

PEMA: An Offsite-Tunable Plug-in External Memory Adaptation for Language Models

Anonymous ACL submission

Abstract

Pre-trained language models (PLMs) show impressive performance in various downstream NLP tasks. However, pre-training large language models demands substantial memory and training compute. Furthermore, due to the substantial resources required, many PLM weights are confidential. Consequently, users are compelled to share their data with model owners for fine-tuning specific tasks. To overcome the limitations, we introduce Plug-in External Memory Adaptation (PEMA), a Parameter-Efficient Fine-Tuning (PEFT) method, enabling PLM fine-tuning without requiring access to all the weights. PEMA integrates with context representations from test data during inference to perform downstream tasks. It uses external memory to store PLM-generated context representations mapped with target tokens. Our method utilizes LoRA-based weight matrices in the PLM’s final layer to enhance efficiency. Our approach also includes Gradual Unrolling, a novel interpolation strategy to improve generation quality. We validate PEMA’s effectiveness through experiments on syntactic and real datasets for machine translation and style transfer. Our findings show that PEMA outperforms other PEFT approaches in memory and latency efficiency for training, and also excels in maintaining sentence meaning and generating appropriate language and styles.

1 Introduction

Pre-trained language models (PLMs) are widely used in downstream NLP tasks (Devlin et al., 2019a). Recent developments in large language models have shown remarkable performance in zero-shot and few-shot learning scenarios (Brown et al., 2020; Hendy et al., 2023; OpenAI, 2023b; Anil et al., 2023; Chowdhery et al., 2022). However, fine-tuning is still required to optimize the performance of the NLP tasks such as machine translation (Üstün and Cooper Stickland, 2022;

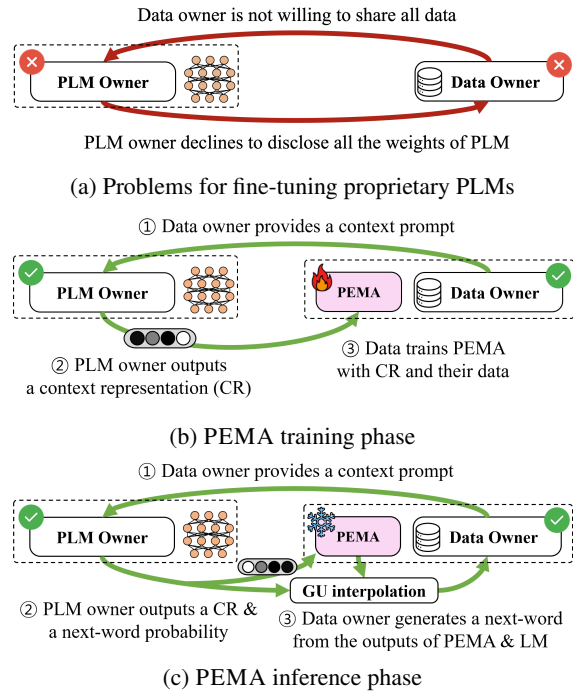


Figure 1: A motivation for PEMA. (a) The data owners who want to fine-tune PLMs encounter a problem when the PLM owner refuses to share all the weights of the PLM. (b) In the PEMA training phase, the data owner takes a CR from the PLM owner by providing a context prompt. They subsequently train their PEMA model with their dataset. (c) At inference, the data owner takes a CR for test data from the PLM owner. Using Gradual Unrolling (GU), they generate the next-token by interpolating between PEMA and PLM next-token probabilities.

Huang et al., 2020; Ding et al., 2022). The most straightforward approach to fine-tuning is full fine-tuning (Raffel et al., 2020; Qiu et al., 2020), which involves fine-tuning all parameters in a PLM. Yet, this approach requires substantial resources regarding memory and training compute (Iyer et al., 2022; Zhang et al., 2022; Touvron et al., 2023). To overcome this limitation, researchers have proposed Parameter-Efficient Fine-Tuning (PEFT) methods

051 to fine-tune a full model efficiently. Adapter tun- 103
052 ing (Pfeiffer et al., 2021; He et al., 2021; Houlsby 104
053 et al., 2019) utilizes small, additional parameters 105
054 known as adapters inserted between layers within 106
055 a PLM. On the other hand, LoRA (Hu et al., 2022) 107
056 uses trainable low-rank matrices that incrementally 108
057 update the pre-trained weights. These fine-tuning 109
058 methods require access to all the weights of PLMs. 110

059 However, proprietary PLMs such as Chat- 111
060 GPT (OpenAI, 2022), Bard (Pichai, 2023), and 112
061 Claude (AnthropicAI, 2023) are confidential. 113
062 Hence, the owners of these PLMs do not reveal 114
063 all the model weights. Consequently, data owners 115
064 possessing their datasets and wishing to fine-tune 116
065 proprietary PLMs for specific downstream tasks 117
066 must provide their datasets to the PLM owners for 118
067 fine-tuning (OpenAI, 2023a). However, this pro- 119
068 cess can be challenging due to the confidential na- 120
069 ture of the datasets, which may involve privacy con- 121
070 cerns (Guinney and Saez-Rodriguez, 2018). Fig- 122
071 ure 1a shows problems for fine-tuning proprietary 123
072 PLMs. To overcome this situation, (Xiao et al., 124
073 2023) proposes the offsite-tuning approach that 125
074 uses one-third of the middle layers of a PLM, re- 126
075 ferred to as the emulator. Nevertheless, this ap- 127
076 proach still needs a large parameter size, and com-
077 pressing the full model into an emulator requires a
078 computationally intensive distillation process.

079 To address the challenges mentioned above, we
080 introduce a novel PEFT method named Plug-in
081 External Memory Adaptation (PEMA) designed
082 for efficient fine-tuning of proprietary PLMs in
083 machine translation tasks. PEMA utilizes LoRA-
084 based weight matrices designed for learning down-
085 stream tasks with accessible features provided by
086 OpenAI API (OpenAI, 2022) and minimal part of
087 PLM’s weight (language model head).

088 In the training phase, the data owner begins the
089 process by providing a prompt with initial input
090 to the PLM owner, which includes an instruction
091 and a source sentence from a parallel corpus. The
092 PLM owner receives this initial input to generate a
093 context representation and predict the next-token.
094 Then, it iteratively processes subsequent inputs con-
095 taining the predicted next-tokens. This approach
096 avoids the need for the full dataset from the data
097 owner. Throughout this process, the data owner
098 builds an external memory comprised of context
099 representations and corresponding desired target
100 tokens. They train PEMA by reconstructing the
101 stored context representations and predicting target
102 tokens based on these representations. Figure 1b

shows the training phase process of PEMA.

