

Improving Low-Resource Machine Translation via Round-Trip Reinforcement Learning

Anonymous ACL submission

Abstract

Low-resource machine translation (MT) has gained increasing attention as parallel data from low-resource language communities is collected, but many potential methods for improving low-resource MT remain unexplored. We investigate a self-supervised reinforcement-learning-based fine-tuning for translation in low-resource settings using round-trip bootstrapping with the No Language Left Behind (NLLB) family of models. Our approach translates English into a target low-resource language and then back into English, using a combination of chrF++ and BLEU as the reward function on the reconstructed English sentences. Using the NLLB-MD dataset, we evaluate both the 600M and 1.3B parameter NLLB models and observe consistent improvements for the following languages: Central Aymara, Friulian, Wolof and Russian. Qualitative inspection of translation outputs indicates increased fluency and semantic fidelity. We argue that our method can further benefit from scale, enabling models to increasingly leverage their pretrained knowledge and continue self-improving.

1 Introduction

Low-resource machine translation remains one of the most persistent challenges in natural language processing. Modern neural translation systems achieve impressive performance for high-resource languages (Vaswani et al., 2017; Conneau et al., 2020; Hendy et al., 2023), where millions of parallel sentence pairs are available for supervised training. However, thousands of languages lack parallel corpora of sufficient scale (Koehn and Knowles, 2017; Haddow et al., 2022), leading to models that perform poorly or fail to generalize outside narrow training domains. As a result, speakers of low-resource languages continue to face limited access to digital information, reduced participation in multilingual technologies, and diminished inclusion in global communication systems.

To address these limitations, we introduce a framework to improve MT system without any parallel corpus that integrates reinforcement learning (RL) via round-trip bootstrapping. Our system operates in a self-supervised manner by translating from English into a low-resource language and then back to English. In this setup, the intermediate low-resource translation is treated as a latent step required to reconstruct the original English sentence. We use RL to optimize this process, rewarding the model based on lexical and semantic metrics such as BLEU and chrF++. This allows the model to improve translation quality in the low-resource language without requiring any human-annotated parallel data.

The necessity for this approach arises from the inherent flaws in standard machine translation training. Most systems rely on maximum likelihood estimation (MLE) (Fisher, 1992), which teaches models to predict the next token given ground-truth history. While effective, MLE suffers from *exposure bias* (Bengio et al., 2015; Wang and Sennrich, 2020; He et al., 2024), arising from discrepancies between training and inference conditions, and an *objective mismatch* (Maksai and Fua, 2018), where the training loss does not directly optimize evaluation metrics like chrF++ (Popović, 2015). In low-resource settings, these issues are amplified because models have fewer "gold" examples to prevent them from drifting during inference.

Furthermore, while back-translation (Sennrich et al., 2016) is currently the dominant method for low-resource MT, its effectiveness is fundamentally constrained by the quality of the reverse model. When initial translations are weak, the synthetic data produced by back-translation can reinforce existing model errors. This creates a feedback loop that limits gains, particularly for morphologically rich or typologically distant languages (McNamee and Duh, 2023). Standard self-supervised approaches lack a mechanism to prioritize high-

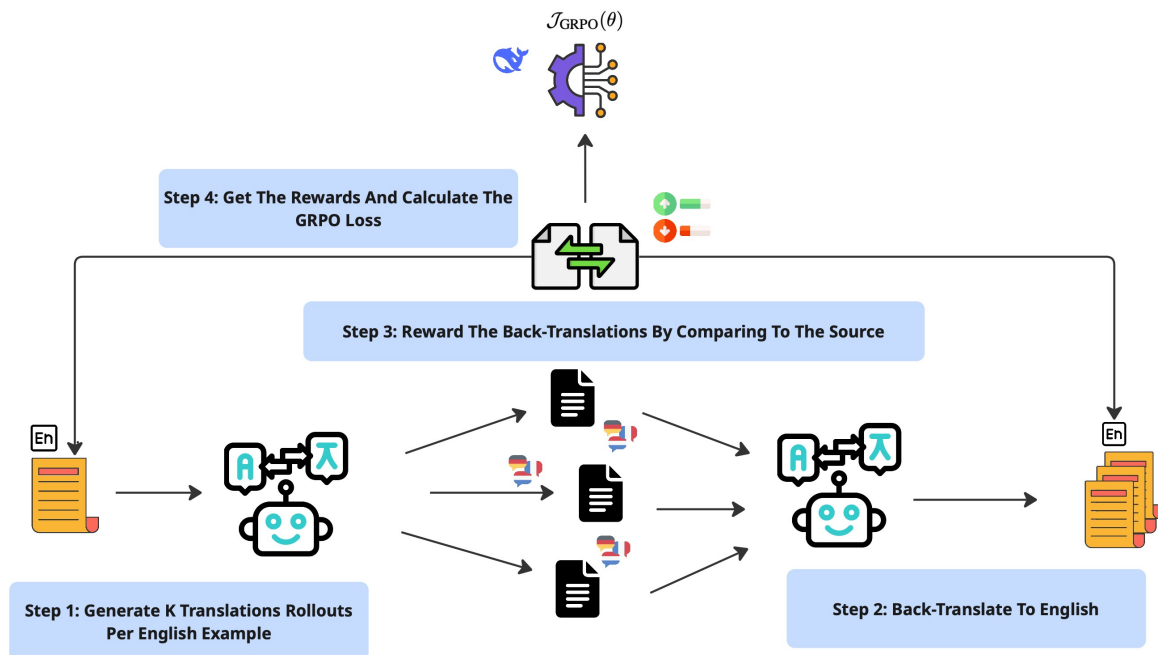


Figure 1: Overview Of Our Low-Resource Machine Translation With Round-Trip Reinforcement Learning Method.

quality synthetic examples or directly optimize task-specific translation metrics.

RL (Sutton and Barto, 2018) provides a natural solution to these bottlenecks. By treating translation as a sequential decision-making process, RL allows models to directly optimize non-differentiable metrics and exposes them to their own predictions during training. While foundational work introduced policy gradient methods (Ranzato et al., 2016; Wu et al., 2018) to reduce exposure bias, and recent studies have applied RL to large-scale models (Feng et al., 2025), the application of RL to bridge the specific gap in round-trip low-resource translation remains relatively under-explored.

In this work, we demonstrate that our RL-based bootstrapping method significantly improves both the adequacy and fluency of translations. Our main contributions are:

- We propose a novel RL-based framework for round-trip translation that operates entirely without parallel data.
- We show that optimizing for reconstruction rewards significantly improves surface-level well-formedness and semantic adequacy in low-resource settings.
- We provide an extensive analysis across several languages.

2 Related Work

2.1 Improving NMT with Monolingual Data and Roundtrip Translations

Back-translation has been a foundational technique for leveraging monolingual data in neural machine translation. Sennrich et al. (2016) introduced back-translation and demonstrated that target-side monolingual data can substantially improve translation quality by generating synthetic parallel data; mixing these synthetic examples with genuine bitext yields strong performance gains and supports effective domain adaptation. Lample et al. (2018) showed that unsupervised machine translation is feasible by combining denoising autoencoding, back-translation, and a cycle-consistency objective that enforces accurate reconstruction of source sentences after a forward-backward translation loop. This line of work treats forward and backward translation models as cooperative components whose outputs must remain mutually consistent. Huang (1990) is an early work that explicitly uses round-trip translation as a quality check for “target language inexpert” users: texts are translated into a foreign language and immediately back so users can spot major errors without knowing the target language. Federmann et al. (2019) extend this idea to paraphrasing, systematically studying “multilingual whispers,” where sentences are translated through one or more pivot languages and back to

generate diverse paraphrases, showing that round-trip translation can serve not only for quality control but also as an effective paraphrase generation method.

