

# LAMPAT: Low-Rank Adaption for Multilingual Paraphrasing Using Adversarial Training

Khoi M. Le<sup>1, 2\*</sup>, Trinh Pham<sup>2\*</sup>, Tho Quan<sup>2</sup>, Anh Tuan Luu<sup>3†</sup>

<sup>1</sup>VinAI Research, Vietnam

<sup>2</sup>Ho Chi Minh City University of Technology (HCMUT), VNU-HCM, Ho Chi Minh City, Vietnam

<sup>3</sup>Nanyang Technological University, Singapore

v.khoilm1@vinai.io, phkhanhtrinh23@gmail.com, qttho@hcmut.edu.vn, anhtuan.luu@ntu.edu.sg

## Abstract

Paraphrases are texts that convey the same meaning while using different words or sentence structures. It can be used as an automatic data augmentation tool for many Natural Language Processing tasks, especially when dealing with low-resource languages, where data shortage is a significant problem. To generate a paraphrase in multilingual settings, previous studies have leveraged the knowledge from the machine translation field, i.e., forming a paraphrase through zero-shot machine translation in the same language. Despite good performance on human evaluation, those methods still require parallel translation datasets, thus making them inapplicable to languages that do not have parallel corpora. To mitigate that problem, we proposed the first unsupervised multilingual paraphrasing model, LAMPAT (Low-rank Adaptation for Multilingual Paraphrasing using Adversarial Training), by which monolingual dataset is sufficient enough to generate a human-like and diverse sentence. Throughout the experiments, we found out that our method not only works well for English but can generalize on unseen languages as well. Data and code are available at <https://github.com/VinAIRResearch/LAMPAT>.

## Introduction

Paraphrase generation involves the transformation of a given sentence or phrase into its equivalent form while preserving its semantic content. By leveraging the power of NLP techniques, paraphrase generation can enhance several applications, such as machine translation (Freitag et al. 2020), information retrieval (Lewis et al. 2020), question-answering systems (Gan and Ng 2019), and text summarization (Cao et al. 2017). However, most existing approaches in paraphrase generation focus on a single language such as English, limiting their effectiveness in multilingual scenarios where accurate and contextually appropriate paraphrases are crucial.

Most of the current approaches for multilingual paraphrasing are built around the mechanism of machine translation. Thompson and Post (2020b) make use of the multilingual neural machine translation (MNMT) model from

\*These authors contributed equally.

†Corresponding author.

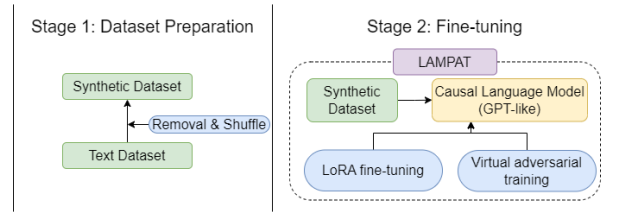


Figure 1: The training process of LAMPAT consists of multiple stages. Firstly, we create a synthetic parallel corpus using unsupervised monolingual data. Next, we utilize LoRA to effectively fine-tune our model. Finally, we obtain the self-supervised model LAMPAT through the utilization of Virtual Adversarial Training.

Thompson and Post (2020a) to translate between the same source and target languages with a special decoding algorithm to reduce the lexical overlaps, creating a paraphrase of the input text. Meanwhile, round-trip translation (Federmann, Elachqar, and Quirk 2019), which creates a pivot language for source language to translate forth and back to create a different wording output, is another approach for paraphrase generation.

While the machine translation approach is widely recognized for its effectiveness in multilingual paraphrasing, it encounters various limitations. One primary hurdle is the necessity of obtaining reliable parallel corpora of high quality for machine translation, which can be challenging to obtain in real-world scenarios. Another drawback lies in the lack of diversity in the generated output, often resulting in output sequences that closely resemble the input and fail to preserve crucial information present in the original input, as depicted in Table 1. Many MNMT models utilize heuristic blocking techniques during inference to avoid generating output sequences that are identical to the input. However, this approach limits the model’s ability to alter word order or employ diverse syntactic structures, such as inversion or active-to-passive transformations. Consequently, the generation lacks diversity. Furthermore, these blocking algorithms heavily rely on the distribution of the vocabulary. For instance, Thompson and Post (2020b) reduce the probability of selecting subsequent subwords in  $n$ -grams to encourage the model to choose different subwords. Neverthe-

less, these alternative subwords may include antonyms or unrelated words, potentially leading to paraphrases with semantic meanings deviating from the intended direction and producing inappropriate results.

Type	Sentence
Input	I like to eat pasta.
Human reference	<b>My favourite food is</b> pasta.
Lack of diversity	<i>I like eating</i> pasta.
Incorrect meaning	I like to eat <u>paste</u> .

Table 1: While human can change the structure of the sentence to create paraphrase (in bold), paraphrasing model usually tends to replace words or slightly modify the syntax (in italic). In the worst case, paraphrasing model even changes the meaning of the sentence by using inappropriate word replacement (in underline).

In this work, we propose LAMPAT (**L**ow-rank **A**daptation for **M**ultilingual **P**araphrasing using **A**dversarial **T**raining) as an approach to mitigate the strict requirements of parallel corpora by learning in an unsupervised manner and alleviate the problem of duplicate generation by using adversarial training objectives. According to Figure 1, to eliminate the need for parallel corpora, we use the monolingual dataset and apply a series of processes (i) identify stop words (ii) remove stop words (iii) randomly shuffle the words to create the corrupted version of the original input. The training of the model focuses on the objective of reconstructing the original sentence from a corrupted version, aiming to recreate the initial sentence. However, by this learning objective, the model is drawn to generate the same sentence compared to the original sentence, which does not create a syntactically diverse paraphrase. To cope with this problem, we further propose using Virtual Adversarial Training (VAT) (Zhang et al. 2019) and noise perturbation added directly to the input embedding to steer the model towards a more diverse paraphrase generation, as in Figure 2. In addition, Large Language Models (LLMs), are known to experience the catastrophic forgetting (Kaushik et al. 2021) during full fine-tuning; therefore, we adapt LoRA (Hu et al. 2021) as a parameter-efficient fine-tuning method to partially update the model’s prior knowledge and preserve all the linguistic knowledge on which the model has been pre-trained. In general, LAMPAT can effectively generate human-like paraphrases while preserving the original semantic meaning and employing different syntactic structures to promote the diversity of the predictions.