103 During the inference phase, the data owner uses 104
105 a prompt to request a context representation for 106
107 test data from the PLM owner. The PLM owner 107
108 then outputs a context representation and a next- 108
109 token probability given the prompt. PEMA also 109
110 outputs a next-token probability based on a con- 110
111 text representation. These probabilities are interpo- 111
112 lated to compute a final next-token probability. We 112
113 propose Gradual Unrolling (*GU*), an interpolation 113
114 strategy that initially emphasizes PEMA’s distri- 114
115 bution, gradually shifts to the PLM’s context-based 115
116 predictions as the sentence progresses. Figure 1c 116
117 illustrates the inference phase process of PEMA.

117 We evaluate PEMA by comparing it with other 117
118 PEFT methods. PEMA shows better resource ef- 118
119 ficiency, consuming less GPU memory and run- 119
120 ning faster. Additionally, PEMA outperforms other 120
121 baselines in translating English sentences into Ger- 121
122 man and paraphrasing informal sentences into for- 122
123 mal ones while preserving the original meaning. 123
124 Lastly, we conduct ablation studies to assess the ef- 124
125 fectiveness of each component of PEMA. PEMA 125
126 is publicly available for further exploration into 126
127 offsite-tunable efficient fine-tuning.¹ 127

2 Related Work 128

2.1 Parameter-Efficient Fine-Tuning 129

129 Parameter-Efficient Fine-Tuning aims to fine-tune 130
131 PLMs to address resource constraints in memory 131
132 and training compute. (Iyer et al., 2022; Zhang 132
133 et al., 2022; Touvron et al., 2023). Several ap- 133
134 proaches have been proposed to overcome this lim- 134
135 itation. Adapter tuning (Pfeiffer et al., 2021; He 135
136 et al., 2021; Houlsby et al., 2019) inserts small 136
137 parameters, known as adapters, between layers 137
138 within a PLM. Prefix and Prompt tuning (Li and 138
139 Liang, 2021; Liu et al., 2021; Lester et al., 2021) 139
140 incorporate additional trainable prefix tokens to a 140
141 PLM’s input or hidden layers. Low-Rank Adap- 141
142 tation (LoRA) (Hu et al., 2022) uses trainable 142
143 low-rank matrices, denoted as B and A , that in- 143
144 crementally update PLM weights. B and A are 144
145 reduced to a low-rank r . This adaptation can be 145
146 mathematically represented as transitioning from 146
147 $h = W_0x$ to $h = W_0x + \Delta Wx = W_0x + BAx$, 147
148 where $W_0 \in \mathbb{R}^{k \times d}$, $B \in \mathbb{R}^{k \times r}$, and $A \in \mathbb{R}^{r \times d}$. 148
149 UniPELT (Mao et al., 2022) combines multiple 149
150 PEFT methods, using a gating mechanism to acti- 150
151 vate the most suitable components for given data or 151

¹The Github repository link will be provided after review.

tasks. We propose a novel adaptation method that leverages LoRA parameters and is offsite-tunable.

2.2 k -Nearest Neighbors Language Model

The k -Nearest Neighbors Language Model (k NN-LM) estimates the next-token distribution by interpolating the output distributions from a pre-trained language model (P_{LM}), and an external memory (P_{kNN}) (Khandelwal et al., 2020). The memory is used to perform a k NN search and to integrate out-of-domain data, thereby enabling a single language model to be adaptive across various domains. Given a context represented as a sequence of tokens $c_i = (w_1, \dots, w_{i-1})$, the k NN-LM utilizes a pre-trained language model $f(\cdot)$ to generate a context representation $f(c_i)$. This representation is then paired with the desired target token y_i to create the external memory (referred to as a datastore in (Khandelwal et al., 2020)) $\{(f(c_i), y_i) | (c_i, y_i) \in \mathcal{E}\}$ from the training dataset \mathcal{E} . The next-token distribution from the external memory, P_{kNN} , is computed using a k -nearest neighborhood approach with the squared L^2 distance. The final next-token distribution is then obtained by interpolating between P_{kNN} and P_{LM} as: $P(y_i|c_i) = \lambda P_{kNN}(y_i|c_i) + (1 - \lambda)P_{LM}(y_i|c_i)$.

We adapt the concept of external memory and interpolation of different next-token distributions to PEMA. Instead of employing a k NN-based approach, we employ a neural network-based model that directly learns to estimate the next-token, which is more effective in mitigating overfitting to the training data. Additionally, we use the Gradual Unrolling interpolation strategy to enhance the quality of interpolation. The k NN-LM method relies on k NN for external memory search to adapt the language model to diverse domains. However, it is well known that the non-parametric model k NN can potentially overfit. Therefore, it often requires a large amount of training data. To address this, we introduce a parametric approach within PEMA to improve its performance on downstream tasks. This approach is better suited for limited training data scenarios. It involves replacing the existing k NN with a parametric model in PEMA, thus enabling effective adaptation to various domains in terms of performance.

3 Plug-in External Memory Adaptation

This section describes Plug-in External Memory Adaptation (PEMA), which aims to fine-tune a

pre-trained language model without requiring a full model during training. PEMA is integrated into the language model during inference to facilitate downstream NLP tasks. It uses external memory to build a context representation $f(c_i)$, mapped with the desired target token y_i . Using the external memory, we train PEMA in two phases. The first phase involves reconstruction training to reconstruct $f(c_i)$ with $B_{rct}A$, resulting in the output of a reconstruction loss. Subsequently, the joint retraining phase focuses on generating the next-token probability P_{PEMA} that predicts target token y_i given $Af(c_i)$ with B_{pd} . Simultaneously, it uses pre-trained B_{rct} to retain the original feature $f(c_i)$. During the inference stage, the next-token probabilities from both the pre-trained generative language model P_{LM} and PEMA P_{PEMA} are interpolated to generate the next-token. Figure 2 shows the structure of PEMA.

3.1 Building an External Memory

The first step of PEMA is to build an external memory. The output $f(c_i)$ represents a context representation obtained from the final layer’s feed-forward network output of a pre-trained language model.