2.2 Reinforcement Learning and GRPO for NMT

Wu et al. (2018) provided a practical recipe for RL in NMT, showing that multinomial sampling, reward shaping, and objective mixing with MLE lead to consistent gains. Feng et al. (2025) proposed MT-R1-Zero, applying R1-Zero-style RL to MT without supervised warm-starts. Their rule-metric mixed reward improved lexical and semantic quality via GRPO (Shao et al., 2024), suggesting that pure RL can rival SFT for LLM-based translation.

GRPO extends Proximal Policy Optimization (PPO) (Schulman et al., 2017) and removes the value critic.

3 Method

We propose a reinforcement learning framework that leverages round-trip translation consistency to fine-tune pretrained multilingual language models for low-resource machine translation. Our approach combines Group Relative Policy Optimization (GRPO) with back-translation-based rewards, enabling effective adaptation without parallel data augmentation.

Preliminary For each input x , GRPO samples a group of G candidate outputs $\{y_1, \dots, y_G\} \sim \pi_{\theta_{\text{old}}}(y | x)$ from the current policy. The update maximizes

$$\mathcal{J}_{\text{GRPO}}(\theta) = E_{q \sim P(Q), \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(O|q)} \left[\frac{1}{G} \sum_{i=1}^G \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \min \left(r_{i,t}(\theta) \hat{A}_{i,t}, \text{clip} \left(r_{i,t}(\theta), 1 - \varepsilon, 1 + \varepsilon \right) \hat{A}_{i,t} \right) - \beta D_{\text{KL}}[\pi_{\theta} \parallel \pi_{\text{ref}}] \right],$$

where $r_{i,t}(\theta) = \frac{\pi_{\theta}(o_{i,t} | q, o_{i,<t})}{\pi_{\theta_{\text{old}}}(o_{i,t} | q, o_{i,<t})}$.

GRPO differs from PPO in how it defines the advantages A_i . Instead of subtracting a critic prediction, GRPO normalizes rewards within the sampled

group. Let $R_i = R(x, y_i)$; then

$$A_i = \frac{R_i - \text{mean}(\{R_1, \dots, R_G\})}{\text{std}(\{R_1, \dots, R_G\})}. \quad (2)$$

When all or none of the candidates solve the problem, every A_i is zero and the policy gradient vanishes except for the KL term.

3.1 Problem Formulation

We consider the setting where *no parallel corpus is available* for the source and target language pair. Instead, we assume access to a pretrained multilingual neural machine translation (NMT) model π_{θ} and monolingual corpora in the source language (English). Our goal is to improve the translation quality in both directions for this low-resource pair *without* relying on any parallel data for training.

We formulate this adaptation problem as a reinforcement learning problem. At each step, the model samples translation candidates and receives rewards that reflect translation quality, defined purely from intrinsic signals. Concretely, starting from a source sentence x , the model generates a target hypothesis $\tilde{y} \sim \pi_{\theta}(\cdot | x)$, then translates back $\hat{x} \sim \pi_{\theta}(\cdot | \tilde{y})$. The agent is rewarded according to the quality of the round-trip reconstruction, for example via a similarity measure (`chrF++`) between x and \hat{x} . In this way, the model learns from self-consistency and monolingual structure, without any parallel supervision.

Let x denote a source sentence and y denote its target translation. The policy $\pi_{\theta}(y|x)$ produces candidate translations by autoregressively sampling tokens $y = (y_1, y_2, \dots, y_T)$ conditioned on the source.

3.2 Round-Trip Translation Reinforcement Training

Our key insight is that translation consistency in both directions provides a stronger learning signal than unidirectional training alone. For each source sentence x , we execute a two-phase generation process:

Phase 1: Forward Translation. Generate K candidate translations from source to target:

$$\hat{y}^{(k)} \sim \pi_{\theta}(\cdot | x), \quad k \in \{1, \dots, K\} \quad (3)$$

Phase 2: Back-Translation. For each forward translation $\hat{y}^{(k)}$, generate K back-translations to the source language:

$$\hat{x}^{(k,j)} \sim \pi_{\theta}(\cdot | \hat{y}^{(k)}), \quad j \in \{1, \dots, K\} \quad (4)$$

The model is trained on the back-translation phase, where rewards measure how well the back-translated sentence \hat{x} matches the original source x . This self-consistency objective encourages the model to produce translations that preserve semantic content bidirectionally.

3.3 Reward Design

We employ a reward function combining character-level and word-level translation quality metrics. For a generated back-translation \hat{x} given the original source x :

$$R(\hat{x}, x) = \lambda_{\text{chrF}} \cdot \text{chrF}++(\hat{x}, x) + \lambda_{\text{BLEU}} \cdot \text{BLEU}(\hat{x}, x) \quad (5)$$

chrF++. (Popović, 2015) is a character n-gram F-score that is particularly suitable for morphologically rich and low-resource languages, as it provides smoother gradients than word-based metrics:

$$\text{chrF}++ = (1 + \beta^2) \cdot \frac{\text{chrP} \cdot \text{chrR}}{\beta^2 \cdot \text{chrP} + \text{chrR}} \quad (6)$$

BLEU. (Nießen et al., 2000) with effective order provides complementary word-level precision signals, helping the model learn appropriate lexical choices. Algorithm 1 summarizes our training procedure.

4 Experiments

We base our experiments on NLLB-200 (No Language Left Behind) (Goyal et al., 2022), a massively multilingual translation model that covers 200 languages. In particular, we use the distilled 600M and 1.3B parameter variants. NLLB is an encoder-decoder Transformer (Vaswani et al., 2017) with language-specific tokens that specify the target language during generation. We use the NLLB-MD dataset (Guzmán et al., 2019; Goyal et al., 2022) in all our experiments. We use the provided split of parallel train-, dev- and test-set with respective sizes of 6000, 1310 and 1600 sentences. More on hyperparameters and implementation details are in appendix A. Our study focuses on four languages: Central Aymara, Friulian, Wolof and Russian. We use GoldFish (Chang et al., 2024) monolingual model family to score the fluency of the forward translations for every language in addition to chrF++.