In summary, the key contributions of this paper are as followed:

- To resolve the requirement of parallel corpora for machine translation, we propose the unsupervised learning method for multilingual paraphrasing.
- To address the issue of predominantly generating identical or highly lexical-similar outputs, we incorporate noise perturbation and a virtual labeling strategy into the adversarial training process, aiming to alleviate this limitation.
- We expand the multilingual paraphrasing evaluation

dataset to include more languages and leverage future research in multilingual paraphrase generation.

## Methodology

The training process of our paraphrasing model is illustrated in Figure 1, comprising three essential elements: Synthetic Parallel Corpora, Parameter-Efficient Fine-Tuning (PEFT), and Virtual Adversarial Training (VAT). We employ the Self-supervised model to generate paraphrases, which not only addresses the data shortage issue using unsupervised learning but also maintains semantic similarity and enhances lexical diversity as well.

### Problem Definition

Given 2 sentences  $x$  and  $y$ , where  $x$  is the original sentence, and  $y$  is the paraphrase reference of  $x$ . Let  $\mathcal{M}(x)$  be the meaning of  $x$  and  $\mathcal{S}(x, y)$  be the lexical or syntactic similarity between  $x$  and  $y$ .  $\hat{y} = \underset{y}{\operatorname{argmax}}[p(y|M(x)) - S(x, y)]$ .

The term  $p(y|M(x))$  is the probability of generating  $y$  that conveys the same meaning as  $x$  (i.e.  $M(x)$ ).  $S(x, y) \in [0, 1]$  measures the lexical similarity of  $x$  and  $y$ , where  $S(z, z) = 1$  for every  $z$ . Based on this formulation, paraphrasing should be a method which not only allows us to convey the same meaning but also enhances lexical diversity in the generated text.

### Synthetic Parallel Corpora

Synthetic Parallel Corpora is a significant component of LAMPAT. First, we corrupt the input by removing the stop words in the sentence, and then randomly shuffling the words in the remaining text. The corrupted sentence is referred to as the source sequence  $S$ , while the original uncorrupted sentence is the target sequence  $T$ . We have a set of stop words  $A$ , which are removed from the sentences. Our goal is to generate the paraphrase by reconstructing  $T$  from the keywords or the corrupted sentence  $S$ , where  $S = \text{Shuffle}(T - A)$ . When we fine-tune the model, we create the input sequence  $X$  by combining the source and target sequences with a special symbol in between. The input sequence  $X$  is represented as  $(x_1, x_2, x_3, \dots, x_k, \backslash n, x_{k+1}, x_{k+2}, \dots, x_h)$ , where the source sequence is denoted as  $S = (x_0, x_1, \dots, x_k)$ , and the target sequence is denoted as  $T = (x_{k+1}, x_{k+2}, \dots, x_h)$ . The special character  $\backslash n$  is included to differentiate between the source and target tokens, and it also serves as a prompt during the inference process.

### Parameter-Efficient Fine-Tuning

The method chosen for parameter-efficient fine-tuning is Low-rank Adaptation (LoRA) (Hu et al. 2021). A major drawback of fine-tuning is that the resulting model contains the same number of parameters as the original model. LoRA overcomes this by indirectly training some dense layers in a neural network by optimizing rank decomposition matrices of the dense layers’ changes during adaptation while keeping the pre-trained weights frozen. This also allows for efficient task-switching by simply replacing the matrices, resulting in reduced storage requirements and task-switching overhead.

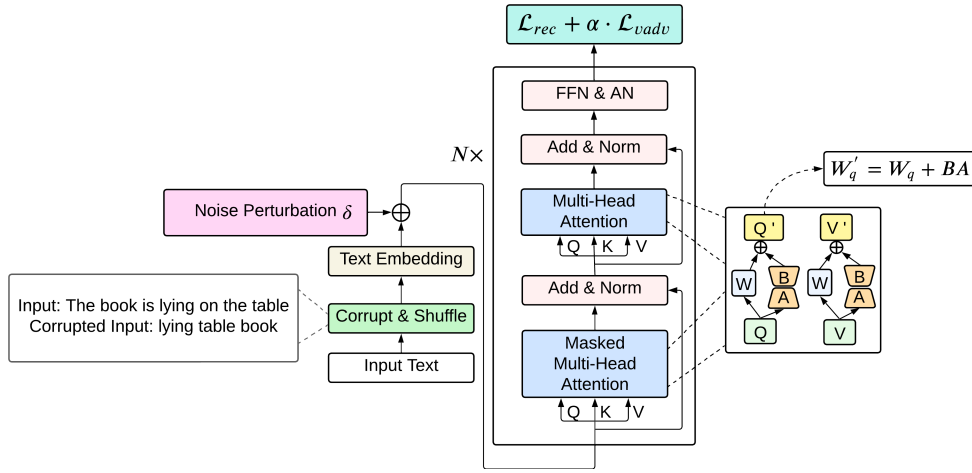


Figure 2: LAMPAT is showcased using actual inputs. Initially, an input text undergoes corruption by removing stopwords and shuffling. Then, a noise perturbation, denoted as  $\delta$ , is introduced into the text embedding to generate a paraphrase that exhibits lexical diversity. The transformer block is replicated  $N$  times, with the Multi-Head Attention component decomposed into low-rank matrices for efficient fine-tuning. Lastly, LAMPAT is trained using virtual adversarial training, incorporating a two-component loss function: the reconstruction loss  $\mathcal{L}_{rec}$  and the virtual adversarial regularizer  $\mathcal{L}_{vadv}$ .