For the i -th token training example in external memory $(c_i, y_i) \in \mathcal{E}$, a paired representation is created by defining an input prompt c_1 and a corresponding target token sequence. Predicted token sequences are generated by sequentially extending the input prompt. ❶ Initially, the input prompt c_1 is fed into the pre-trained language model, resulting in the predicted next-token \hat{w}_1 and ❷ the corresponding context representation $f(c_1)$. ❸ Including \hat{w}_1 in the input prompt extends it to the next context $c_2 = \{c_1, \hat{w}_1\}$, subsequently producing the next predicted token \hat{w}_2 and its context representation $f(c_2)$. This iterative process yields a sequence of context representations $(f(c_1), f(c_2), \dots, f(c_t = \{c_1, \hat{w}_1, \dots, \hat{w}_{t-1}\}))$ for training, with each context c_i corresponding to the i -th position in the token sequence and t denoting the total number of tokens in a token sequence of one sentence training example.

We map the context representation $f(c_i) \in \mathbb{R}^{1 \times d}$, where d is the size of the context representation with the target token y_i , resulting in the pair $(f(c_i), y_i)$. The external memory $(f(C), Y)$ is formed by collecting all such context and token pairs constructed from the training set \mathcal{E} as below:

$$(f(C), Y) = \{(f(c_i), y_i) | (c_i, y_i) \in \mathcal{E}\} \quad (1)$$

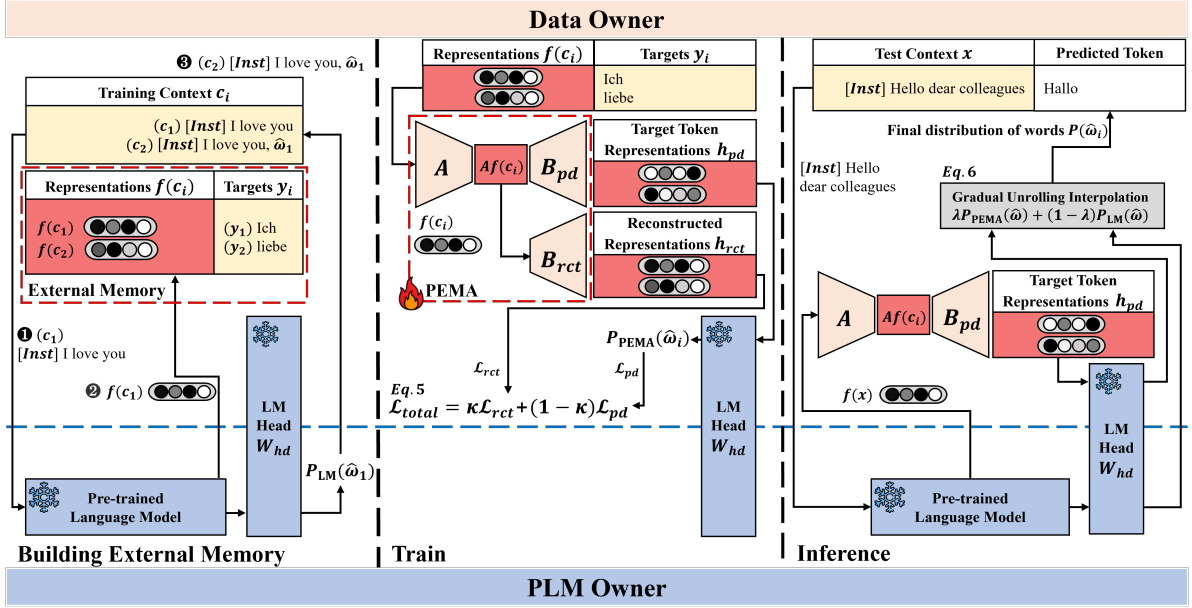


Figure 2: An illustration of PEMA. The areas of the PLM owner and the data owner are separated by the blue horizontal line. The data owner can train and infer using only the PLM’s LM head. PEMA builds an external memory from the training context with an instruction $[Inst]$ given to a PLM. The PLM outputs the representation $f(c_i)$ and predicts the next-token distribution $P_{LM}(\hat{w}_i)$. The representation $f(c_i)$ is then aligned with its target y_i . In the training phase, PEMA uses external memory for two tasks: preserving the original representation via reconstruction training with B_{rct} and generating a target token probability distribution using B_{pd} . For inference, the model inputs a test data representation to generate two probability distributions: $P_{LM}(\hat{w}_i)$ and $P_{PEMA}(\hat{w}_i)$. These are then interpolated using Gradual Unrolling to obtain the final token distribution.

3.2 PEMA Adaptation Model

We incorporate LoRA (Hu et al., 2022), a low-rank parameterization adaptation known for its effectiveness in various adaptation tasks, into PEMA for adapting to multiple text generation tasks.

The PEMA consists of three weight matrices: $A \in \mathbb{R}^{r \times d}$, $B_{rct} \in \mathbb{R}^{d \times r}$, and $B_{pd} \in \mathbb{R}^{d \times r}$ where d is the size of the context representation and r is a rank-size that $r < d$. Given $Af(c_i)$ where $f(c_i) \in \mathbb{R}^{1 \times d}$, B_{rct} is used to reconstruct the context representation input $f(c_i)$, with the goal of approximating $h_{rcti} \approx f(c_i)$. Additionally, B_{pd} is used to produce a representation h_{pd_i} that maximizes target token prediction when fed into the frozen weight of a language model head (LM head) $W_{hd} \in \mathbb{R}^{v \times d}$ where v is the vocabulary size that outputs the predicted next-token \hat{w}_i .

$$\begin{aligned} h_{rct_i} &= \Delta W_{rct} f(c_i) = B_{rct} f(c_i) \\ h_{pd_i} &= \Delta W_{pd} f(c_i) = B_{pd} f(c_i) \\ P_{PEMA}(\hat{w}_i | c_i) &= \text{softmax}(W_{hd} h_{pd_i}) \end{aligned} \quad (2)$$

3.3 Model Training

The training process consists of two distinct phases: initial reconstruction training to preserve the gen-

eral knowledge within the context representation of PLM and subsequent joint retraining, encompassing both the reconstruction of context representations and the prediction of next-tokens.