First, we compare the base and trained models on these four languages using chrF++ and additionally translation fluency using Goldfish log-probabilities

Algorithm 1 Round-Trip Translation Reinforcement Training with GRPO

Require: Pretrained model π_θ , source-language corpus \mathcal{D} , learning rate η , group size K , KL coefficient β , clip parameter ϵ_c

- 1: Initialize reference policy $\pi_{\text{ref}} \leftarrow \pi_\theta$
- 2: **for** epoch = 1 to E **do**
- 3: **for** batch (x) in $\mathcal{D}_{\text{source lang}}$ **do**
- 4: // Forward translation
- 5: $\{\hat{y}^{(k)}\}_{k=1}^K \sim \pi_\theta(\cdot | x)$ with temperature sampling
- 6: // Back-translation
- 7: $\{\hat{x}^{(k)}\}_{k=1}^K \sim \pi_\theta(\cdot | \hat{y}^{(k)})$
- 8: // Compute rewards
- 9: $R^{(k)} \leftarrow \lambda_{\text{chrF}} \cdot \text{chrF}++(\hat{x}^{(k)}, x) + \lambda_{\text{BLEU}} \cdot \text{BLEU}(\hat{x}^{(k)}, x)$
- 10: // Update policy
- 11: $\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}_{\text{GRPO}}(\theta)$
- 12: **if** step mod $N_{\text{update}} = 0$ **then**
- 13: $\pi_{\text{ref}} \leftarrow \pi_\theta$
- 14: **end if**
- 15: **end for**
- 16: **end for**
- 17: **return** Fine-tuned policy π_θ

(Chang et al., 2024). Next, we present reward trajectories for forward and backward translation directions. Also, we present a reward function ablation study. Finally, we provide a qualitative analysis of model outputs on examples that are not seen in the training data.

5 Results

5.1 Forward Translations Evaluation

Translation Quality with chrF++ Following the original NLLB work (Goyal et al., 2022) in terms of low-resource language translation evaluation, Table 1 summarizes chrF++ scores on the forward translations before and after training. Across all languages and both model sizes, round-trip reinforcement training improves translation quality and fluency, yielding consistent absolute gains. We generally avoid model-based metrics like COMET (Rei et al., 2020) or BERTScore (Zhang et al., 2020) for forward translations due to the lack of reliable scoring models for many of the languages we evaluate. However, for Russian—where a reliable COMET model is available—we additionally report COMET scores, which change from 0.84 before training to 0.86 after training.

Language	600M		1.3B	
	Before	After	Before	After
Central Aymara	24.38	27.96 (+3.58)	26.01	28.70 (+2.69)
Friulian	45.75	47.50 (+1.75)	49.41	50.74 (+1.33)
Wolof	20.73	23.78 (+3.05)	24.73	26.81 (+2.08)
Russian	48.80	52.68 (+3.88)	50.10	54.21 (+4.11)

Table 1: Forward Translation chrF++ scores on the test set before and after training on each language for the 600M and 1.3B models using the composite reward function (Bleu + chrF++). Green values indicate absolute improvement in parentheses.

Language	Before	After	Δ
Central Aymara	-20.03	-19.87	0.16
Friulian	-19.23	-19.09	0.14
Russian	-15.64	-15.37	0.27
Wolof	-18.26	-18.07	0.19

Table 2: Goldfish model log-probability scores (natural log) before (vanilla) and after training, with improvement $\Delta = \text{trained} - \text{vanilla}$.

The largest gain is observed on Russian (+3.88), suggesting that the intrinsic round-trip signal strongly improves adequacy and consistency even when the target language is comparatively well supported in multilingual pretraining. It’s also evident that the more knowledge the model has on the language the better it is able to exploit it with the chrF++ reward signal.

For the **1.3B** model, we observe gains of **+1.33** on Friulian, **+2.08** on Wolof, **+4.11** on Russian, and **+2.69** on Central Aymara. Notably, the largest improvement again occurs on Russian (+4.11), indicating that increased model capacity does not saturate the benefits of our self-consistency objective; instead, larger models appear to better exploit the reward signal.

Translation Fluency Since we only reward the round-trip translation, it might be possible that the model generates nonsensical forward translation text, despite having an improved translation quality according to chrF++. To measure this, we check the fluency of the generated translation under the *Goldfish* language model, reporting average log-probability (natural log) before and after training (Table 2). Across all four languages, training increases the log-probability (i.e., makes it less negative), consistent with more fluent and model-preferred output distributions, confirming that our method do not damage the fluency.

The most significant improvement is again observed in Russian, mirroring the chrF++ trends.



Figure 2: forward translations chrF++ gain (after – before) for different reward functions across model/language settings.

This suggests that the round-trip reinforcement objective enhances not only adequacy, as captured by chrF++, but also surface-level well-formedness and overall likelihood according to an external scorer. These results are further supported by the qualitative analysis of the translation output discussed in a later section.

5.2 Validation Curves For Forward Translations

Figure 3 shows the chrF++ scores curves on the out-of-distribution evaluation set of NLLB-MD for the four languages and both the 1.3B and 600M NLLB models. In all four languages, both models improve steadily over the course of training. The curves are smooth and approximately monotonic, without large oscillations, which indicates that the optimization procedure is stable. Across all languages and throughout training, the 1.3B model consistently outperforms the 600M model. The 1.3B curves start from higher chrF++ scores and maintain a similar margin over the 600M curves at every step. The gap remains roughly stable over time, which implies that scaling model size gives better improvements. Also, The overall shape of the curves is similar across languages, but the absolute gains differ. Central Aymara and Wolof (top-left and bottom-right) start from relatively low chrF++ scores, and both models show substantial gains over the 8K steps. This indicates that continued training is particularly helpful for low-resource languages where the baseline is under-trained. Friulian and Russian (top-right and bottom-left) have higher initial chrF++ scores and still benefit from continued training, but the gains are smaller, especially for Friulian.

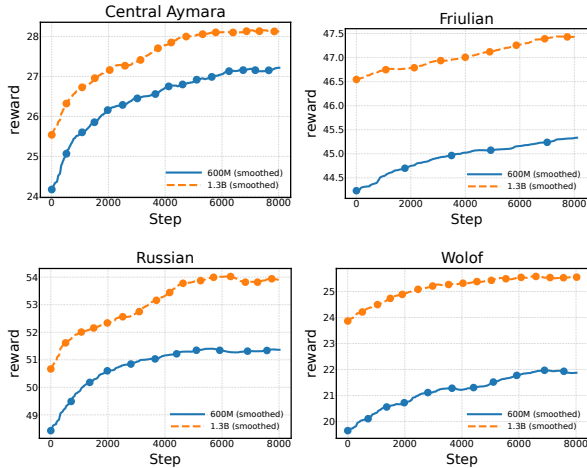


Figure 3: Validation curves showing forward translation chrF++ scores for the four languages after 8K optimization steps.

5.3 Reward Function Ablation Study

We next study how the choice of intrinsic reward impacts training results. We compare three reward functions: (i) a uniform weighted combination of chrF++ and BLEU (“BLEU+chrF++”), (ii) BLEU only, and (iii) chrF++ only. In all cases, we keep the training procedure and hyperparameters fixed and report the *forward translation* chrF++ score before and after adaptation on each language and model size.

Figure 2 summarizes the chrF++ *gain* (after minus before). Across all six settings, BLEU-only rewards consistently improve performance but yield substantially smaller gains than rewards that include chrF++. Averaged over all conditions, BLEU-only increases chrF++ by +1.60, whereas BLEU+chrF++ yields +2.41 and chrF++-only yields +2.47.

This trend is especially pronounced for the most challenging pairs (e.g., wol and aym), where BLEU-only rewards provide a weaker learning signal under sparse n-gram matching. By contrast, chrF++ provides a smoother character-level similarity objective that appears better aligned with low-resource and morphologically rich settings, leading to the strongest and most consistent improvements (best in 4/6 conditions). Adding BLEU to chrF++ is competitive and occasionally beneficial (notably for 1.3b-wol), but does not uniformly outperform chrF++ alone, suggesting that word-level precision signals can sometimes be redundant or slightly misaligned with the round-trip consistency objective.