LoRA (Hu et al. 2021) introduces a constraint on the update of a pre-trained matrix  $W_x \in \mathbb{R}^{d \times k}$  using a low-rank matrix  $W_\beta \in \mathbb{R}^{d \times k}$ . Instead of directly updating  $W_x$ , LoRA modifies it as  $W'_x = W_x + W_\beta$  and focuses on updating the parameters involved in the construction of  $W_\beta$ . The construction of  $W_\beta$  involves the multiplication of two matrices  $B \in \mathbb{R}^{d \times r}$  and  $A \in \mathbb{R}^{r \times k}$ , where  $r \ll \min(d, k)$ , resulting in a low-rank matrix  $W'_x = W_x + BA$ . By using LoRA, we reduce the number of parameters to tune from  $d \times k$  to  $r \times (d + k)$ . Specifically, we follow the application of LoRA to the query and value transformation matrices in the multi-head attention sublayers, as in Hu et al. (2021).

### Virtual Adversarial Training

Consider a standard classification task with an underlying data distribution  $D$  over examples  $x \in \mathbb{R}^d$  and corresponding labels  $y$ . Assume that we are given a suitable loss function  $\mathcal{L}_{rec}$ , for example, the cross-entropy loss for a neural network. As usual,  $\theta \in \mathbb{R}^p$  is the set of model parameters. Our goal is then to find model parameters  $\theta$  that satisfy:

$$\min_{\theta} E_{(x,y) \sim \mathcal{D}} [\mathcal{L}_{rec}(f(x, \theta), y)] \quad (1)$$

Once we have created the synthetic parallel data, we define a set of permissible perturbations  $\mathcal{C} \subseteq \mathbb{R}^d$  for each data point  $x$ , which represents the manipulative ability of the adversary. Instead of directly using samples from the distribution  $\mathcal{D}$  in the loss  $\mathcal{L}$ , we allow the adversary  $\delta \in \mathcal{C}$  to first perturb the input embedding. Specifically,  $x$  denotes the sub-word embedding in  $f(x, \theta)$ . We notice that by perturbing the embedding space  $x + \delta$ , rather than the input space, adversarial training may unintentionally favour on-manifold perturbations over regular perturbations, leading to improved generalization. Hence, we apply perturbations to the embedding space. On the other hand, complete label information

may not always be available, and especially in this unsupervised manner, we aim to output virtual labels other than the label  $y$ . Consequently, we adopt a strategy to replace the label  $y$  with its current approximation,  $f(x, \theta)$ . This approximation is not necessarily naive, as  $f(x, \theta)$  tends to be close to  $y$  when the number of labelled training samples is large. This rationale also explains the use of the term ‘‘virtual’’ in Miyato et al. (2018). Essentially, we employ virtual labels generated from  $f(x, \theta)$  in place of paraphrasing labels, and compute the adversarial direction based on these virtual labels. As a result, we replace the Equation 1 by:

$$\min_{\theta} E_{(x,y) \sim \mathcal{D}} \left[ \max_{\delta} [\mathcal{L}_{rec}(f(x + \delta, \theta), y) + \alpha \mathcal{L}_{vadv}(f(x + \delta, \theta), f(x, \theta))] \right] \quad (2)$$

In reference to Algorithm 1, our method can reduce the computational cost of adversarial training using projection over constraints algorithms. It achieves similar levels while conducting fewer sweeps of forward and backward propagation, making it faster and less computationally expensive. Our method takes advantage of every propagation to update weights and allows multiple updates per iteration, potentially leading to faster convergence. According to Zhang et al. (2019), by combining these factors, the  $K$  ascent steps significantly accelerate standard adversarial training. After that, the model’s parameter  $\theta$  is updated all at once with the accumulated gradients. By taking a descent step along the  $K$  gradients, we can approximately calculate the following objective function:

$$\min_{\theta} E_{x,y \sim \mathcal{D}} \left[ \frac{1}{K} \sum_{t=0}^{K-1} \max_{\delta} [\mathcal{L}_{rec}(f(x + \delta, \theta), y) + \alpha \mathcal{L}_{vadv}(f(x + \delta, \theta), f(x, \theta))] \right] \quad (3)$$

---

**Algorithm 1: Low-rank Adaptation Multilingual Paraphrasing using Adversarial Training.**

---

**Input:**  $X$ : Training samples,  $f(x; \theta)$ : the machine learning model parameterized by  $\theta$ ,  $\epsilon$ : the perturbation bound,  $\tau$ : the global learning rate,  $\alpha$ : the smoothing proportion of the adversarial training,  $\eta$ : the ascent step size,  $\mathcal{H}$ : the Hessian gradient matrix,  $N$ : the number of epochs,  $K$ : the number of ascent steps,  $e^*$ : the number of epochs trained with PGD algorithm,  $\gamma$ : the scaling factor.

**Output:**  $\theta$ .

```
1: for epoch = 1... $N$  do
2:   for minibatch  $B \in X$  do
3:      $\delta \sim \gamma \cdot \mathcal{N}(0, \sigma^2 I)$ 
4:     for  $m = 1...K$  do
5:       Accumulate gradient of parameter  $\theta$ :
6:        $g_m \leftarrow g_{m-1} + \frac{1}{K} E_{(x,y) \in B} [\nabla_{\theta} \mathcal{L}_{rec}(f(x + \delta, \theta), y) + \alpha \nabla_{\theta} \mathcal{L}_{adv}(f(x + \delta, \theta), f(x, \theta))]$ 
7:       Calculate the gradient of the perturbation  $\delta$ :
8:        $g_{adv} \leftarrow \nabla_{\delta} \mathcal{L}_{adv}(f(x + \delta, \theta), f(x, \theta))$ 
9:       Update the perturbation  $\delta$  through gradient ascent:
10:      if epoch  $\leq e^*$  then
11:         $\delta \leftarrow \prod_{\|\delta\| \leq \epsilon} (\delta + \eta \frac{g_{adv}}{\|g_{adv}\|_F})$ 
12:      else
13:         $\delta \leftarrow \prod_{\mathcal{H}} (\delta + \eta [\mathcal{H}]^{-1} \frac{g_{adv}}{\|g_{adv}\|_F})$ 
14:      end if
15:    end for
16:    Update the parameter  $\theta$  through gradient descent:
17:     $\theta \leftarrow \theta - \tau g_K$ 
18:  end for
19: end for
```

