinement steps or iterations. For translate\_again and step\_by\_step, they adopt the self-refinement (using gpt-4o-mini and gemini-2.0-flash refinement), while Remedy-R utilizes 32B model Remedy-R model for refinement.

Method	$\theta$	TTS/Turns	Acc	$acc_{eq}^*$	Avg
KIWI-XXL	ensemble	1	91.1	54.6	72.9
MetricX-23	13B	1	90.7	56.9	73.8
MetricX-23-QE	13B	1	89.0	56.1	72.6
XCOMET-XXL	ensemble	1	92.8	57.7	75.3
XCOMET-XXL-QE	ensemble	1	91.6	55.8	73.7
ReMedy <sub>9B-23</sub>	9B	1	94.1	58.2	76.2
EAPrompt (GPT4o-mini)	-	1 turn	90.3	45.9	68.1
M-MAD (GPT4o-mini)	-	3 turns	94.5	53.7	74.1
GEMBA-MQM (GPT4)	-	1 turn	94.5	55.2	74.9
Remedy-R	7B	1 turn	93.7	49.9	71.8
Remedy-R	7B	TTS=2	94.1	51.5	72.8
Remedy-R	7B	TTS=3	93.7	52.6	73.1
Remedy-R	7B	TTS=4	93.7	53.1	73.4
Remedy-R	7B	TTS=5	93.7	53.4	73.5
Remedy-R	7B	TTS=6	93.7	53.6	73.7
Remedy-R	14B	1 turn	92.4	53.6	73.0
Remedy-R	14B	TTS=2	93.3	55.0	74.1
Remedy-R	14B	TTS=3	93.7	55.1	74.4
Remedy-R	14B	TTS=4	93.7	55.4	74.5
Remedy-R	14B	TTS=5	94.1	55.4	74.8
Remedy-R	14B	TTS=6	94.1	55.7	74.9
Remedy-R	32B	1 turn	94.5	51.5	73.0
Remedy-R	32B	TTS=2	94.1	52.6	73.3
Remedy-R	32B	TTS=3	94.5	53.9	74.2
Remedy-R	32B	TTS=4	94.9	54.1	74.5
Remedy-R	32B	TTS=5	94.5	54.3	74.4
Remedy-R	32B	TTS=5	95.0	54.3	74.6

Table 12: Evaluation on WMT MQM23 Metric Shared task. Both KIWI-XXL and XCOMET-XXL are identical ensembles of  $2 \times 10.7B$  and  $1 \times 3.5B$  models.

## Remedy-R (14B) ref-free Evaluation Example (Chinese-English)

**Source:** 麦乐鸡, 没有麦乐鸡酱

**Translation:** Wheat chicken, no maillot chicken sauce

### Remedy-R's Evaluation:

Let's evaluate the translation step by step based on the given criteria:

#### • Accuracy

- "麦乐鸡" (màilè jī) is a brand name for McDonald's Chicken McNuggets. The translation "Wheat chicken" is incorrect and does not convey the intended meaning.
- "麦乐鸡酱"(màilè jī jiàng) is the brand name for McDonald's McNugget sauce. The translation "maillot chicken sauce" is also incorrect and does not convey the intended meaning.

#### • Fluency

- "Wheat chicken" is not a natural or common English phrase. It sounds awkward and is not idiomatic.
- "Maillot" is a French word meaning "swimsuit" or "jersey," which is completely unrelated to the source text and makes no sense in this context.
- The phrase "no maillot chicken sauce" is not grammatically correct and is not a natural way to express the idea in English.

#### • Completeness

- The translation fails to accurately convey the meaning of both "麦乐鸡" and "麦乐鸡酱." It does not preserve the intended meaning of the source text.

Given the above evaluation, the translation is very poor in all three criteria: accuracy, fluency, and completeness. #####  
Score: 20 (0-100).