Initial Reconstruction Training. First, we train the decoder B_{rct} by reconstructing the i -th original context representation of the n -th sentence training example $f(c_i)^n$. We use a mean-square error loss between original input $f(c_i)^n$ and the output $h_{rct_i}^n$ as below:

$$\mathcal{L}_{rct} = \frac{1}{|\mathcal{E}|} \sum_{n=1}^{|\mathcal{E}|} \sum_{i=1}^{t_n} (f(c_i)^n - h_{rct_i}^n)^2 \quad (3)$$

where t_n is the number of tokens in a token sequence of n -th sentence training example and $|\mathcal{E}|$ is the size of the training dataset.

Joint Retraining After completing the initial reconstruction training, we proceed to the joint retraining phase, using the pre-trained B_{rct} and randomly initialized A . Our first objective is to acquire a representation $h_{pd_i}^n$ that is optimized for predicting the target token y_i^n . We utilize a cross-entropy loss based on the softmax function of the output of

$W_{hd}h_{pd_i}^n$ given the target token y_i^n as below:

$$\mathcal{L}_{pd} = -\frac{1}{|\mathcal{E}|} \sum_{n=1}^{|\mathcal{E}|} \sum_{i=1}^{t_n} y_i^n \log P_{PEMA}(y_i^n | W_{hd}h_{pd_i}^n) \quad (4)$$

The second objective is to reconstruct the input context representation x_i using the randomly initialized A and pre-trained B_{rct} with the reconstruction loss function as depicted in Equation 3. The reconstruction loss intends to retain the general knowledge obtained from the pre-trained language model while maximizing the target token prediction. We introduce a parameter κ that can be fine-tuned to adjust the emphasis on the objectives as below:

$$\mathcal{L}_{total} = \kappa \mathcal{L}_{rct} + (1 - \kappa) \mathcal{L}_{pd} \quad (5)$$

3.4 Model Inference

To generate the next-token \hat{w} , we exclude B_{rct} and use B_{pd} and A . The PLM receives the input context x from the test dataset, and generates $f(x)$, which serves as input for two pathways. One pathway uses PEMA’s A and B_{pd} to create h_{pd} for x . Subsequently, it is passed through W_{hd} to produce a distribution of the next-token $P_{PEMA}(\hat{w}|x)$. The other pathway directly feeds r into W_{hd} to produce the next-token distribution $P_{LM}(\hat{w}|x)$. Finally, these two distributions are blended using a tuned parameter λ to produce the final distribution of tokens for the desired task as below:

$$P(\hat{w}|x) = \lambda P_{PEMA}(\hat{w}|x) + (1 - \lambda) P_{LM}(\hat{w}|x) \quad (6)$$

4 Gradual Unrolling Interpolation

Given that an adaptation model trained with only a limited number of parameters may lack the context-awareness and language-generation capabilities of pre-trained language models, it is more effective to use the adaptation model to guide the generation of tokens of the desired task at the beginning of the sentence, and rely on a pre-trained language model to provide context for the rest of the sentence. To achieve this, we suggest the Gradual Unrolling strategy, which aims for strong $P_{PEMA}(\hat{w}|x)$ interpolation at the beginning of generation and gradually decreases the interpolation. As the sentence progresses, the pre-trained language model increasingly contributes to providing the necessary context, as shown in Figure 3.

In the context of sentence generation, we define SL as the input sentence length, excluding instruction and user-defined variables λ_{max} . λ represents the proportion of the adaptation model’s

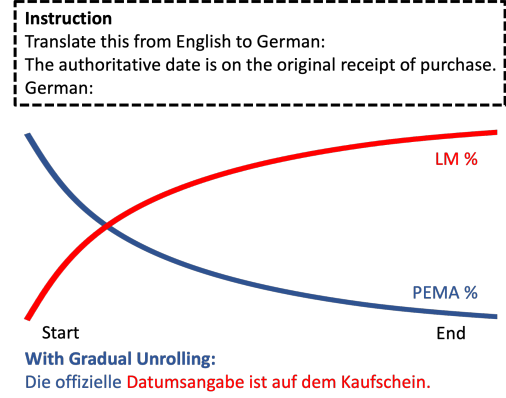


Figure 3: The intuition of Gradual Unrolling. Given the input sentence (Black), the interpolation percentage of the adaptation model (Blue) decreases gradually while that of the language model (Red) increases as the sentence is being generated. This strategy ensures that the adaptation model generates tokens trained for the desired task at the beginning of the sentence, and the language model provides the necessary context in the remaining part of the sentence.

interpolation ($0 \leq \lambda \leq 1$). We also have the dependent variables of the current step (CS) and the step size (SS). The step size is computed as $SS = \lambda_{max}/SL$, and CS is initialized to λ_{max} at the start of sentence generation. At each token generation step, CS decreases by SS until the end of the sentence (i.e., $CS_{cur} = CS_{past} - SS$ where CS_{past} is the latest token’s CS variable). Then, we calculate the current interpolation proportion λ_{cur} (i.e., λ at Equation 6) as $\lambda_{cur} = CS_{cur}^2$.

5 Experiments

This section describes the experiments and results to show both the computational efficiency and performance in downstream tasks of PEMA. First, we perform an experiment on the computational efficiency of PEMA. Subsequently, we evaluate PEMA across two downstream tasks: the WMT22 EN→DE machine translation task (Kocmi et al., 2022) and the GYAFC formal style transfer task (Rao and Tetreault, 2018). Lastly, we conduct an ablation study to show the gradual improvement by incorporating each idea of PEMA.

5.1 Computational Efficiency

To evaluate the computational efficiency of PEMA, we conduct a comparison of different fine-tuning methods based on their resource utilization during both training and inference. We follow the approach of previous work (Pope et al., 2023) that

employs a fixed size of input tensors. We use input tensors with the size [1, 10], equivalent to sequences of 10 tokens with OPT-IML-MAX-1.3B. The resource utilization metrics encompass training memory consumption, training latency, inference memory consumption, inference latency, and floating point operations per token.

The evaluation involves several steps. First, we clear the CUDA cache to compute the memory and ensure no background GPU processes. GPU memory utilization is determined using the `memory_summary` function provided by PyTorch (Paszke et al., 2019). We calculate the time difference before inputting the data into the model and after obtaining the output. For training latency, we consider the time encompassing the entire back-propagation process. To ensure the accuracy of latency, we compute the mean and variance based on ten trials of inputs for each fine-tuning method. We conducted a comparative analysis with the offsite-tuning baseline approach, Offsite-Tuning (Xiao et al., 2023). Offsite-Tuning involves knowledge distillation (OT Emulator) and downstream task training using the OT Emulator (OT Plug-in). Subsequently, it utilizes the OT Plug-in to interact with the PLM during the inference phase.