---

The rationale behind computing  $g(\delta)$  with respect to the virtual adversarial regularizer  $\nabla_{\delta} \mathcal{L}_{adv}$  instead of  $\nabla_{\delta} \mathcal{L}_{rec}$  in Algorithm 1 is due to the unsupervised nature of the model training and the objective to guide the model towards the virtual labels rather than reconstructing the original sentence, which could result in duplication. Equation 3 is essentially replacing the original batch  $x$  with a virtual batch that is  $K$  times larger, comprising samples with embeddings of  $X + \delta_0, \dots, X + \delta_{K-1}$ . While the original adversarial training Equation 2 minimizes the maximum risk at a single estimated point in the vicinity of each training sample, Equation 3 minimizes the maximum risk at each ascent step and guides the model towards the virtual labels with minimal additional overhead. Moreover, we use both Projected Gradient Descent (PGD) and Projected-Newton Method (PNM) in Algorithm 1. The Gradient Descent step involves descending along the linear estimate of the function, while Newton’s step involves moving the point towards the minimum of the parabola that approximates the function, which can lead to faster convergence.

## Experimental Setup

### Primary Model

**Pre-trained model:** We utilize the mGPT model, a pre-trained multilingual GPT-like model, with a 1.3B checkpoint consisting of 100K tokens and supporting 61 languages. The LoRA implementation from PEFT<sup>1</sup>.

**Dataset:** To assess the fine-tuning, we choose to use the latest version WMT19 (Foundation 2019) to train the model. This dataset covers a wide range of 15 languages including

Arabic, Czech, German, English, Spanish, French, Hindi, Indonesian, Italian, Japanese, Kazakh, Dutch, Portuguese, Russian, and Chinese. The WMT19 dataset we use is in its latest version, which is just released in 2023. To balance language resources, we employ a uniform distribution to sample sentences, creating a training set of nearly 600k sentences and a validation set of around 100k sentences. For training, we use the monolingual version of WMT19 and corrupt the input sentence by removing all of the stop words<sup>2</sup>, further we randomly shuffle the words 33% of the time. The goal is to reconstruct a sentence from its keywords, or its corrupted sentence. For the machine translation approach, we use the available bilingual version available in WMT19, we sample with the same strategy as LAMPAT’s training with 600k samples for training and 100k for validation.

### Baseline Model

We compare our method with other approaches such as multilingual machine translation (MMT) proposed in Thompson and Post (2020b) and denoising auto-encoder (DAE) in Guo et al. (2019). Initially, these methods are trained on different datasets, thus, we re-train them, following the procedures proposed in each paper, on the WMT19 dataset.

### Evaluation Dataset

We follow Guo et al. (2019), randomly select 10k sentences respectively from each language of English, Spanish, Russian, and Chinese to construct the test set. However, the number of languages covered in Guo et al. (2019) is relatively small compared to the number of languages around

---

<sup>1</sup><https://github.com/huggingface/peft>

<sup>2</sup><https://github.com/stopwords-iso/stopwords-iso>

the world, over 7000 languages<sup>3</sup>. Therefore, we expand the languages covered in the task of multilingual paraphrasing to 13 languages, including some figurative languages such as Japanese, and Chinese or accented languages such as Vietnamese. The proposed evaluation dataset has two types to assess different aspects of the model: **Input-only** and **Input-reference**.

**Input-only:** We follow Thompson and Post (2020b) and use the validation set from the WMT19 dataset, which is released in 2019 and available on HuggingFace<sup>4</sup>. Since the WMT19 dataset is used for the task Machine Translation, we thus, extract one side of the dataset in order to build the **Input-only** evaluation dataset. The languages we extracted from the WMT19 are Czech (cs), German (de), English (en), Finnish (fi), French (fr) and Chinese (zh). Most of the samples in this evaluation dataset are of the news domains, which is used to test the model’s ability on producing a paraphrase that conveys the same meaning.

**Input-reference:** We use in total 3 datasets to construct this evaluation set:

- **PAWS-X** (Yang et al. 2019) is the cross-lingual paraphrase identification dataset, thus, we extract only the sentence pairs with label 1 (indicating paraphrase) for 6 languages: Japanese (ja), Chinese (zh), German (de), French (fr), Spanish (es) and English (en).
- **Opusparcus** (Creutz 2019): is a paraphrase corpus for six European languages: German (de), English (en), Finnish (fi), French (fr), Russian (ru), and Swedish (sv). We extract the test set of Opusparcus with a score of 4 to ensure high-quality sentence pairs.
- **STAPLE** (Duolingo 2020): is a multi-reference machine translation dataset in which each reference could be viewed as the paraphrase. Since STAPLE does not have the validation or test set, we randomly extract 1000 samples, each with 5 reference texts, covering 3 languages: Vietnamese (vi), Portuguese (pt) and Hungarian (hu), to construct the first multi-reference paraphrase generation corpus with 1000 samples and 4 references each.

### Automatic Evaluation

We follow Chowdhury, Zhuang, and Wang (2022) to report the result on BLEU (Papineni et al. 2002), Self-BLEU, Self-TER, which adapted TER (Snover et al. 2006) to the input instead of the reference, BERTScore (Zhang et al. 2020) with `bert-base-multilingual-cased` checkpoint, iBLEU with  $\alpha = 0.7$  following Hosking and Lapata (2021). In addition, we further use two latest paraphrase metrics: ParaScore (Shen et al. 2022) and BERT-iBLEU (Niu et al. 2021).

