

Reinforcement Learning for Machine Translation: Fine-Tuning the MT5 LLM for English-Indonesian Translation

Arkan Abuyazid

Department of Computer Science

Texas A&M University

College Station, Texas

ata757@tamu.edu

Abstract—This proposal outlines an exploratory approach employing reinforcement learning (RL) to fine-tune the MT5-base model for English-Indonesian translation. Our objective is to investigate the feasibility of leveraging RL-based techniques, with translation quality metrics (such as BLEU, ROUGE-L, and METEOR) serving as reward signals. By framing translation as a sequential decision-making process within a Markov Decision Process (MDP), we aim to establish a flexible optimization framework and compare RL algorithms, notably REINFORCE with baseline subtraction and Proximal Policy Optimization (PPO), against more conventional fine-tuning practices. The focus of this work is to explore the potential benefits and underlying mechanisms of reward-driven fine-tuning in low-resource settings.

Index Terms—Neural Machine Translation, Reinforcement Learning, MT5, Low-Resource Languages, Preference Optimization

Source Code: <https://github.com/aabuyazid/MT5-EN-ID-Translator>

Video: <https://youtu.be/0DjFmuV0yrs>

I. PROPOSAL

The central aim of this proposal is to explore how reinforcement learning can enhance machine translation, specifically addressing challenges faced when working with low-resource language pairs like English-Indonesian. We provide an initial framework and motivation for applying RL to an existing multilingual text-to-text transformer model (MT5-base) with the goal of improving translation quality using a combination of automatic evaluation metrics. We select MT5-base (580M parameters) over larger alternatives like mT5-XXL (13B) due to its demonstrated balance between multilingual capability and practical deployability. Our approach is designed to further understand the dynamics and potential improvement pathways offered by reward-driven methods.

A. Sequential Nature

Machine translation is by nature a sequential generation task where each output token is influenced by previously generated tokens. Recognizing this inherent sequentiality, our work frames the translation process as a partially observable Markov Decision Process (POMDP). Each decision made by the decoder is critical, as it influences the semantic and

syntactic construction of the final sentence. Our exploration focuses on how the sequential dependencies can be more effectively captured by integrating long-term reward signals into token generation decisions, and thus addresses issues of coherence and fluency across lengthy translations.

B. State Space

The state at each time step, denoted as s_t , encapsulates several components:

- **Encoder’s Source Representation:** The semantic and syntactic information extracted from the English sentence.
- **Decoder’s Hidden State:** The internal state reflecting the decoded sequence up to token t , which carries forward context.
- **Generated Target Tokens:** The sequence of previously generated Indonesian tokens serves both as context and a historical path for current decisions.
- **Attention Weights History:** A record of attention distributions to potentially modulate rewards based on alignment quality.

Expanding on these components, we intend to analyze how each aspect contributes to the overall translation quality and how the RL agent might best utilize this state information for improved decision making.

C. Action Space

In our formulation, the action space A is defined as the full MT5 vocabulary, comprising 250,112 tokens. Each action represents the selection of a next token, with decisions made under a temperature-controlled softmax sampling strategy to balance exploration and exploitation. The extensive vocabulary challenges the RL algorithm to manage a high-dimensional space effectively. Our investigation will shed light on how this combinatorial complexity affects learning dynamics and convergence.

D. Transition Function

Given the transformer architecture, the transition function is modeled as deterministic:

$$s_{t+1} = \text{MT5}_\theta(s_t, a_t) \quad (1)$$

Here, the function encapsulates the state update based on the current state and the token chosen. In our proposed exploration, we will consider how this deterministic mapping interacts with stochastic policy choices, as well as potential modifications that could account for uncertainty and long-term dependencies in the generated sequences.

E. Reward Function

The reward function R_t is a weighted combination of several quality metrics:

$$R = \lambda_1 \text{BLEU} + \lambda_2 \text{ROUGE-L} + \lambda_3 \text{METEOR} - \lambda_4 \text{RepPenalty} \quad (2)$$

Each coefficient λ_i will be tuned via a proposed grid search method to balance the contributions of precise translation match (BLEU), sequence-level overlap (ROUGE-L), and semantic alignment (METEOR), while penalizing undesirable repetitive patterns. Our aim is to explore various weighting schemes to understand their impact on the translation policy, laying the groundwork for RL strategies that can optimally balance multiple metrics.