As shown in Table 1, PEMA demonstrates the efficiency by utilizing one-tenth of the training memory consumption compared to LoRA. In addition, PEMA shows the fastest training latency among all the methods. This is because PEMA uses external memory to store context representations and does not require access to a pre-trained language model during the training phase, as illustrated in Figure 2. These results highlight the significance of PEMA’s reduced training memory consumption and improved training latency, making it an appealing choice for efficient natural language generation tasks.

5.2 Performance of Downstream Tasks

We present a comprehensive analysis of the performance of PEMA and baseline models on two downstream tasks: the WMT22 (EN→DE) translation task and the GYAFC task involving Family & Relationships and Entertainment & Music.

For the machine translation task, we use the EN→DE news-commentary dataset to address the limitation noted in (Brown et al., 2020), where translations into English tend to be stronger than those from English due to training set biases. We evaluate our model using the latest test set provided

Method	Tr-MC	Tr-Lat	Inf-MC	Inf-Lat	FLOPs
FT	20,082	250.4 \pm 140.6	5,021	17.1 \pm 1.0	2.41e9
FT-top2	7,355	70.3 \pm 108.6	5,021	17.3 \pm 1.3	2.41e9
k NN-LM	None	20.3 \pm 567.2	5,021	37.5 \pm 1.4	FT+6.29e6
LoRA	5,056	21.6 \pm 0.4	5,031	20.5 \pm 1.5	FT+4.19e6
UniPELT	5,138	30.3 \pm 0.1	5,047	21.3 \pm 0.6	FT+1.49e7
OT Emulator	11,713	88.4 \pm 309.4	None	None	FT+8.03e8
OT Plug-in	5,267	59.6 \pm 107.8	5,269	21.3 \pm 0.1	FT+4.82e8
PEMA	478	18.5 \pm 1.0	5,043	18.2 \pm 0.5	FT+4.19e6

Table 1: Comparison of various training and inference resource utilization methods with OPT-IML-MAX-1.3B. We evaluate memory consumption (MC) and latency (Lat) for training (Tr) and inference (Inf), as well as FLOPs per token, using 10-token length sequences. Memory size is measured in megabytes, and latency is measured in milliseconds. PEMA stands out by using only one-tenth of the training memory utilized by LoRA. Furthermore, PEMA demonstrates the fastest training latency among the methods.

by (Hendy et al., 2023).

For the formality style transfer task, we use the GYAFC dataset (Rao and Tetreault, 2018), which consists of a parallel training set of informal and formal sentences. The test set comprises four reference sentences paired with one informal sentence. In this task, our objective is to transfer the style of informal sentences into formal ones.

We use three pre-trained language models: OPT-IML-MAX-1.3B, LLaMA-7B, and OPT-IML-MAX-30B (Iyer et al., 2022; Touvron et al., 2023). We compare PEMA with the following methods:

Full fine-tuning (FT) updates all pre-trained model parameters, including weights and biases.

Fine-tuning top-2 (FT-Top2) updates the last two layers while the remaining layers are frozen.

k -Nearest Neighbors Language Model (k NN-LM) (Khandelwal et al., 2020) uses k NN search within an external memory to derive a next-token distribution P_{kNN} , which is then interpolated with P_{LM} to produce an adapted next-token distribution.

LoRA (Hu et al., 2022) uses two additional trainable matrices. We apply LoRA at the last layer output projection matrices in the self-attention module.

UniPELT (Mao et al., 2022) is a state-of-the-art PEFT method that combines Adapter tuning (Houlsby et al., 2019), Prefix tuning (Li and Liang, 2021), and LoRA (Hu et al., 2022) with a gating mechanism to select the optimal approaches. We apply UniPELT at the last layer.

Offsite-Tuning (Xiao et al., 2023) is an offsite-tunable method that uses a distilled PLM emulator with an adapter, which includes multiple copies at

Model	Tr-MC (MB)	WMT22 (EN→DE)			GYAFC (F&R)			GYAFC (E&M)		
		sBLEU	PPL	COMET	sBLEU	PPL	FormImp	sBLEU	PPL	FormImp
OPT-1.3B	None	9.55	51.30	57.24	55.00	18.98	11.05	53.98	<u>20.89</u>	10.67
OPT-1.3B (FT)	20,082	<u>10.15</u>	<u>40.83</u>	<u>61.44</u>	29.17	24.82	<u>52.28</u>	31.50	27.99	46.82
OPT-1.3B (FT-Top2)	7,355	3.57	51.36	38.35	21.60	24.33	59.00	23.94	27.07	<u>51.52</u>
OPT-1.3B (kNN-LM)	None	8.07	91.37	41.75	56.69	20.87	16.26	54.74	23.15	14.46
OPT-1.3B (LoRA)	5,025	4.28	61.25	39.32	20.98	<u>19.07</u>	45.71	15.57	19.71	46.32
OPT-1.3B (UniPELT)	5,138	9.15	47.09	56.30	51.38	44.43	52.22	46.67	22.08	53.31
OPT-1.3B (Offsite-Tuning)	5,267	7.65	36.91	52.85	<u>59.01</u>	20.70	24.82	<u>57.01</u>	23.25	23.76
OPT-1.3B (PEMA)	478	12.87	<u>42.62</u>	64.16	64.82	23.15	41.90	61.24	24.28	36.28
LLaMA-7B	None	2.78	78.49	39.49	20.18	34.53	42.81	24.14	37.33	44.81
LLaMA-7B (kNN-LM)	None	0.07	85.09	38.53	1.72	41.50	55.13	1.94	46.31	<u>68.61</u>
LLaMA-7B (LoRA)	13,237	<u>11.46</u>	<u>51.36</u>	<u>67.48</u>	52.67	22.42	72.23	52.15	24.74	71.28
LLaMA-7B (UniPELT)	13,810	9.13	46.62	56.31	<u>59.81</u>	<u>22.95</u>	<u>71.69</u>	<u>58.07</u>	<u>25.35</u>	68.33
LLaMA-7B (PEMA)	996	14.50	54.26	70.31	63.99	23.19	61.40	60.88	26.00	60.94
OPT-30B	None	<u>18.22</u>	45.81	<u>77.41</u>	60.41	20.04	29.33	57.60	21.97	23.88
OPT-30B (kNN-LM)	None	16.65	74.06	62.98	61.02	<u>20.86</u>	29.80	58.58	<u>22.75</u>	23.39
OPT-30B (LoRA)	58,083	8.26	46.97	69.41	61.39	22.00	73.10	<u>59.76</u>	<u>23.97</u>	68.29
OPT-30B (UniPELT)	59,028	15.57	47.34	73.42	<u>64.54</u>	21.72	47.14	56.86	23.77	34.08
OPT-30B (PEMA)	1,909	19.22	<u>46.62</u>	79.21	70.84	22.04	<u>52.35</u>	65.43	25.53	<u>44.63</u>