### Human Evaluation

In addition to machine evaluation, we also conduct the human evaluation of the paraphrase generated by our model. For each of the following languages: English, Vietnamese, German, French and Japanese, we randomly extract 200

sentence triples of the input sentence, our model prediction and the output from the model of Thompson and Post (2020b). For each mentioned language, we ask 5 annotators to score 200 sentence pairs independently. Each annotator is instructed to rate on the 1-5 scale (with 5 being the highest) based on 3 criteria: **(i) Semantic preservation**, evaluating how much information is preserved in the output; **(ii) Lexical similarity**, evaluating how much similar in term of syntax or word choices of the output compared to input; and **(iii) Fluency**, assessing the fluency and coherence of the generated output. The annotators’ agreement is measured using Krippendorff’s alpha (Krippendorff 1970), which provides a measure of inter-annotator reliability.

## Results

### Main Results

Method	en	es	zh	ru
<i>BERTScore</i> ↑				
DAE	79.25	80.56	78.91	77.83
MMT	84.93	82.68	<b>84.79</b>	81.32
LAMPAT	<b>86.86</b>	<b>84.35</b>	83.26	<b>84.01</b>
<i>Self-BLEU</i> ↓				
DAE	20.35	30.49	20.38	20.61
MMT	19.89	28.57	<b>10.32</b>	15.19
LAMPAT	<b>19.46</b>	<b>20.16</b>	14.95	<b>12.46</b>
<i>Self-TER</i> ↑				
DAE	50.48	45.19	45.92	48.31
MMT	52.45	43.16	41.25	45.68
LAMPAT	<b>61.32</b>	<b>55.43</b>	<b>55.28</b>	<b>50.67</b>
<i>BERT-iBLEU</i> ↑				
DAE	79.33	78.08	79.05	78.14
MMT	83.92	80.16	<b>85.72</b>	81.99
LAMPAT	<b>85.52</b>	<b>83.41</b>	83.61	<b>84.69</b>
<i>ParaScore</i> ↑				
DAE	88.75	90.46	89.37	88.50
MMT	89.45	91.56	90.05	88.47
LAMPAT	<b>92.95</b>	<b>92.96</b>	<b>90.64</b>	<b>91.92</b>

Table 2: Multilingual paraphrase generation test results over 4 languages English, Spanish, Chinese and Russian from the work of DAE (Guo et al. 2019).

Method	fr	cs	fi	de	en	zh
<i>BERTScore</i> ↑						
DAE	70.39	78.45	80.96	69.32	65.68	<b>81.35</b>
MMT	74.40	80.20	81.20	71.20	65.48	80.29
LAMPAT	<b>85.42</b>	<b>83.92</b>	<b>84.91</b>	<b>86.16</b>	<b>82.59</b>	79.18
<i>Self-TER</i> ↑						
DAE	45.53	39.97	54.50	46.18	30.42	28.49
MMT	48.40	37.12	56.74	48.42	33.76	31.56
LAMPAT	<b>49.53</b>	<b>40.49</b>	<b>63.92</b>	<b>53.76</b>	<b>39.31</b>	<b>40.80</b>
<i>BERT-iBLEU</i> ↑						
DAE	71.10	79.56	81.38	71.24	67.36	<b>81.05</b>
MMT	75.48	78.81	82.15	74.16	66.08	69.70
LAMPAT	<b>84.63</b>	<b>85.22</b>	<b>86.03</b>	<b>85.38</b>	<b>84.08</b>	79.90
<i>ParaScore</i> ↑						
DAE	89.91	88.46	75.93	87.43	88.56	88.12
MMT	89.95	90.10	81.15	89.47	91.49	90.35
LAMPAT	<b>93.24</b>	<b>92.45</b>	<b>86.82</b>	<b>93.21</b>	<b>94.84</b>	<b>92.24</b>

Table 3: Multilingual paraphrase generation test results on our input-only evaluation dataset.