The translation is largely inaccurate, awkward, and fails to convey the intended meaning of the source text. A more appropriate translation would be "McNuggets, no McNugget sauce." However, since the translation provided is so far off, the score is quite low. A perfect translation would score 100, while a completely incorrect translation would score 0. Given the significant errors, a score of 20 is appropriate.

## Remedy-R (14B) ref-free Evaluation Example (English-Japanese)

**Source:** "Which notebook is that?" Ivory asked, sitting down next to Kari on her bed.

**Translation:** "「哪个笔记本是你的?」アイボリがカリの横に座ってした。「きみのnotebookはどれ?」"

### Remedy-R's Evaluation:

Let's evaluate the translation step by step based on the given criteria:

- **Accuracy: The translation attempts to convey the meaning of the source text, but there are some issues.**

- "Which notebook is that?" is translated as "哪个笔记本是你的?" which is in Chinese, not Japanese.
- "Ivory asked" is correctly translated as "アイボリがした".
- "sitting down next to Kari" is translated as "カリの横に座って", which is accurate.
- "on her bed" is missing in the translation.

- **Fluency**

- The translation is mostly fluent, but the inclusion of "notebook" in English ("notebook") is awkward and unnatural in Japanese.
- The phrase "哪个笔记本是你的?" is not in Japanese, which makes the translation less fluent.

- **Completeness**

- The translation is incomplete as it does not convey "on her bed".
- The inclusion of "notebook" in English is unnecessary and adds confusion.

- **Language Correctness**

- The translation mixes Chinese and Japanese, which is incorrect.
- "notebook" in English is not appropriate in a Japanese translation.

Given these points, the translation is not very good. Let's assign a score:

- Accuracy: 50 (partially accurate but with significant errors)
  - Fluency: 50 (awkward and unnatural due to language mix)
  - Completeness: 50 (missing "on her bed" and unnecessary "notebook" in English)
  - Language Correctness: 30 (mixing Chinese and Japanese)
- #### Score: 45 (0-100).

## Remedy-R (14B) ref-free Evaluation Example (English-German)

**Source:** "And with hooker Dan Frost sin-binned on the stroke of half-time, their 13 men were on the ropes and Munster reacted by engineering tries for debutant Patrick Campbell and Andrew Conway.

**Translation:** Und mit der Hure Dan Frost, die auf den Schlag der Halbzeit sündigte, waren ihre 13 Männer in den Seilen und Munster reagierte mit technischen Versuchen für den Debütanten Patrick Campbell und Andrew Conway.

### Remedy-R's Evaluation:

Let's evaluate the translation step by step according to the given criteria:

- **Accuracy: The translation attempts to convey the meaning of the source text, but there are some issues.**
  - "hooker" in rugby refers to a position, not a prostitute ("Hure"). The correct term would be "Flanker" or "Ruckspitz" in German
  - "sin-binned" is a rugby term meaning "sent off for a temporary suspension." The correct translation would be "gesperrt" or "gesperrt für die Zeitung."
  - "on the ropes" is an idiomatic expression meaning "in a difficult situation." A more appropriate translation would be "in Schwierigkeiten" or "am Ende ihrer Kräfte."
  - "engineering tries" is a metaphorical way of saying "scoring tries." The correct translation would be "Tore für" or "Versuche für."
- **Fluency: The translation is mostly fluent but contains several awkward phrases and incorrect vocabulary choices.**
  - "sündigte" is not the correct translation for "sin-binned."
  - "in den Seilen" is an incorrect translation for "on the ropes."
  - "technische Versuche" is not a standard rugby term and is confusing.
- **Completeness: The translation does not accurately convey the meaning of the source text.**
  - "sin-binned" is mistranslated, which changes the meaning significantly.
  - "on the ropes" is mistranslated, which also changes the meaning.
  - "engineering tries" is mistranslated, which is confusing.