Table 2: Comparison of various models across different tasks. The evaluated tasks include WMT22 (EN→DE) translation and GYAFC Family & Relationships (F&R) and GYAFC Entertainment & Music (E&M) style transfer. The models considered for evaluation are OPT-IML-MAX-1.3B, LLaMA-7B, and OPT-IML-MAX-30B, each with specific adaptations and configurations.

the PLM’s beginning and end. We use four adapter layers for training and inference.

We use widely used evaluation metrics to assess the performance of PEMA as follows:

Sacre-Bleu (sBLEU) (Post, 2018) is a commonly used metric to calculate the n-gram accuracy between the source and target sentences. It evaluates how well the generated sentence preserves the meaning of the reference and captures target domain distribution. Higher scores are better.

Perplexity (PPL) (Jelinek et al., 1977) is to assess the fluency of generated sentences. We use pre-trained GPT-2 large (Radford et al., 2019) to calculate the exponential of the negative log-likelihood of a current token given the previous context. Lower scores are better.

COMET (Rei et al., 2020) is a neural network-based metric for assessing machine translation quality. It shows a positive correlation with human judgments. We utilize the default, pre-trained COMET model² for the WMT22. Higher scores are better.

Formality Improvement (FormImp) measure formality improvement based on XFORMAL (Briakou et al., 2021a). To measure the formality score of a sentence, we train a BERT-Large (Devlin et al., 2019b) on an external formality dataset consisting of 4K human-annotated examples (Pavlick and Tetreault, 2016). We compute the formality score for each formal reference sentence (FR), informal

input sentence (I), and generated sentence (G). Then, we measure the relative distance using the formula: $\frac{G}{FR-I} \times 100$. We employ this metric for the GYAFC task. Higher scores are better.

5.2.1 Results

For the WMT22 (EN→DE) translation task, we evaluated sBLEU, PPL, and COMET metrics. As Table 2 shows, PEMA outperforms baselines in sBLEU and COMET. Offsite-Tuning, LoRA, and UniPELT perform slightly better than a naive pre-trained language model and PEMA in terms of PPL. However, they require more memory consumption for training than PEMA. Finally, PEMA generates more appropriate translated sentences than other baselines for sBLEU with relatively small memory consumption.

For the GYAFC style transfer task, we evaluated sBLEU, PPL, and Formality Improvement (FormImp) metrics. As Table 2 shows, PEMA consistently achieves favorable performance. PEMA shows the highest sBLEU scores, effectively maintaining meaning preservation across different domains and models. PEMA performs slightly better than a naive pre-trained language model and is comparable to other baselines in terms of FormImp. Furthermore, we observe a trade-off between sBLEU and formality. These findings support previous observations in the same formality style transfer task with multilingual formality (Briakou et al., 2021b).

²<https://github.com/Unbabel/COMET>

WMT22 (EN→DE)	sBLEU	PPL	COMET
OPT-30B	18.22	45.81	77.41
OPT-30B+ B_{pd}	18.74	48.05	77.76
OPT-30B+ B_{pd} + GU	19.17	48.60	78.57
OPT-30B+ B_{pd} + GU + B_{rct} (PEMA)	19.22	46.62	79.21
GYAFC (F&R)	sBLEU	PPL	FormImp
OPT-30B	60.41	<u>20.04</u>	29.33
OPT-30B+ B_{pd}	70.00	20.38	47.38
OPT-30B+ B_{pd} + GU	<u>70.29</u>	16.95	<u>51.24</u>
OPT-30B+ B_{pd} + GU + B_{rct} (PEMA)	70.84	22.04	52.35
GYAFC (E&M)	sBLEU	PPL	FormImp
OPT-30B	57.60	21.97	23.88
OPT-30B+ B_{pd}	64.37	26.76	38.80
OPT-30B+ B_{pd} + GU	<u>64.82</u>	25.62	<u>42.61</u>
OPT-30B+ B_{pd} + GU + B_{rct} (PEMA)	65.43	<u>25.53</u>	44.63

Table 3: Ablation results of PEMA over our proposed approaches. The techniques include a token prediction decoder (B_{pd}), Gradual Unrolling (GU), and a reconstruction decoder (B_{rct}). We use OPT-IML-MAX-30B as a baseline. Implementing all techniques together enhances overall performance.

λ/λ_{max}	0.9	0.8	0.7	0.6	0.5	0.4	0.3	0.2	0.1
With GU	47.45	46.61	46.62	46.18	46.12	46.03	45.85	45.89	45.84
Without GU	54.29	51.87	50.22	49.70	49.45	48.09	47.76	47.67	47.52

Table 4: Impact of Gradual Unrolling (GU) on perplexity across different λ/λ_{max} values. Using GU consistently outperforms the approach without GU for all λ/λ_{max} values, ranging from 0.1 to 0.9.

5.3 Ablation Study

To assess the effectiveness of PEMA, we conduct ablation studies to demonstrate the incremental improvement achieved by incorporating each component of PEMA. We utilize a token prediction decoder (B_{pd}) to predict the target token based on the context representation obtained from the pre-trained language model. As shown in Table 3, the token prediction decoder enhances task performance. Building on this, we incorporated Gradual Unrolling (GU) and the Reconstruction Decoder (B_{rct}) to further improve performance. The inclusion of these three methods yields the highest performance gains, as shown in the results.