<sup>3</sup><https://www.ethnologue.com/statistics>

<sup>4</sup><https://huggingface.co/datasets/wmt19>

Method	sv	fi	en	fr	de	ru	ja	zh	es	hu	pt	vi
BERTScore ↑												
DAE	80.20	79.50	84.00	85.20	89.44	89.90	74.50	<b>83.00</b>	88.80	75.30	78.40	80.50
MMT	83.48	84.92	75.86	77.01	82.53	81.49	72.61	71.93	82.45	73.11	83.44	75.76
LAMPAT	<b>85.47</b>	<b>85.47</b>	<b>94.87</b>	<b>89.94</b>	<b>90.10</b>	<b>92.16</b>	<b>90.77</b>	81.00	<b>92.87</b>	<b>86.16</b>	<b>91.99</b>	<b>87.70</b>
BLEU ↑												
DAE	5.22	4.20	10.50	15.51	14.44	5.46	20.62	<b>47.80</b>	18.58	8.10	16.01	15.05
MMT	1.22	0.39	11.79	6.10	11.98	0.77	9.55	16.69	20.49	0.84	11.09	3.58
LAMPAT	<b>6.04</b>	<b>4.90</b>	<b>23.07</b>	<b>20.55</b>	<b>21.65</b>	<b>6.48</b>	<b>29.13</b>	30.52	<b>26.95</b>	<b>8.16</b>	<b>16.65</b>	<b>20.24</b>
Self-BLEU ↓												
DAE	24.58	17.56	45.50	30.20	25.75	36.82	50.78	50.22	40.69	25.82	<b>30.44</b>	27.26
MMT	25.72	16.76	50.33	23.47	35.65	39.67	45.62	63.32	36.65	27.82	34.70	27.24
LAMPAT	<b>23.68</b>	<b>14.53</b>	<b>43.47</b>	<b>19.55</b>	<b>24.54</b>	<b>30.25</b>	<b>40.56</b>	<b>47.58</b>	<b>29.00</b>	<b>22.94</b>	31.43	<b>25.51</b>
iBLEU ↑												
DAE	0.05	0.01	-0.04	-0.05	-0.01	0.01	0.08	0.10	0.04	0.05	0.01	-0.01
MMT	0.01	0.02	0.02	<b>0.04</b>	0.04	-0.03	0.06	0.12	0.09	0.06	0.11	<b>0.14</b>
LAMPAT	<b>0.08</b>	<b>0.04</b>	<b>0.04</b>	0.02	<b>0.05</b>	<b>0.04</b>	<b>0.09</b>	<b>0.15</b>	<b>0.13</b>	<b>0.15</b>	<b>0.15</b>	0.05
Self-TER ↑												
DAE	50.18	<b>60.57</b>	24.55	55.34	57.44	33.14	40.27	30.75	31.55	40.57	47.81	50.55
MMT	47.16	53.59	23.90	47.07	43.48	29.27	40.88	26.85	30.97	36.88	27.10	39.95
LAMPAT	<b>57.20</b>	59.45	<b>56.56</b>	<b>56.49</b>	<b>61.33</b>	<b>47.80</b>	<b>42.62</b>	<b>35.34</b>	<b>38.55</b>	<b>41.51</b>	<b>55.11</b>	<b>57.91</b>
BERT-iBLEU ↑												
DAE	70.20	78.20	66.60	56.20	62.70	68.20	52.00	67.20	57.20	67.50	68.80	77.20
MMT	61.43	68.44	65.25	66.33	63.07	69.06	65.33	62.67	63.95	72.39	72.34	69.91
LAMPAT	<b>73.30</b>	<b>78.60</b>	<b>67.60</b>	<b>82.33</b>	<b>80.25</b>	<b>69.98</b>	<b>78.04</b>	<b>67.72</b>	<b>86.46</b>	<b>76.32</b>	<b>78.92</b>	<b>81.39</b>
ParaScore ↑												
DAE	82.00	72.80	85.00	85.30	88.90	90.20	84.50	87.60	90.10	80.50	80.52	80.02
MMT	83.17	81.72	76.83	77.62	82.59	81.69	73.92	73.07	82.70	74.54	83.84	77.01
LAMPAT	<b>86.42</b>	<b>86.00</b>	<b>94.51</b>	<b>90.73</b>	<b>90.29</b>	<b>91.47</b>	<b>91.33</b>	<b>93.01</b>	<b>92.96</b>	<b>86.89</b>	<b>91.74</b>	<b>89.78</b>

Table 4: Multilingual paraphrase generation test results on our input-reference evaluation dataset.

For the test dataset from Guo et al. (2019) and our **input-only** evaluation, we evaluate using BERTScore, Self-BLEU, Self-TER, BERT-iBLEU and ParaScore, as these evaluation metrics do not require the reference text. The test results are depicted by Table 2 and 3. For **input-reference**, we report BLEU and iBLEU, in addition, as depicted by Table 4. LAMPAT can generate diverse output compared to the input, which is indicated by the low score in Self-BLEU and high score in Self-TER. LAMPAT also preserves better information from the input as demonstrated by the BERTScore results. We have manually examined the text generated by all four methods by randomly selecting 100 samples per language. Although these sentences convey the same meaning, the word choices and syntax structures are largely different from both input and reference, leading to the low score of iBLEU. Overall, BERT-iBLEU and ParaScore metrics, which are the metrics that grade both lexical and semantic aspects, show that all three methods could generate comprehensive sentences. However, our method, LAMPAT, still achieves the highest score over 13 languages we have tested. Even though our model learned in an unsupervised manner, it can still outperform the supervised counterpart. Overall, LAMPAT has demonstrated that integrating adversarial training into unsupervised learning could improve the performance of multilingual paraphrase generation.

## Human Evaluation Results

The results of our human evaluation can be found in Table 5, where we present the assessments provided by human evaluators. The evaluations conducted by human experts provide valuable insights into the performance of our model.

Method	SP ↑	LS ↓	F ↑
MMT	3.2	4.8	4.8
LAMPAT	4.2	3.2	4.8
Human generated paraphrase	4.4	2.4	4.8
Krippendorff’s alpha	0.68	0.7	0.78

Table 5: Human evaluation results. SP: Semantic Preservation; LS: Lexical Similarity; F: Fluency. All the scores reported are the average value of 5 chosen languages.

## Ablation Study

### Parameter-Efficient Fine-Tuning

In order to examine which ingredients help improve the performance of LAMPAT, we experiment with Adapter (Houlsby et al. 2019), Prefix Tuning (Li and Liang 2021), Prompt Tuning (Lester, Al-Rfou, and Constant 2021), LoRA (Hu et al. 2021), P-Tuning (Liu et al. 2022) and Full fine-tuning to find out which PEFT methods result in a more stable and better result. For each of the methods, we train on the

same number of epochs with 3 random seeds and report the mean and standard deviation of ParaScore of all languages.

Method	ParaScore
Prefix Tuning	77.45±8.5
Prompt Tuning	82.25±7.9
P-Tuning	82.94±5.5
Adapter	89.56±3.2
LoRA	<b>90.67±2.5</b>
Full	89.46±1.5

Table 6: ParaScore of different fine-tuning methods on 13 languages. The mean and standard deviation are the weighted mean and standard deviation in all 13 languages.

According to Table 6, LoRA experiences to be the most stable method and achieves the highest scores, especially for generation tasks.

## Adversarial Optimization

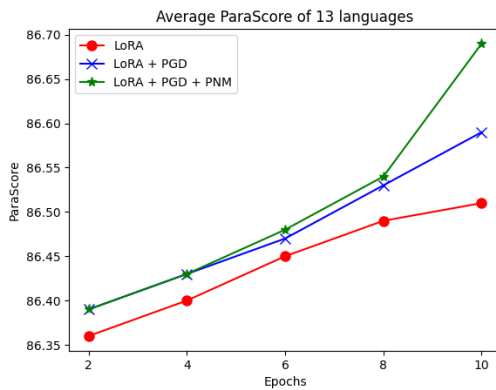


Figure 3: The average ParaScore of each technique over 13 languages.