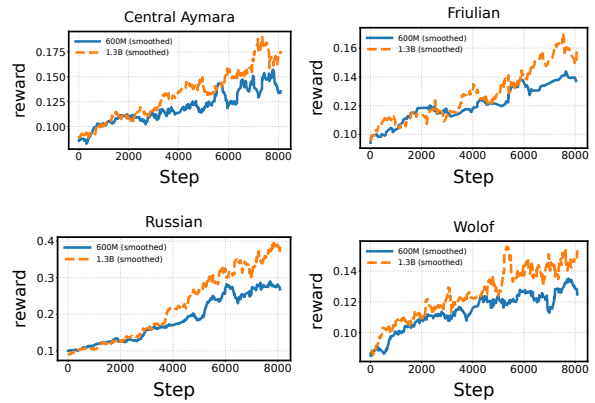


Figure 4: Training curves showing the English backward translation scores during training for the 1.3B and 600M parameter models on the 4 Languages.

Target	Before			After			Δ F1
	P	R	F1	P	R	F1	
Central Aymara	0.48	0.49	0.48	0.54	0.57	0.55	+0.07
Friulian	0.78	0.73	0.75	0.78	0.76	0.77	+0.02
Russian	0.77	0.69	0.73	0.80	0.77	0.79	+0.06
Wolof	0.53	0.48	0.51	0.54	0.55	0.55	+0.04

Table 3: Mean Precision/Recall/F1 for vanilla vs. trained models English Back-Translations.

5.4 Target-to-English Back-Translation Evaluation

Figure 4 shows smoothed backward-translation rewards (target \rightarrow English) over training steps for the 600M and 1.3B models. As in the forward direction, both models exhibit a clear upward trend across all four languages, indicating that continued training consistently improves backward translation quality rather than leading to early overfitting. The curves are noisier than in the forward setting, which is expected given the higher variance of the backward reward signal, but the overall trajectory is monotonic: rewards increase steadily over the first several thousand steps and then continue to grow more slowly toward the end of training. Across all languages, the 1.3B model (orange) either matches or exceeds the performance of the 600M model (blue) throughout training. The larger model typically starts from slightly higher rewards and gradually widens the gap as training progresses, especially for Russian and Wolof. This pattern suggests that scaling model capacity yields consistent benefits for backward translation, and that the adaptation procedure is able to exploit the additional capacity rather than saturating at the smaller model’s performance level.

Output Source	Language	chrF++	Text
Source (English)	English	N/A	That’s so awesome! Which one is it?
Reference	Central Aymara	N/A	¡Ukax wali musphkayawa! ¿Kawkirisá?
Base Model	Central Aymara	35.32	¡Wali sumäskiwa ! ¿Kawkiris uka?
Trained Model	Central Aymara	64.09	¡Ukax wali musparkañawa ! ¿Kawkiris ukaxa?
Source (English)	English	N/A	However, it was also of great importance "for a country from which we are separated by only two nations," she said.
Reference	Friulian	N/A	Dut câs, e jere ancje di grande impuartance "par une nazion de cuâl o sin separâts dome di dôs nazions," e à dite.
Base Model	Friulian	16.61	"O vin di fâ alc par fâ cressi la nestre identitât", e à dite la ministra.
Trained Model	Friulian	70.46	Dut câs, e jere ancje di grande impuartance "par un paîs di dulà che o sin separâts dome di dôs nazions", e à dit jê.
Source (English)	English	N/A	Rheumatic illnesses in particular rheumatic arthritis and ankylosing spondylitis (so called Bechterew’s disease)
Reference	Wolof	N/A	Jàngoroy Rheumatik rawatina arthritis rheumatik ak ankylosing spondytilis (nu duppe ko jàngoro Bechterew)
Base Model	Wolof	5.84	Jàmm yuy néew-ji-doole yuy Jàmm yuy néew-ji-doole yuy
Trained Model	Wolof	42.22	Yàmm yu réwmatik, rawatina aarthritis réwmatik ak spondylitis ankylosing (yite bi ñuy wax jàngoroy Bechterew)
Source (English)	English	N/A	Sounds great to me, I’m so excited! What is the horse’s name?
Reference	Russian	N/A	Zvuchit zdorovo, ya tak rad! Kak zovut konya?
Base Model	Russian	21.32	- Kak zovut loshadey?
Trained Model	Russian	74.52	Zvuchit zdorovo, ya tak rada! Kak zovut loshadey?

Table 4: Example outputs for Source (English), Reference translations, Base Model, and Trained model across four languages. Tokens in green match the reference, tokens in red differ from the reference, and tokens in blue mark segments where the trained model is closer to the reference than the base model. chrF++ scores are shown for model outputs with respect to the Reference sentence.

BERTScore To validate our findings beyond chrF++ and bleu score in which are used as a reward function, we also show the improvement with BERTScore (Zhang et al., 2020). Table 3 shows the BERTScores on the back-translated English sentences before and after training. Across all four target languages, the trained model improves BERTScore F1 relative to the vanilla baseline, with gains ranging from +0.02 (Friulian) to +0.07 (Central Aymara). The improvements are generally driven by higher recall after training (notably Central Aymara and Russian), while precision also increases or remains stable, indicating the trained back-translations are both more complete and at least as faithful under this semantic metric. Our findings show consistent improvement upon different metrics.

5.5 Qualitative Analysis On Translation Outputs

Finally, in Table 4 we provide a qualitative comparison of translations produced by the 1.3B base model and the trained model across four evaluation

languages (one example each; additional samples are reported in Appendix B). Overall, the trained model produces outputs that are close to the human reference.

For **Central Aymara**, the base model captures the positive tone but diverges in lexical choice and phrasing. The trained model matches the reference opener (“*Ukax wali . . .*”) and better recovers the intended predicate meaning, yielding a question form that is closer to the reference. This suggests that training primarily improves lexical selection and common conversational templates in this low-resource setting.

For **Friulian**, the base model fails to translate the source sentence and instead generates unrelated quoted text (English: “We have to do something to grow our identity,” the minister said.). In contrast, the trained model restores the intended discourse structure and recovers most of the reference content (“*Dut câs, e jere ancje di grande impuartance . . .*”), with only a minor mismatch in the final pronoun choice and the final translation to English is faithful to the reference. The chrF++ jump (16.61 → 70.46)

459 indicates substantially improved faithfulness and
460 reduced semantic drift.

461 For **Wolof**, the baseline output degenerates
462 into repetition, a common failure mode for low-
463 resource and domain-specific inputs. The trained
464 model largely recovers the technical terminology
465 and coordination structure (e.g., *aarthritis réw-*
466 *matik*, *spondylitis ankylosing*, and the parenthet-
467 ical *Bechterew*), though some surrounding func-
468 tion words remain imperfect. The corresponding
469 chrF++ gain (5.84 → 42.22) suggests improved ro-
470 bustness (avoiding repetition) and better retention
471 of specialized terminology.

472 Finally, for **Russian**, the base model produces
473 a partially correct question but uses an incorrect
474 lexical/inflectional form for “horse.” The trained
475 model matches the reference more closely in the
476 preceding clause and improves morphosyntactic
477 realization (e.g., correct gender marking in “*rada*”),
478 while the noun error persists. The chrF++ increase
479 (21.32 → 74.52) highlights improved fluency and
480 morphology, with remaining challenges in precise
481 lexical choice.