F. Initial Framework and Baseline

At this stage, our plan is to establish the RL framework by implementing the REINFORCE algorithm enhanced with a baseline subtraction. This will serve as a preliminary method, benchmarked against the standard supervised fine-tuning approach. Our evaluation strategy—though conceptual at this point—will include:

- **Automatic Metrics:** BLEU-4, ROUGE-L, and METEOR scores computed in a post-hoc analysis to quantify translation quality.
- **Human Evaluation:** A plan to incorporate human judgments on translation adequacy and fluency using the Flores-200 evaluation benchmark.
- **Training Dynamics:** Analysis of training stability and convergence characteristics within the RL framework.

These measures are proposed to validate the conceptual benefits of the RL approach for translation tasks, while noting that our current focus is on methodology and potential rather than empirical validation.

G. Stretch Goal

Beyond the initial RL implementation, we propose a stretch goal that involves exploring advanced preference optimization techniques. This includes:

- **Proximal Preference Optimization (PPO):** An investigation into PPO as a scalable method for aligning translation policies with human preferences.

- **Contrastive Preference Optimization (CPO):** Experimenting with pairwise ranking of translations to incorporate nuanced human feedback mechanisms.

The stretch goal is designed to expand our understanding of how complex reward structures can guide large language models (LLMs) towards not only higher quality translations but also more human-aligned outputs. This part of the proposal will explore theoretical considerations and the potential methodological integration of these techniques into the RL framework.

II. LITERATURE REVIEW

Our research builds on a robust set of prior works that have investigated reinforcement learning in natural language processing, preference optimization, and neural machine translation. The following subsections expand on the existing literature and illustrate why our proposed approach is both relevant and promising.

A. Reinforcement Learning Foundations

The REINFORCE algorithm [1] laid the groundwork for policy gradient methods in sequence generation tasks. Despite its pioneering role, the algorithm faces challenges with delayed rewards and large action spaces—issues we aim to address by adapting baseline subtraction and entropy regularization techniques. Our proposal is motivated by the potential to extend these classical methods to handle the intricacies of modern transformer architectures.

B. Preference Optimization and Policy Alignment

Recent advancements in preference optimization, including Direct Preference Optimization (DPO) [4] and Proximal Preference Optimization (PPO) [5], have shown significant promise in aligning language model outputs with human values. The introduction of Contrastive Preference Optimization (CPO) [6] further demonstrates the value of incorporating pairwise comparisons into the learning process. These methodologies highlight the potential of explicitly modeling human preferences alongside traditional reward metrics, and our work proposes to integrate such approaches into the context of machine translation.

C. Evaluation Metrics in Machine Translation

Evaluation metrics such as BLEU [2], ROUGE-L [9], and METEOR [10] have long served as proxies for human judgment in machine translation. The literature reveals both strengths and limitations in these metrics—while BLEU effectively captures n-gram overlap, ROUGE-L provides insights into longer sequence coherence, and METEOR offers semantic alignment considerations. By combining these metrics into a unified reward function, our approach seeks to leverage the complementary strengths of each metric to provide a more holistic assessment of translation quality.

D. Model Architecture and Low-Resource Challenges

MT5 [7] represents a crucial step toward building multilingual models that can handle a broad spectrum of languages, including those with relatively low resources. While the transformer architecture has achieved remarkable success in high-resource settings, its direct application to low-resource language pairs often exposes limitations in translation fidelity and coverage. Our research is driven by the hypothesis that an RL-based fine-tuning strategy can better adapt large models like MT5 to the unique challenges present in low-resource translation scenarios, bridging gaps that conventional supervised fine-tuning might not address.

E. Dataset Considerations

The Flores-200 benchmark [3] offers a standardized evaluation framework specifically designed for low-resource languages. By aligning our evaluation strategy with the methodologies established in Flores-200, we intend to create a consistent and comparable measure of translation quality. Our current proposal focuses on the conceptual design of these evaluation metrics as integral components of the reward function within our RL framework.