Since LoRA is a stable method for partially fine-tuning LLMs, we adapt LoRA as the PEFT method for fine-tuning mGPT. In order to study the effect of adversarial training in unsupervised learning, we employ two settings which are Projected Gradient Descent (PGD) and Projected Newton Method (PNM) (Bertsekas 1982) together with LoRA. According to Figure 3, we hypothesized that the main driving force that makes LoRA + PGD + PNM better compared to PGD is because it uses the Taylor expansion, which has a better approximation of the objective function. In general, when we are near the local optima (for example, at the last 2 epochs), we can take a few more Newton’s steps to reach the optimum point instead of taking many small Gradient Descent steps.

## Related Works

### Multilingual Paraphrasing

Numerous techniques for paraphrasing in multiple languages employ Machine Translation-based models (MNMT). To illustrate, Thompson and Post (2020b) applied a pretrained MNMT model introduced by Thompson and

Post (2020a), along with a customized decoding algorithm aimed at reducing repetition of words and encouraging diverse vocabulary usage. Another instance is the work by Guo et al. (2019), which leveraged a language model pretraining task adapted from Conneau and Lample (2019). During inference, the same language code is provided to the model, and the sequence is generated sequentially in an autoregressive manner. Despite the fact that the translation-based approach produces high-quality and fluent paraphrases, it faces certain inherent challenges. Firstly, there is a potential for bias stemming from the dominance of certain languages used for training the MNMT model. Secondly, there may be instances of incorrect synthetic paraphrases due to the inherent ambiguity in the pivot language, as pointed out by Thompson and Post (2020b). These issues need careful consideration in the development of multilingual paraphrasing methods.

### Adversarial Training

Adversarial training is a powerful technique employed in the development of resilient neural networks. While the computer vision community, as highlighted by Goodfellow, Shlens, and Szegedy (2015), has generally accepted that adversarial training can have a detrimental impact on model generalization, the scenario appears to be quite different for language models, as evidenced by studies such as Pereira et al. (2020) and Dong et al. (2021a). Incorporating adversarial training into large language models (LLMs), as explored by Miyato et al. (2018) and further elaborated upon by Dong et al. (2021b), has been found to yield improvements in both model generalization and robustness. An innovative training algorithm, denoted as YOPO (You Only Propagate Once), was proposed by Zhang et al. (2019). YOPO takes advantage of the “free” training strategies advocated by Shafahi et al. (2019) to diversify the training data by incorporating various adversarial samples while imposing different norm constraints. Also, Miyato et al. (2018) proposes a new training method that regularizes the training objective by using virtual labels in adversarial training. These approaches collectively showcase the effectiveness of adversarial training in enhancing both the robustness and generalization capabilities of neural networks.

## Conclusion and Future Work

In this research, we introduce an efficient method for generating paraphrases in multiple languages using Low-Rank Adaptation combined with virtual labeling during adversarial training. Importantly, our approach delivers satisfactory results without relying on supervised learning. Additionally, we contribute to the field by creating a novel multilingual multi-domain evaluation dataset. While LAMPAT has demonstrated proficiency in generating human-like paraphrases across various languages, it still requires improvements in handling idiomatic expressions. Furthermore, our evaluation dataset covers only 13 languages, leaving out many, especially those with limited resources. This highlights the ongoing need for research to enhance and expand the capabilities of multilingual paraphrase generation models.

## Acknowledgement

This research/project is supported by the National Research Foundation, Singapore under its AI Singapore Programme, AISG Award No: AISG2-TC-2022-005.