482 Taken together, these examples suggest that self-
483 supervised training improves both (i) content cov-
484 erage and faithfulness (especially evident in Friu-
485 lian and Wolof) and (ii) surface-form realization,
486 including idiomatic openers and inflectional mor-
487 phology (Central Aymara and Russian). Resid-
488 ual mismatches point to remaining difficulties in
489 fine-grained lexical selection and domain-specific
490 function-word control. Nevertheless, despite rely-
491 ing only on monolingual reward signals, the model
492 is able to self-improve on the target low-resource
493 languages without any direct supervision.

494 6 Conclusion

495 We presented a self-supervised RL fine-tuning ap-
496 proach for low-resource machine translation based
497 on round-trip bootstrapping with NLLB. By trans-
498 lating English into a target language and rewarding
499 the quality of the reconstructed English using in-
500 trinsic signals (chrF++ and BLEU), GRPO training
501 yields consistent chrF++ improvements on an out-
502 of-distribution test set for Central Aymara, Friulian,
503 Wolof, and Russian across both 600M and 1.3B
504 models. Qualitatively, the trained models reduce
505 common low-resource failure modes such as repe-
506 tition and semantic drift, and produce outputs that
507 are more fluent and faithful to the reference.

508 Our ablations show that chrF++ provides the

509 most reliable learning signal in this setting, while
510 adding BLEU is competitive but not uniformly bet-
511 ter. Overall, these results indicate that, even with
512 only monolingual reward signals, pretrained multi-
513 lingual models can self-improve translation quality
514 for low-resource languages without direct parallel
515 supervision.

516 7 Limitations

517 Our study has several limitations that point to clear
518 directions for future work.

519 **Scaling to larger NLLB checkpoints.** Due to
520 resource constraints, we only evaluate the distilled
521 600M and 1.3B NLLB variants, and do not test
522 whether the same round-trip GRPO training re-
523 mains stable and beneficial for substantially larger
524 models (e.g., 3.3B NLLB). This is because empiri-
525 cal scaling laws suggest that multilingual MT qual-
526 ity often improves predictably with model scale,
527 and that scaling can change the effective capac-
528 ity allocated to individual language pairs, including
529 out-of-distribution behavior (Fernandes et al., 2023;
530 Kaplan et al., 2020).

531 **Broader low-resource language coverage within**
532 **NLLB-MD.** We focus on four languages (Cen-
533 tral Aymara, Friulian, Wolof, and Russian) and
534 do not evaluate on the remaining low-resource
535 languages included in NLLB-MD (notably Bho-
536 jpuri and Dyula) (Goyal et al., 2022). Extending
537 experiments to these additional languages would
538 strengthen conclusions about robustness across
539 scripts, typology, and domain mismatch, and would
540 help identify when round-trip rewards are most ef-
541 fective versus when they induce overly paraphrastic
542 solutions.

543 **Richer reward functions for low-resource MT.**
544 Our intrinsic reward is a simple linear combina-
545 tion of chrF++ and BLEU computed on the re-
546 constructed English, which is effective but limited.
547 Future work should explore richer and more ro-
548 bust reward design, for example: (i) adding learned
549 semantic metrics (e.g., COMET/BLEURT-style sig-
550 nals) to reduce reliance on lexical overlap, (ii) us-
551 ing target-side fluency rewards (monolingual LM
552 scoring) to better control grammaticality.

553 8 Ethical Considerations

554 Our approach fine-tunes multilingual MT mod-
555 els using self-supervised, round-trip reinforce-
556 ment learning (English→target→English) with au-

557	automatic rewards (chrF++/BLEU). While this avoids	<i>Machine Learning</i> , volume 202 of <i>Proceedings of</i>	606
558	collecting new human-labeled parallel data, it in-	<i>Machine Learning Research</i> , pages 10053–10071.	607
559	troduces several potential risks:	PMLR.	608
560	Amplification of model errors and biases. Be-	R. A. Fisher. 1992. <i>Statistical Methods for Research</i>	609
561	cause training relies on the model’s own genera-	<i>Workers</i> , pages 66–70. Springer New York, New	610
562	tions, systematic errors (e.g., gender/role stereo-	York, NY.	611
563	types, offensive terms, dialectal mismatches) can	Naman Goyal and 1 others. 2022. No language left	612
564	be reinforced. Biases present in pretraining or in	behind: Scaling human-centered machine translation.	613
565	monolingual English data may propagate into low-	<i>arXiv preprint arXiv:2207.04672</i> .	614
566	resource outputs.	Francisco Guzmán, Peng-Jen Chen, Myle Ott, Juan	615
567	Hallucinations and harmful content translation.	Pino, Guillaume Lample, Philipp Koehn, Vishrav	616
568	RL may increase fluency while still hallucinating	Chaudhary, and Marc’Aurelio Ranzato. 2019. Two	617
569	facts, altering named entities, or mistranslating	new evaluation datasets for low-resource machine	618
570	safety-critical content (medical/legal). Addition-	translation: Nepali-english and sinhala-english.	619
571	ally, improved MT quality could facilitate trans-	Barry Haddow and 1 others. 2022. Survey of machine	620
572	lation of toxic or hateful content into additional	translation for low-resource languages. <i>Computing</i>	621
573	languages.	<i>Research Repository</i> .	622
574	Mitigations. We mitigate these risks by (i) moni-	Jianfei He, Shichao Sun, Xiaohua Jia, and Wenjie Li.	623
575	toring for degeneracy (repetition/copying) and re-	2024. <i>Recovery should never deviate from ground</i>	624
576	porting qualitative examples, (ii) evaluating with	<i>truth: Mitigating exposure bias in neural machine</i>	625
577	multiple metrics (chrF++, BLEU, BERTScore on	<i>translation</i> . In <i>Proceedings of the 25th Annual</i>	626
578	back-translations, and target-side fluency), (iii) lim-	<i>Conference of the European Association for Ma-</i>	627
579	iting training data to vetted corpora and document-	<i>chine Translation (Volume 1)</i> , pages 68–79, Sheffield,	628
580	ing compute settings for transparency.	UK. European Association for Machine Translation	629
581	References	(EAMT).	630
582	Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam	Amr Hendy, Mohamed Abdelrehim, Amr Sharaf,	631
583	Shazeer. 2015. <i>Scheduled sampling for sequence</i>	Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita,	632
584	<i>prediction with recurrent neural networks</i> . <i>Preprint</i> ,	Young Jin Kim, Mohamed Afify, and Hany Hassan	633
585	arXiv:1506.03099.	Awadalla. 2023. How good are gpt models at ma-	634
586	Tyler A. Chang, Catherine Arnett, Zhuowen Tu, and	chine translation? a comprehensive evaluation. <i>arXiv</i>	635
587	Benjamin K. Bergen. 2024. <i>Goldfish: Monolingual</i>	<i>preprint arXiv:2302.09210</i> .	636
588	<i>language models for 350 languages</i> . <i>Preprint</i> .	Xiuming Huang. 1990. <i>A machine translation system</i>	637
589	Alexis Conneau and 1 others. 2020. Unsupervised cross-	<i>for the target language inexpert</i> . In <i>COLING 1990</i>	638
590	lingual representation learning at scale. <i>ACL</i> .	<i>Volume 3: Papers presented to the 13th International</i>	639
591	Christian Federmann, Oussama Elachqar, and Chris	<i>Conference on Computational Linguistics</i> .	640
592	Quirk. 2019. <i>Multilingual whispers: Generating</i>	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B.	641
593	<i>paraphrases with translation</i> . In <i>Proceedings of the</i>	Brown, Benjamin Chess, Rewon Child, Scott Gray,	642
594	<i>5th Workshop on Noisy User-generated Text (W-NUT</i>	Alec Radford, Jeffrey Wu, and Dario Amodei. 2020.	643
595	<i>2019)</i> , pages 17–26, Hong Kong, China. Association	<i>Scaling laws for neural language models</i> . <i>Preprint</i> ,	644
596	for Computational Linguistics.	arXiv:2001.08361.	645
597	Zhaopeng Feng, Shaosheng Cao, Jiahao Ren, Jiayuan	Diederik P. Kingma and Jimmy Ba. 2017. <i>Adam:</i>	646
598	Su, Ruizhe Chen, Yan Zhang, Zhe Xu, Yao Hu, Jian	<i>A method for stochastic optimization</i> . <i>Preprint</i> ,	647
599	Wu, and Zuozhu Liu. 2025. <i>Mt-r1-zero: Advanc-</i>	arXiv:1412.6980.	648
600	<i>ing llm-based machine translation via r1-zero-like</i>	Philipp Koehn and Rebecca Knowles. 2017. Six chal-	649
601	<i>reinforcement learning</i> . <i>Preprint</i> , arXiv:2504.10160.	<i>lenges for neural machine translation</i> . In <i>Proceedings</i>	650
602	Patrick Fernandes, Behrooz Ghorbani, Xavier Garcia,	<i>of the First Workshop on Neural Machine Transla-</i>	651
603	Markus Freitag, and Orhan Firat. 2023. <i>Scaling laws</i>	<i>tion</i> .	652
604	<i>for multilingual neural machine translation</i> . In <i>Pro-</i>	Guillaume Lample, Alexis Conneau, Ludovic Denoyer,	653
605	<i>ceedings of the 40th International Conference on</i>	and Marc’Aurelio Ranzato. 2018. <i>Unsupervised ma-</i>	654
		<i>chine translation using monolingual corpora only</i> .	655
		<i>Preprint</i> , arXiv:1711.00043.	656
		Andrii Maksai and Pascal Fua. 2018. <i>Eliminating expo-</i>	657
		<i>sure bias and loss-evaluation mismatch in multiple</i>	658
		<i>object tracking</i> . <i>Preprint</i> , arXiv:1811.10984.	659