Interpolation Parameter (λ_{max}) We propose the Gradual Unrolling (GU) interpolation strategy, where PEMA initially guides the generation of a new task and subsequently leverages the language model for contextual completion of sentences. Table 3 shows the effectiveness of GU in enhancing performance by enabling the language model to provide context completion. We further compare with and without GU by adjusting the λ_{max} hyperparameter in the WMT22 task. As shown in Figure 4, with GU maintains better performance

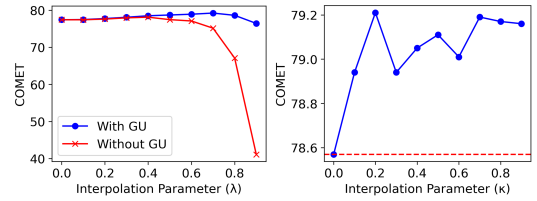


Figure 4: Performance variations on the WMT22 task with interpolation values λ_{max} (left) and κ (right). For λ_{max} , using Gradual Unrolling (GU) prevents performance degradation and enhances results, unlike without GU , where performance drops sharply. With κ when λ_{max} is set at 0.7, combining reconstruction loss with next-token prediction loss improves performance over excluding reconstruction loss (red dotted line), as indicated by better results when κ is above zero.

stability at higher λ_{max} values while achieving noticeable performance improvement over without GU . We also report details on the impact of incorporating λ_{max} in Figure 5 in the appendix. Additionally, we conduct an experiment to measure perplexity. Table 4 shows that GU consistently outperforms across λ/λ_{max} values from 0.1 to 0.9.

Interpolation Parameter (κ) We investigate the effectiveness of the reconstruction decoder, which reconstructs the original vector $f(c_i)$. Table 3 and Figure 4 demonstrate that incorporating the reconstruction decoder improves performance across desired tasks, demonstrating its efficacy in enhancing generation quality. We also report details on the impact of incorporating κ in Figure 6 in the appendix.

6 Conclusion

In this paper, we present PEMA, a novel parameter-efficient fine-tuning approach for language models. Unlike existing PEFT methods, PEMA utilizes minimal pre-trained model parameters during training, making it an efficient and adaptable method for offsite-tuning. PEMA includes a token prediction decoder, Gradual Unrolling, and a reconstruction decoder to improve model performance. Our comprehensive evaluations on translation and style transfer tasks demonstrate PEMA's effectiveness in generating text that more closely follows target domain distributions. Additionally, PEMA proves its computational efficiency by utilizing minimal training memory and achieving faster training latency with a syntactic dataset. Overall, PEMA offers efficient fine-tuning and presents a promising direction for an offsite-tunable PEFT approach in downstream NLP tasks.

568 Limitations

569 PEMA introduces a novel Parameter-Efficient
570 Fine-Tuning (PEFT) method for privacy-preserving
571 offsite-tuning. However, this process requires data
572 owners to share predicted next-tokens with PLM
573 owners during inference, which raises potential pri-
574 vacy concerns. These concerns necessitate further
575 investigation of effective mitigation strategies.

576 Additionally, sharing the W_{hd} weight between
577 PLM owners and data owners poses challenges re-
578 lated to model privacy. In our experiments, we used
579 open-source PLMs due to the confidentiality issues
580 associated with proprietary PLMs. Our future work
581 will explore enabling data owners to generate a
582 new Language Model (LM) head using a shared
583 tokenizer from the PLM owner, enhancing privacy
584 between the PLM and the data owner.

585 Finally, through PEMA, data and PLM owners
586 can fine-tune efficiently and effectively with mini-
587 mal communication. However, the way data own-
588 ers use PEMA could unintentionally lead to data
589 leakage issues. Subsequent research will explore
590 solutions to address this challenge.

591 While our research has been focused on machine
592 translation tasks, it can be applied to various NLP
593 tasks depending on the initial input. Consequently,
594 future studies will investigate the application of our
595 method across a range of NLP tasks.

596 Ethics Statement

597 The results of our research are based on existing
598 studies, and all generation models and datasets used
599 are publicly available and used for their intended
600 use with no ethical concerns.

601 References

602 Rohan Anil, Andrew M Dai, Orhan Firat, Melvin John-
603 son, Dmitry Lepikhin, Alexandre Passos, Siamak
604 Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng
605 Chen, et al. 2023. Palm 2 technical report. *arXiv*
606 *preprint arXiv:2305.10403*.

607 AnthropicAI. 2023. Introducing claude. [https://www.
608 anthropic.com/index/introducing-claude](https://www.anthropic.com/index/introducing-claude).
609 Accessed: 2023-08-15.

610 Eleftheria Briakou, Di Lu, Ke Zhang, and Joel Tetreault.
611 2021a. *Olá, bonjour, salve! XFORMAL: A bench-
612 mark for multilingual formality style transfer*. In
613 *Proceedings of the 2021 Conference of the North
614 American Chapter of the Association for Computa-
615 tional Linguistics: Human Language Technologies*,
616 pages 3199–3216, Online. Association for Computa-
617 tional Linguistics.

618 Eleftheria Briakou, Di Lu, Ke Zhang, and Joel Tetreault.
619 2021b. *Olá, bonjour, salve! XFORMAL: A bench-
620 mark for multilingual formality style transfer*. In
621 *Proceedings of the 2021 Conference of the North
622 American Chapter of the Association for Computa-
623 tional Linguistics: Human Language Technologies*,
624 pages 3199–3216, Online. Association for Computa-
625 tional Linguistics.

626 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
627 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
628 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
629 Askell, et al. 2020. Language models are few-shot
630 learners. *Advances in neural information processing
631 systems*, 33:1877–1901.

632 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,
633 Maarten Bosma, Gaurav Mishra, Adam Roberts,
634 Paul Barham, Hyung Won Chung, Charles Sutton,
635 Sebastian Gehrmann, et al. 2022. Palm: Scaling
636 language modeling with pathways. *arXiv preprint
637 arXiv:2204.02311*.

638 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
639 Kristina Toutanova. 2019a. *BERT: Pre-training of
640 deep bidirectional transformers for language under-
641 standing*. In *Proceedings of the 2019 Conference of
642 the North American Chapter of the Association for
643 Computational Linguistics: Human Language Tech-
644 nologies, Volume 1 (Long and Short Papers)*, pages
645 4171–4186, Minneapolis, Minnesota. Association for
646 Computational Linguistics.