Overall, the translation is quite poor and does not meet the criteria for a high-quality translation. #### Score: 50 (0-100). The translation has significant issues with accuracy and fluency, and it does not fully convey the meaning of the source text. A score of 50 reflects these major shortcomings.

## Remedy-R (14B) ref-free Evaluation Example (English-Czech)

**Source:** So either there's something wrong with this panel or my multimeter, or they changed the listing.

**Translation:** Takže buď je něco špatně s tímto panel  $\circ$  M, nebo s mým virometrem, nebo změnili popis.

### Remedy-R's Evaluation:

Let's evaluate the translation step by step according to the given criteria:

- **Accuracy: The translation attempts to convey the meaning of the source text, but there are some inaccuracies.**
  - "panel" is correctly translated as "panel".
  - "multimeter" is incorrectly translated as "virometrem" (which means "viro-meter" in Czech). The correct translation should be "multimetro".
  - "listing" is translated as "popis", which is a reasonable approximation of "listing" in this context, but "označení" or "popis produktu" might be more precise.
- **Fluency**
  - The translation is mostly fluent and reads naturally, but "virometrem" is a non-standard term and makes the sentence sound awkward.
  - The sentence structure is mostly correct, but "virometrem" disrupts the flow.
- **Completeness**
  - The translation attempts to convey all the information from the source, but the use of "virometrem" instead of "multimetro" is a significant inaccuracy.
  - The translation is mostly complete but loses some precision due to the incorrect term.
- **Language Correctness**
  - The translation is mostly correct in Czech, but "virometrem" is not a standard term and is incorrect.
  - The use of "panel  $\circ$  M" (with a Cyrillic "o") is a typographical error. It should be "panel".

Given these points, let's assign a score:

- **Accuracy:** 75 (due to the incorrect translation of "multimeter")
  - **Fluency:** 85 (mostly fluent but disrupted by "virometrem")
  - **Completeness:** 80 (mostly complete but loses precision)
  - **Language Correctness:** 85 (mostly correct but "virometrem" is incorrect)
- #### Score: 78 (0-100).

## Remedy-R (14B) ref-free Evaluation Example (English-Spanish)

**Source:** According to the office of the commissioner for the protection of the state language, 3,692 language law violations were reported in 2023. Kiev, Kharkov, and Odessa remain "leaders in terms of the number of complaints of violations. Therefore, strict controls over the observance of the language rights of citizens in all spheres of public life should be among the priorities of community leaders this year," Taras Kremen, the country's Language Ombudsman said.

**Translation:** Según la oficina del comisionado para la protección del idioma oficial del estado, se reportaron 3.692 violaciones de la ley lingüística en 2023. Kiev, Kharkov y Odessa continúan siendo los "líderes" en términos del número de quejas sobre violaciones. Por lo tanto, el control estricto sobre el cumplimiento de los derechos lingüísticos de los ciudadanos en todas las esferas de la vida pública debe ser una prioridad para los líderes comunitarios este año, dijo Taras Kremen, el Defensor del Lenguaje del país.

### Remedy-R's Evaluation:

Let's evaluate the translation step by step based on the given criteria:

#### • Accuracy:

- The translation accurately conveys the meaning of the source text, including the number of reported violations, the cities mentioned, and the statement by Taras Kremen.
- The phrase "leaders in terms of the number of complaints" is translated as "líderes en términos del número de quejas," which is a very accurate translation.

#### • Fluency

- The translation is mostly fluent and reads naturally.
- The sentence structure is clear and follows Spanish grammar rules.
- There are no awkward phrasings or unnatural constructions.

#### • Completeness

- All the information from the source text is conveyed accurately.
- There are no unnecessary additions or omissions.

#### • Language Correctness

- The translation is entirely in Spanish without any mixing of languages.
- The use of "Defensor del Lenguaje" for "Language Ombudsman" is a good choice, as it is a more formal and accurate term in Spanish.

Overall, the translation is very good, with only minor room for improvement in fluency and accuracy. ##### Score: 95 (0-100).