660 Paul McNamee and Kevin Duh. 2023. An extensive
661 exploration of back-translation in 60 languages. In
662 *Findings of the Association for Computational Lin-*
663 *guistics: ACL 2023*, pages 8166–8183.

664 Sonja Nießen, Franz Josef Och, Gregor Leusch, Her-
665 mann Ney, and 1 others. 2000. An evaluation tool for
666 machine translation: Fast evaluation for mt research.
667 In *LREC*.

668 Adam Paszke, Sam Gross, Francisco Massa, Adam
669 Lerer, James Bradbury, Gregory Chanan, Trevor
670 Killeen, Zeming Lin, Natalia Gimelshein, Luca
671 Antiga, Alban Desmaison, Andreas Köpf, Edward
672 Yang, Zach DeVito, Martin Raison, Alykhan Te-
673 jani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang,
674 and 2 others. 2019. *Pytorch: An imperative style,*
675 *high-performance deep learning library*. *Preprint*,
676 arXiv:1912.01703.

677 Maja Popović. 2015. chrF: character n-gram f-score for
678 automatic mt evaluation. In *Proceedings of the tenth*
679 *workshop on statistical machine translation*, pages
680 392–395.

681 Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli,
682 and Wojciech Zaremba. 2016. *Sequence level*
683 *training with recurrent neural networks*. *Preprint*,
684 arXiv:1511.06732.

685 Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon
686 Lavie. 2020. *Comet: A neural framework for mt*
687 *evaluation*. *Preprint*, arXiv:2009.09025.

688 John Schulman, Filip Wolski, Prafulla Dhariwal,
689 Alec Radford, and Oleg Klimov. 2017. *Prox-*
690 *imal policy optimization algorithms*. *Preprint*,
691 arXiv:1707.06347.

692 Rico Sennrich, Barry Haddow, and Alexandra Birch.
693 2016. *Improving neural machine translation models*
694 *with monolingual data*. *Preprint*, arXiv:1511.06709.

695 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,
696 Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan
697 Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024.
698 *Deepseekmath: Pushing the limits of mathematical*
699 *reasoning in open language models*. *Preprint*,
700 arXiv:2402.03300.

701 Richard S. Sutton and Andrew G. Barto. 2018. *Rein-*
702 *forcement Learning: An Introduction*, second edition.
703 The MIT Press.

704 Ashish Vaswani and 1 others. 2017. Attention is all you
705 need. In *Advances in Neural Information Processing*
706 *Systems*.

707 Weiyue Wang and Rico Sennrich. 2020. On exposure
708 bias, hallucination and domain shift in neural ma-
709 chine translation. In *Proceedings of the 58th Annual*
710 *Meeting of the Association for Computational Lin-*
711 *guistics*, pages 3544–3552.

Hyperparameter	Value
Learning rate	2×10^{-6}
Batch size	2
Group size K	4
Generation temperature	1.8
Top- k sampling	100
Top- p sampling	0.95
KL coefficient β	0.04
Clip parameter ϵ_c	0.2
Reference update frequency	16 steps
Maximum sequence length	64 tokens
Training epochs	2

Table 5: Hyperparameters used for training.

712 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien
713 Chaumond, Clement Delangue, Anthony Moi, Pier-
714 ric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz,
715 Joe Davison, Sam Shleifer, Patrick von Platen, Clara
716 Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le
717 Scao, Sylvain Gugger, and 3 others. 2020. *Trans-*
718 *formers: State-of-the-art natural language processing*.
719 In *Proceedings of the 2020 Conference on Empirical*
720 *Methods in Natural Language Processing: System*
721 *Demonstrations*, pages 38–45, Online. Association
722 for Computational Linguistics.

723 Lijun Wu, Fei Tian, Tao Qin, Jianhuang Lai, and
724 Tie-Yan Liu. 2018. *A study of reinforcement*
725 *learning for neural machine translation*. *Preprint*,
726 arXiv:1808.08866.

727 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q.
728 Weinberger, and Yoav Artzi. 2020. *Bertscore:*
729 *Evaluating text generation with bert*. *Preprint*,
730 arXiv:1904.09675.