## References

- Bertsekas, D. P. 1982. Projected Newton Methods for Optimization Problems with Simple Constraints. *SIAM Journal on Control and Optimization*, 20(2): 221–246.
- Cao, Z.; Luo, C.; Li, W.; and Li, S. 2017. Joint Copying and Restricted Generation for Paraphrase. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, AAAI’17, 3152–3158. AAAI Press.
- Chowdhury, J. R.; Zhuang, Y.; and Wang, S. 2022. Novelty Controlled Paraphrase Generation with Retrieval Augmented Conditional Prompt Tuning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(10): 10535–10544.
- Conneau, A.; and Lample, G. 2019. Cross-lingual Language Model Pretraining. In Wallach, H.; Larochelle, H.; Beygelzimer, A.; d’Alché-Buc, F.; Fox, E.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Creutz, M. 2019. Open subtitles paraphrase corpus for six languages. *LREC 2018 - 11th International Conference on Language Resources and Evaluation*, (2005): 1364–1369.
- Dong, X.; Luu, A. T.; Ji, R.; and Liu, H. 2021a. Towards robustness against natural language word substitutions. *arXiv preprint arXiv:2107.13541*.
- Dong, X.; Luu, A. T.; Lin, M.; Yan, S.; and Zhang, H. 2021b. How should pre-trained language models be fine-tuned towards adversarial robustness? *Advances in Neural Information Processing Systems*, 34: 4356–4369.
- Duolingo. 2020. Data for the 2020 Duolingo Shared Task on Simultaneous Translation And Paraphrase for Language Education (STAPLE).
- Federmann, C.; Elachqar, O.; and Quirk, C. 2019. Multilingual Whispers: Generating Paraphrases with Translation. In *Proceedings of the 5th Workshop on Noisy User-generated Text (W-NUT 2019)*, 17–26. Hong Kong, China: Association for Computational Linguistics.
- Foundation, W. 2019. ACL 2019 Fourth Conference on Machine Translation (WMT19), Shared Task: Machine Translation of News.
- Freitag, M.; Foster, G.; Grangier, D.; and Cherry, C. 2020. Human-Paraphrased References Improve Neural Machine Translation. In *Proceedings of the Fifth Conference on Machine Translation*, 1183–1192. Online: Association for Computational Linguistics.
- Gan, W. C.; and Ng, H. T. 2019. Improving the Robustness of Question Answering Systems to Question Paraphrasing. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 6065–6075. Florence, Italy: Association for Computational Linguistics.
- Goodfellow, I. J.; Shlens, J.; and Szegedy, C. 2015. Explaining and Harnessing Adversarial Examples. arXiv:1412.6572.
- Guo, Y.; Liao, Y.; Jiang, X.; Zhang, Q.; Zhang, Y.; and Liu, Q. 2019. Zero-Shot Paraphrase Generation with Multilingual Language Models. arXiv:1911.03597.
- Hosking, T.; and Lapata, M. 2021. Factorising Meaning and Form for Intent-Preserving Paraphrasing. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 1405–1418. Online: Association for Computational Linguistics.
- Houlsby, N.; Giurgiu, A.; Jastrzebski, S.; Morrone, B.; De Laroussilhe, Q.; Gesmundo, A.; Attariyan, M.; and Gelly, S. 2019. Parameter-Efficient Transfer Learning for NLP. In Chaudhuri, K.; and Salakhutdinov, R., eds., *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, 2790–2799. PMLR.
- Hu, E. J.; Shen, Y.; Wallis, P.; Allen-Zhu, Z.; Li, Y.; Wang, S.; Wang, L.; and Chen, W. 2021. LoRA: Low-Rank Adaptation of Large Language Models. arXiv:2106.09685.
- Kaushik, P.; Gain, A.; Kortylewski, A.; and Yuille, A. 2021. Understanding Catastrophic Forgetting and Remembering in Continual Learning with Optimal Relevance Mapping. arXiv:2102.11343.
- Krippendorff, K. 1970. Estimating the Reliability, Systematic Error and Random Error of Interval Data. *Educational and Psychological Measurement*, 30: 61 – 70.
- Lester, B.; Al-Rfou, R.; and Constant, N. 2021. The Power of Scale for Parameter-Efficient Prompt Tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 3045–3059. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics.
- Lewis, M.; Ghazvininejad, M.; Ghosh, G.; Aghajanyan, A.; Wang, S.; and Zettlemoyer, L. 2020. Pre-training via Paraphrasing. In Larochelle, H.; Ranzato, M.; Hadsell, R.; Balcan, M.; and Lin, H., eds., *Advances in Neural Information Processing Systems*, volume 33, 18470–18481. Curran Associates, Inc.
- Li, X. L.; and Liang, P. 2021. Prefix-Tuning: Optimizing Continuous Prompts for Generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 4582–4597. Online: Association for Computational Linguistics.
- Liu, X.; Ji, K.; Fu, Y.; Tam, W.; Du, Z.; Yang, Z.; and Tang, J. 2022. P-Tuning: Prompt Tuning Can Be Comparable to Fine-tuning Across Scales and Tasks. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, 61–68. Dublin, Ireland: Association for Computational Linguistics.
- Miyato, T.; Ichi Maeda, S.; Koyama, M.; and Ishii, S. 2018. Virtual Adversarial Training: A Regularization



- Method for Supervised and Semi-Supervised Learning. arXiv:1704.03976.
- Niu, T.; Yavuz, S.; Zhou, Y.; Keskar, N. S.; Wang, H.; and Xiong, C. 2021. Unsupervised Paraphrasing with Pretrained Language Models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 5136–5150. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics.
- Papineni, K.; Roukos, S.; Ward, T.; and Zhu, W.-J. 2002. Bleu: a Method for Automatic Evaluation of Machine Translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, 311–318. Philadelphia, Pennsylvania, USA: Association for Computational Linguistics.
- Pereira, L.; Liu, X.; Cheng, F.; Asahara, M.; and Kobayashi, I. 2020. Adversarial Training for Commonsense Inference. In *Proceedings of the 5th Workshop on Representation Learning for NLP*, 55–60. Online: Association for Computational Linguistics.
- Shafahi, A.; Najibi, M.; Ghiasi, M. A.; Xu, Z.; Dickerson, J.; Studer, C.; Davis, L. S.; Taylor, G.; and Goldstein, T. 2019. Adversarial training for free! In Wallach, H.; Larochelle, H.; Beygelzimer, A.; d'Alché-Buc, F.; Fox, E.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Shen, L.; Liu, L.; Jiang, H.; and Shi, S. 2022. On the Evaluation Metrics for Paraphrase Generation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 3178–3190. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics.
- Snover, M.; Dorr, B.; Schwartz, R.; Micciulla, L.; and Makhoul, J. 2006. A Study of Translation Edit Rate with Targeted Human Annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, 223–231. Cambridge, Massachusetts, USA: Association for Machine Translation in the Americas.
- Thompson, B.; and Post, M. 2020a. Automatic Machine Translation Evaluation in Many Languages via Zero-Shot Paraphrasing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 90–121. Online: Association for Computational Linguistics.
- Thompson, B.; and Post, M. 2020b. Paraphrase Generation as Zero-Shot Multilingual Translation: Disentangling Semantic Similarity from Lexical and Syntactic Diversity. In *Proceedings of the Fifth Conference on Machine Translation*, 561–570. Online: Association for Computational Linguistics.
- Yang, Y.; Zhang, Y.; Tar, C.; and Baldridge, J. 2019. PAWS-X: A Cross-lingual Adversarial Dataset for Paraphrase Identification. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 3687–3692. Hong Kong, China: Association for Computational Linguistics.
- Zhang, D.; Zhang, T.; Lu, Y.; Zhu, Z.; and Dong, B. 2019. You Only Propagate Once: Accelerating Adversarial Training via Maximal Principle. In Wallach, H.; Larochelle, H.; Beygelzimer, A.; d'Alché-Buc, F.; Fox, E.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Zhang, T.; Kishore, V.; Wu, F.; Weinberger, K. Q.; and Artzi, Y. 2020. BERTScore: Evaluating Text Generation with BERT. In *International Conference on Learning Representations*.