III. IMPLEMENTATION

We first obtained the Google MT5-small model and performed supervised fine-tuning on the Asian Language Treebank (ALT) dataset [11]. After the supervised stage, we trained the model using REINFORCE, comparing the effectiveness of the model at each stage of training (i.e. base vs. supervised fine-tuned vs. REINFORCE). We also used Google Colab to conduct all of our experiments. We have made extensive efforts to implement PPO; however, we were not able to do so. As such, all analysis regarding reinforce learning will be directed towards the REINFORCE algorithm.

A. Supervised Fine-Tuning

The ALT dataset was introduced by Riza et al. [11] and provides parallel English–Indonesian sentence pairs. We fine-tuned MT5-small on this dataset under a standard cross-entropy objective to create an initial translation model.

B. REINFORCE

The Flores-200 dataset [3] was used for the REINFORCE portion of training as it contains high-quality, yet sparse, English-Indonesian sentence pairs. The baseline is determined by calculating the moving average of the rewards. We used a learning rate of 5×10^{-6} and trained for 15 epochs.

C. Reward Function

The reward function R_t is defined as a weighted combination of translation quality metrics and a repetition penalty:

$$R = 0.4 \cdot \text{BLEU} + 0.3 \cdot \text{ROUGE-L} + 0.2 \cdot \text{METEOR} - 0.1 \cdot \text{RepPenalty}$$

The coefficients were set as follows: $\lambda_{\text{BLEU}} = 0.4$, $\lambda_{\text{ROUGE-L}} = 0.3$, $\lambda_{\text{METEOR}} = 0.2$, and $\lambda_{\text{RepPenalty}} = 0.1$. These values were decided empirically.

IV. EVALUATION

A. Metrics

We have utilized the BLEU, ROUGE-L, and METEOR metrics to evaluate the model at each stage of training. It turns out that REINFORCE does improve model performance, albeit rather insignificantly.

Stage	BLEU		ROUGE-L		METEOR	
	Mean	Std	Mean	Std	Mean	Std
Supervised	17.65	15.68	42.25	19.12	38.82	20.39
REINFORCE	17.69	15.66	42.33	18.87	38.96	20.15

TABLE I
COMPARISON OF SUPERVISED AND REINFORCE STAGES

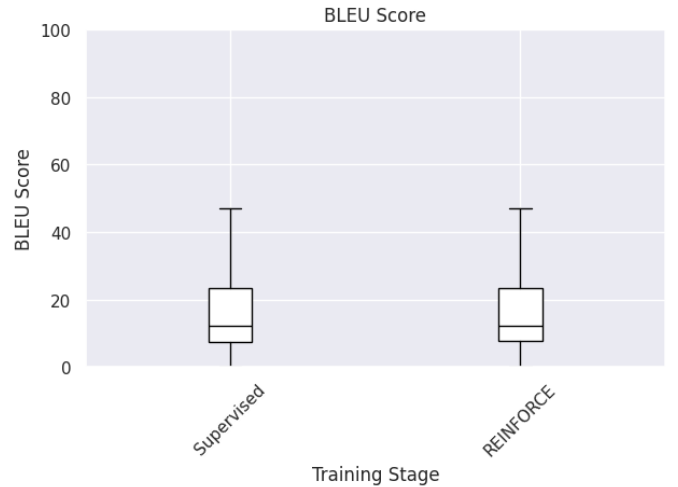


Fig. 1. BLEU scores on Flores200 test dataset at each stage of training.

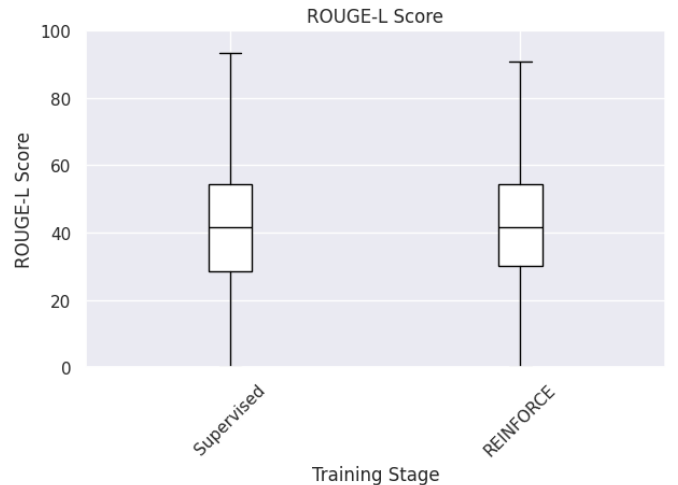


Fig. 2. ROUGE-L scores on Flores200 test dataset at each stage of training.