A Implementation Details 731

We implement our method using PyTorch (Paszke et al., 2019) and the Hugging Face Transformers (Wolf et al., 2020) library. Table 5 summarizes the key hyperparameters used in our experiments. 732 733 734 735

We use AdamW (Kingma and Ba, 2017) for opti- 736 mization. During generation, we employ nucleus 737 sampling with temperature scaling to encourage 738 diverse candidate generation. 739

B Qualitative Translation Examples 740

B.1 Central Aymara 741

Reference: Ukax jichhax amuyt’añjamaxiwa, 742 kipka irnaqawinakjamakiwa. Jaqinakax kunti 743 munapki ch’amañañawa. ¡kayak anatañ lu- 744 rañ yant’asma! Wali jasakiw apnaqañaxa, qutan 745 smart’añjamaspawa jan ukasti qhuru umanakan 746

747	chuyma jalsuyasirjamarakispa.		
748	Source (EN): Ah that makes a ton of sense, they		
749	seem like similar jobs at least in some respects.		
750	Encouraging people and helping them meet their		
751	goals. You should totally try kayaking! It's super		
752	versatile, it can be relaxing on a lake or thrilling		
753	doing whitewater.		
754	Vanilla (system output): Ah, al à sens, a somein		
755	lavoris similârs almancul in cualchi sens. incoragjâ		
756	lis personis e judâlis a rivâ ai lôr obietîfs. Tu va-		
757	ressis di provâ il kayak! al è une vore versatile, al		
758	pues jessi un lavôr di		
759	Vanilla (EN gloss): Ah, al à sens, a somin lavoris		
760	similarârs at least in some sens. incoragjâ lis per-		
761	sonis e judâlis a rivâ ai lôr obietîfs. Tu varessis di		
762	provâ il kayak! al è une vore versatile, al pues jessi		
763	un lavôr di		
764	Trained (system output): Ah, al à sens, a somein		
765	lavoris similis almancul in cualchi sens. Incoragjâ		
766	lis personis e judâlis a rivâ ai lôr obietîfs. Tu vares-		
767	sis di provâ a fâ il kayak! Al è une vore versatile,		
768	al pues jessi ri		
769	Trained (EN gloss): Ah, in a sense, some people		
770	work similarly at least in some senses. Incoragjâ		
771	the people and Jews to achieve their objectives. You		
772	have a variety of ways to try to do the kayaking!		
773	Al è une vore versatile, you can be very versatile,		
774	and you can be the best in the world.		
775	Reference: O ai apene sfofnât i biscos PB cu la		
776	forme dai mei cjans!		
777	Source (EN): I baked PB cookies in the shape of		
778	my dogs!		
779	Vanilla (system output): O ai fat biscots di PB a		
780	forme dai miei cans!		
781	Vanilla (EN gloss): Or PB biscuits shaped like my		
782	dogs!		
783	Trained (system output): O ài fat biscuits PB a		
784	forme dai miei cans!		
785	Trained (EN gloss): I made PB cookies in the		
786	shape of my dogs!		
787	Reference: O soi sigur che al divente mat, ma		
788	tu asreson, il mangja co cusini no l'è just par lôr.		
789	Perciò o ten simpri sot man al di bon par lôr quant		
790	co cusini.		
791	Source (EN): I am sure the smells drive him crazy,		
792	but you are right. The food I bake is not good for		
793	them. I have dog treats on hand anytime I bake for		
794	this reason.		
795	Vanilla (system output): O soi sigure che i odôrs		
796	lu fasin diventâ mat, ma tu âs reson. Il mangjâ che		
	o fâs no je buine par lôr. Par chest motif o ai simpri		797
	lis cjasis a man.		798
	Vanilla (EN gloss): I'm sure the smells drive them		799
	crazy, but you're right. The food I make isn't good		800
	for them. That's why I always keep the houses		801
	handy.		802
	Trained (system output): O soi sigure che i odôrs		803
	lu fasin diventâ mat, ma tu âs reson. Il mangjâ che		804
	o fâs jo nol è bon par lôr. Par cheste reson o ai		805
	simpri lis cjampanis pai cans in man.		806
	Trained (EN gloss): I am sure the smells make		807
	them go crazy, but you are right. The food I make		808
	is not good for them. That is why I always have the		809
	dogs' bells in my hand.		810
	B.2 Friulian		811
	Reference: O ai apene sfofnât i biscos PB cu la		812
	forme dai mei cjans!		813
	Source (EN): I baked PB cookies in the shape of		814
	my dogs!		815
	Vanilla (system output): O ai fat biscots di PB a		816
	forme dai miei cans!		817
	Vanilla (EN gloss): Or PB biscuits shaped like my		818
	dogs!		819
	Trained (system output): O ài fat biscuits PB a		820
	forme dai miei cans!		821
	Trained (EN gloss): I made PB cookies in the		822
	shape of my dogs!		823
	Reference: O soi sigur che al divente mat, ma		824
	tu asreson, il mangja co cusini no l'è just par lôr.		825
	Perciò o ten simpri sot man al di bon par lôr quant		826
	co cusini.		827
	Source (EN): I am sure the smells drive him crazy,		828
	but you are right. The food I bake is not good for		829
	them. I have dog treats on hand anytime I bake for		830
	this reason.		831
	Vanilla (system output): O soi sigure che i odôrs		832
	lu fasin diventâ mat, ma tu âs reson. Il mangjâ che		833
	o fâs no je buine par lôr. Par chest motif o ai simpri		834
	lis cjasis a man.		835
	Vanilla (EN gloss): I'm sure the smells drive them		836
	crazy, but you're right. The food I make isn't good		837
	for them. That's why I always keep the houses		838
	handy.		839
	Trained (system output): O soi sigure che i odôrs		840
	lu fasin diventâ mat, ma tu âs reson. Il mangjâ che		841
	o fâs jo nol è bon par lôr. Par cheste reson o ai		842
	simpri lis cjampanis pai cans in man.		843
	Trained (EN gloss): I am sure the smells make		844
	them go crazy, but you are right. The food I make		845
	is not good for them. That is why I always have the		846