B. Training Performance

It is evident that the loss is rather noisy as well as the rewards. There is an argument to be made that the rewards

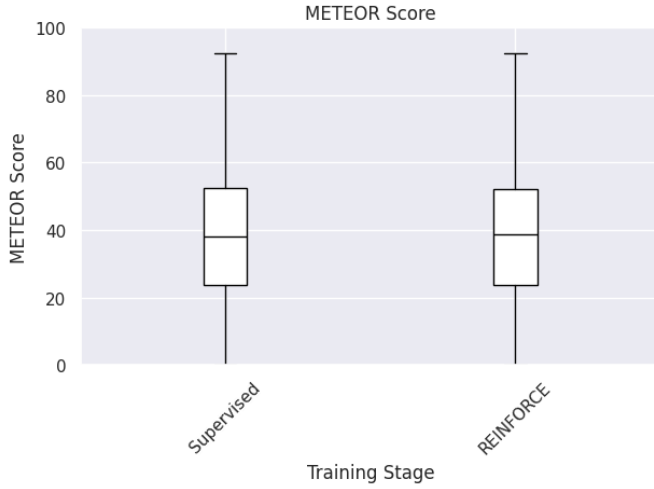


Fig. 3. METEOR scores on Flores200 test dataset at each stage of training.

have a slight upward trajectory, implying that perhaps training the model with a larger number of epochs could yield better results. However, considering that the reward function is within the intervals $[0, 100]$ and the changes within the rewards do not exceed 1, it is unlikely increasing the number of epochs will improve model performance.

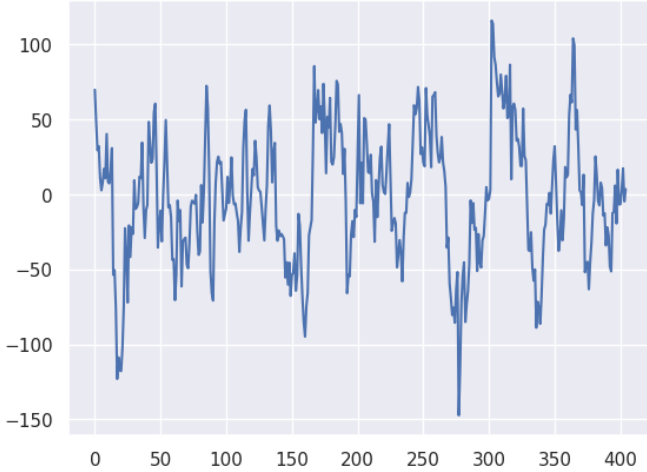


Fig. 4. The REINFORCE training loss with a smoothing window of 25.

V. TRAINING CHALLENGES AND UNSTABLE BEHAVIOR

Implementing the REINFORCE algorithm for LLM training was not an easy task, and implementing PPO proved to be insurmountable for us. We will discuss the challenges we have faced in implementing the two.

A. REINFORCE Training Instability

The REINFORCE algorithm revealed several critical instability issues that required iterative debugging. Initial attempts to use a learning rate of 5×10^{-4} catastrophically destabilized the model, causing it to generate empty strings for all inputs.

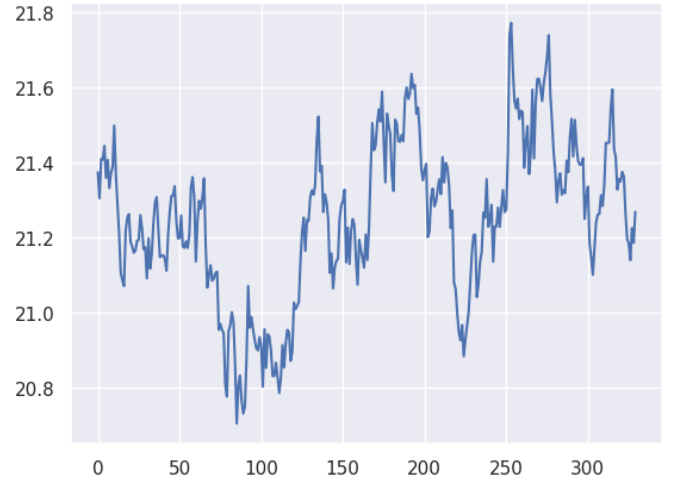


Fig. 5. The REINFORCE training rewards with a smoothing window of 25.