847	dogs' bells in my hand.	
848	Reference: No tant. Stât impegnât a ciri plas gnûs.	
849	Soi sigû co tacarai a vore a mae quat co sarai stac	
850	di cusinà.	
851	Source (EN): Not recently. I have been consumed	
852	by trying to find recipes. I tend to hobby jump quite	
853	a bit. I am sure I will start knitting like crazy once	
854	I get tired of baking.	
855	Vanilla (system output): No dome di pôc, mi soi	
856	consumade a cirî ricetis, o soi propit un pôc a saltâ	
857	fûr di hobby, o soi sigure che o fasarai il tricot come	
858	une matete cuant che mi stufarà di fâ il pan.	
859	Vanilla (EN gloss): Not only that, I've been con-	
860	sumed with finding recipes, or I've been kind of	
861	jumping out of my hobby, or I'm sure I'll be knit-	
862	ting like a mat when I get tired of baking.	
863	Trained (system output): No di pôc timp. O soi	
864	stât consumât a cirî ricetis. O ten a saltâ un grum	
865	di hobbies. O soi sigûr che o tacarà a fâ il tricot	
866	come un mat cuant che o varès stufât di fâ il pan.	
867	Trained (EN gloss): Not recently. I have been	
868	consumed looking for recipes. I want to skip a lot	
869	of hobbies. I am sure that I will start knitting like	
870	crazy when I get tired of baking bread.	
871	B.3 Russian	
872	Reference: Da, no ya podumyvayu provesti vikend	
873	v etom kan'yone ryadom so mnoy. Eto ne plyazh,	
874	no i ne les, dolzhno byt' zdorovo.	
875	Source (EN): I do, but I am thinking of doing a	
876	weekend in this canyon near me. It's not the beach,	
877	but it's nto a forest either, should be cool.	
878	Vanilla (system output): Ya dumayu provesti	
879	vykhodnye v kan'one nedaleko ot menya.	
880	Vanilla (EN gloss): I'm thinking of spending the	
881	weekend in a canyon near me.	
882	Trained (system output): Da, no ya dumayu	
883	provesti vykhodnye v etom kan'one ryadom so	
884	mnoy. Eto ne plyazh, no eto tozhe les, dolzhno byt'	
885	kruto.	
886	Trained (EN gloss): Yes, but I'm thinking of	
887	spending the weekend in that canyon next to me.	
888	It's not the beach, but it's also the forest, so that	
889	should be cool.	
890	Reference: Tochno. Kak tvoya sem'ya? Skol'ko	
891	vas?	
892	Source (EN): That's right. How is your family?	
893	how many of you are there?	
894	Vanilla (system output): Kak vasha sem'ya?	
895	Vanilla (EN gloss): How's your family?	
	Trained (system output): Tochno. Kak vasha	896
	sem'ya? Skol'ko vas tam?	897
	Trained (EN gloss): Exactly. How is your family?	898
	How many are you there?	899
	Reference: Nikakikh problem! Skol'ko i kakikh?	900
	Source (EN): No problem at all! How many and	901
	which ones?	902
	Vanilla (system output): – Skol'ko i kakie?	903
	Vanilla (EN gloss): - How many and what?	904
	Trained (system output): Nikakikh problem!	905
	Skol'ko i kakie?	906
	Trained (EN gloss): No problem! How many and	907
	what?	908
	B.4 Wolof	909
	Reference: Kurélu réewum Suède bu wérgu-	910
	yaram (Rättsmedicinalverket) tàmbali woon na di	911
	def tests yi ci lu teel ci at mi.	912
	Source (EN): Sweden's national Forensic	913
	Medicine Agency (Rättsmedicinalverket) started	914
	carrying out the tests earlier year.	915
	Vanilla (system output): Bèsum-bèsum-bèsum-	916
	bèsum-bèsum-bèsum-bèsum-bèsum-bèsum-	917
	bèsum-bèsum-bèsum-bèsum-bèsum-bèsum-bés	918
	Vanilla (EN gloss): Day-day-day-day-day-day-	919
	day-day-day-day-day-day-day-day-day-day-	920
	day-day-day-day-day-day-day-day-day-day-	921
	Trained (system output): Bàngaasu Farās gu réew	922
	mi ci Suweed (Rättsmedicinalverket) tàmbalee	923
	jëfandikoo ay seetlu ci njëlbéenug at mi.	924
	Trained (EN gloss): The French national bank of	925
	Sweden (Rättsmedicinalverket) began using the	926
	studies at the beginning of the year.	927
	Reference: Moom ak mbiipkat bu dëkk Bulgari	928
	Miroslav Petkov, dafay def ab konseer yu am bàkk	929
	biiGéeg gu ñuul gi tasena ak Àll bu ñuul.; wala sax,	930
	Ëalkan folkloređaje naak ab njàng mu Allëmaañ	931
	bu kalasik ak taxawaay bu dëggër ci misiku brass.	932
	Source (EN): Together with the Bulgarian trum-	933
	peter Miroslav Petkov, he is giving a concert that	934
	carries the slogan “The Black Sea meets the Black	935
	Forest”, or even, “Balkan folklore” meets a classic	936
	German education with roots in brass music.	937
	Vanilla (system output): Ci biir li muy ànd ak	938
	mbëjug buum buum buum buum buum buum buum	939
	buum buum buum buum buum buum buum buum	940
	buum buum buum buum buum buum buum buum	941
	buum buum buum	942
	Vanilla (EN gloss): In addition to its use as	943
	an anti-inflammatory, it is also used as an anti-	944

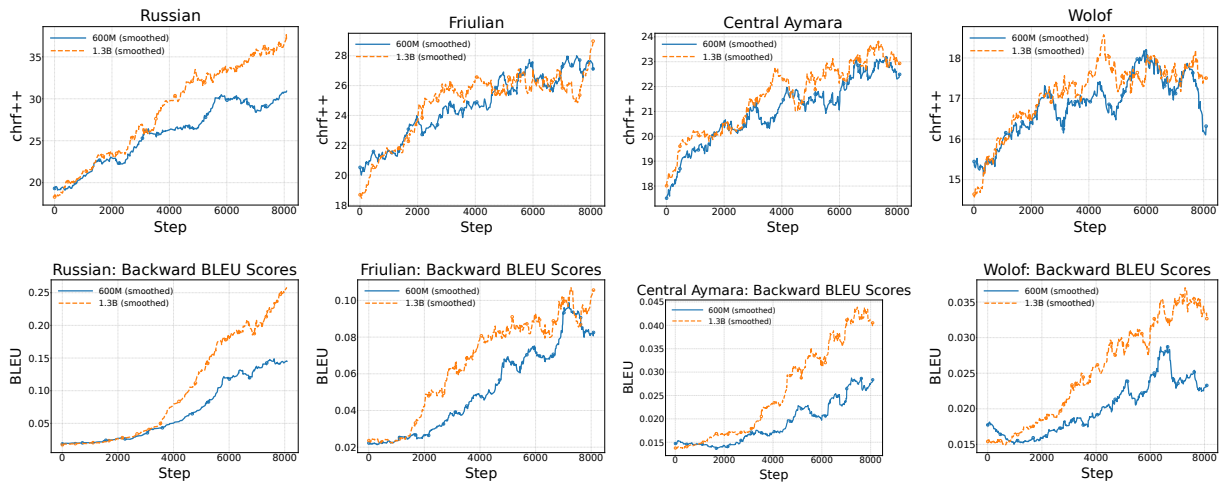


Figure 5: Validation curves after 8K optimization steps. Top row: forward translation chrF++ scores. Bottom row: backward translation BLEU scores.

inflammatory and anti-inflammatory agent.

Trained (system output): Ci biir àndandoo ak mbëjjug buulgar bi di Miroslav Petkov, mu ngi doon jox ab kaaraange bu am kàdduy waxtaan bi di “Géej gu Màgg gi daje ak Géej gu Màgg gi”, walla sax, “Balkan folkl

Trained (EN gloss): In collaboration with the Bulgarian writer Miroslav Petkov, she was developing a cover with the words of the speech “The Old Sea meets the Great Sea”, or even, “Balkan folklore”

Reference: Ñeen lañu! Amna sama benn makk bu goor ak samay waajur. Ñun goor yi deñoo bëgg bannexuloo ci benn sigaar ci ngoon gi.

Source (EN): There are four of us! I have an older brother and my two parents. Us men like to enjoy a cigar together in the evenings.

Vanilla (system output): Am nanu ñeenti nit, am naa magam ak samay maam ak sama baay.

Vanilla (EN gloss): We had four children, one of whom died in infancy.

Trained (system output): Am nanu ñeenti nit! Am naa mag bu góor ak samay ñaari waajur. Nun góor yi, bëgg nanu a lekk cigare ci ngoon ci ngoon.

Trained (EN gloss): We have four people! I have a older brother and my two parents. We men like to smoke in the afternoon in the afternoon.

C Additional Training Plots

Figure 5 reports validation curves for the forward and backward translation models over 8K optimization steps. The top row shows chrF++ for forward translation, and the bottom row shows BLEU for backward translation, for Russian (Cyril), Friulian (Latin), Aymara (Latin), and Wolof (Latin). These

curves are included for completeness and to illustrate the round-trip reinforcement learning training dynamics across languages and directions. For both the forward and backward translations during training, it’s evident that there is no plateau and the models can further benefit from continued RL training.

978
979
980
981
982
983
984