This suggested that large parameter updates were erasing the knowledge gained during supervised fine-tuning. Reducing the learning rate to 5×10^{-6} alleviated the immediate collapse but introduced new concerns. Training metrics show high noise in both loss and rewards, implying that the model may not be converging to the optimal policy.

B. PPO Implementation Pains

The PPO exhibited fundamental incompatibilities with the task constraints. Early on, the KL divergence between successive policies became negative, a mathematical impossibility under standard PPO theory. This is most likely caused by conflicting gradients from the multi-component reward function interacting poorly with the policy clipping mechanism. Such divergence occurred consistently across multiple random seeds, suggesting either a flawed reward scaling approach or inherent limitations in applying PPO's trust region constraints to sequence generation tasks with composite rewards. Ultimately, the training loop prematurely exits due to errors; therefore, no meaningful data can be drawn from the experiment.

VI. CONCLUSION

This exploratory work investigated reinforcement learning as a mechanism for improving English-Indonesian machine translation in low-resource settings, framing the task as a Markov Decision Process optimized through reward signals derived from automated metrics (BLEU, ROUGE-L, METEOR). While the REINFORCE algorithm demonstrated tentative feasibility for fine-tuning MT5-base, its practical utility remains questionable due to fundamental limitations. The high variance of policy gradient updates, coupled with the delicate balance required to prevent catastrophic forgetting of supervised fine-tuning knowledge, resulted in unstable training trajectories and minimal convergence guarantees. Notably, attempts to implement Proximal Policy Optimization (PPO) failed due to irrecoverable negative KL divergence, suggesting that trust-region constraints may be inherently incompatible

with composite reward structures in neural machine translation.

These findings underscore the challenges of applying RL to low-resource language pairs like English-Indonesian, where sparse reward signals and limited parallel data amplify optimization instabilities. While REINFORCE offers theoretical appeal as a gradient estimator for sequence generation, its sample inefficiency and sensitivity to hyperparameters make it unlikely to scale effectively for large language model fine-tuning. The persistent noise floor in reward signals and inconsistent correlation between automated metrics and human judgments further question the viability of pure metric-driven RL for this domain. Future work should prioritize hybrid approaches that combine RL’s flexibility with more stable optimization techniques, particularly for languages lacking the robust evaluation frameworks available in high-resource settings.

REFERENCES

- [1] R. J. Williams, “Simple statistical gradient-following algorithms for connectionist reinforcement learning,” *Machine Learning*, 1992.
- [2] K. Papineni et al., “BLEU: a method for automatic evaluation of machine translation,” *ACL* 2002.
- [3] N. Goyal et al., “The Flores-101 Evaluation Benchmark for Low-Resource and Multilingual Machine Translation,” arXiv:2106.03193, 2022.
- [4] R. Rafailov et al., “Direct Preference Optimization: Your Language Model is Secretly a Reward Model,” arXiv:2305.18290, 2023.
- [5] A. Rafailov et al., “Proximal Preference Optimization: A New Approach to Scalable Policy Alignment,” Technical Report, 2023.
- [6] Z. Sun et al., “Contrastive Preference Optimization: Pushing the Boundaries of LLM Utility in Alignment,” arXiv:2401.08417, 2024.
- [7] L. Xue et al., “mT5: A Massively Multilingual Pre-trained Text-to-Text Transformer,” *NAACL* 2021.
- [8] A. Radford et al., “Improving Language Understanding by Generative Pre-Training,” OpenAI Technical Report, 2018.
- [9] C.-Y. Lin, “ROUGE: A Package for Automatic Evaluation of Summaries,” *ACL* 2004.
- [10] S. Banerjee et al., “METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments,” *ACL* 2005.
- [11] H. Riza et al., “Introduction of the asian language treebank,” *2016 Conference of The Oriental Chapter of International Committee for Coordination and Standardization of Speech Databases and Assessment Techniques (O-COCOSDA)*, 2016.