647 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
648 Kristina Toutanova. 2019b. *BERT: Pre-training of
649 deep bidirectional transformers for language under-
650 standing*. In *Proceedings of the 2019 Conference of
651 the North American Chapter of the Association for
652 Computational Linguistics: Human Language Tech-
653 nologies, Volume 1 (Long and Short Papers)*, pages
654 4171–4186, Minneapolis, Minnesota. Association for
655 Computational Linguistics.

656 Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zong-
657 han Yang, Yusheng Su, Shengding Hu, Yulin Chen,
658 Chi-Min Chan, Weize Chen, et al. 2022. Delta tuning:
659 A comprehensive study of parameter efficient meth-
660 ods for pre-trained language models. *arXiv preprint
661 arXiv:2203.06904*.

662 Justin Guinney and Julio Saez-Rodriguez. 2018. Al-
663 ternative models for sharing confidential biomedical
664 data. *Nature biotechnology*, 36(5):391–392.

665 Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-
666 Kirkpatrick, and Graham Neubig. 2021. Towards a
667 unified view of parameter-efficient transfer learning.
668 In *International Conference on Learning Representa-
669 tions*.

670 Amr Hendy, Mohamed Abdelrehim, Amr Sharaf,
671 Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita,
672 Young Jin Kim, Mohamed Afify, and Hany Hassan
673 Awadalla. 2023. How good are gpt models at ma-
674 chine translation? a comprehensive evaluation. *arXiv
675 preprint arXiv:2302.09210*.

676	Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski,	Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Lam Tam,	734
677	Bruna Morrone, Quentin De Laroussilhe, Andrea	Zhengxiao Du, Zhilin Yang, and Jie Tang. 2021. P-	735
678	Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019.	tuning v2: Prompt tuning can be comparable to fine-	736
679	Parameter-efficient transfer learning for NLP . In <i>Pro-</i>	tuning universally across scales and tasks. <i>arXiv</i>	737
680	<i>ceedings of the 36th International Conference on</i>	<i>preprint arXiv:2110.07602</i> .	738
681	<i>Machine Learning</i> , volume 97 of <i>Proceedings of Ma-</i>		
682	<i>chine Learning Research</i> , pages 2790–2799. PMLR.		
683	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan	Yuning Mao, Lambert Mathias, Rui Hou, Amjad Alma-	739
684	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and	hairi, Hao Ma, Jiawei Han, Scott Yih, and Madian	740
685	Weizhu Chen. 2022. LoRA: Low-rank adaptation of	Khabsa. 2022. Unipelt: A unified framework for	741
686	large language models . In <i>International Conference</i>	parameter-efficient language model tuning. In <i>Pro-</i>	742
687	<i>on Learning Representations</i> .	<i>ceedings of the 60th Annual Meeting of the Associa-</i>	743
688		<i>tion for Computational Linguistics (Volume 1: Long</i>	744
689	Changwu Huang, Yuanxiang Li, and Xin Yao. 2020. A	<i>Papers)</i> , pages 6253–6264.	745
690	survey of automatic parameter tuning methods for	OpenAI. 2022. Chatgpt: Optimizing language mod-	746
691	metaheuristics . <i>IEEE Transactions on Evolutionary</i>	els for dialogue. https://online-chatgpt.com/ .	747
692	<i>Computation</i> , 24(2):201–216.	Accessed: 2023-08-15.	748
693	Srinivasan Iyer, Xi Victoria Lin, Ramakanth Pasunuru,	OpenAI. 2023a. Fine-tuning - openai api.	749
694	Todor Mihaylov, Daniel Simig, Ping Yu, Kurt Shus-	https://platform.openai.com/docs/guides/	750
695	ter, Tianlu Wang, Qing Liu, Punit Singh Koura, et al.	fine-tuning . Accessed: 2023-08-15.	751
696	2022. Opt-impl: Scaling language model instruc-	OpenAI. 2023b. Gpt-4 technical report .	752
697	tion meta learning through the lens of generalization.		
698	<i>arXiv preprint arXiv:2212.12017</i> .	Adam Paszke, Sam Gross, Francisco Massa, Adam	753
699	Fred Jelinek, Robert L Mercer, Lalit R Bahl, and	Lerer, James Bradbury, Gregory Chanan, Trevor	754
700	James K Baker. 1977. Perplexity—a measure of the	Killeen, Zeming Lin, Natalia Gimelshein, Luca	755
701	difficulty of speech recognition tasks. <i>The Journal of</i>	Antiga, Alban Desmaison, Andreas Kopf, Edward	756
702	<i>the Acoustical Society of America</i> , 62(S1):S63–S63.	Yang, Zachary DeVito, Martin Raison, Alykhan Te-	757
703	Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke	jani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang,	758
704	Zettlemoyer, and Mike Lewis. 2020. Generalization	Junjie Bai, and Soumith Chintala. 2019. Pytorch:	759
705	through Memorization: Nearest Neighbor Language	An imperative style, high-performance deep learning	760
706	Models. In <i>International Conference on Learning</i>	library . In <i>Advances in Neural Information Process-</i>	761
707	<i>Representations (ICLR)</i> .	<i>ing Systems 32</i> , pages 8024–8035. Curran Associates,	762
708	Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton	Inc.	763
709	Dvorkovich, Christian Federmann, Mark Fishel,	Ellie Pavlick and Joel Tetreault. 2016. An Empiri-	764
710	Thamme Gowda, Yvette Graham, Roman Grund-	cal Analysis of Formality in Online Communication .	765
711	kiewicz, Barry Haddow, Rebecca Knowles, Philipp	<i>Transactions of the Association for Computational</i>	766
712	Koehn, Christof Monz, Makoto Morishita, Masaaki	<i>Linguistics</i> , 4:61–74.	767
713	Nagata, Toshiaki Nakazawa, Michal Novák, Martin	Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé,	768
714	Popel, and Maja Popović. 2022. Findings of the 2022	Kyunghyun Cho, and Iryna Gurevych. 2021.	769
715	conference on machine translation (WMT22) . In	AdapterFusion: Non-destructive task composition for	770
716	<i>Proceedings of the Seventh Conference on Machine</i>	transfer learning . In <i>Proceedings of the 16th Con-</i>	771
717	<i>Translation (WMT)</i> , pages 1–45, Abu Dhabi, United	<i>ference of the European Chapter of the Association</i>	772
718	Arab Emirates (Hybrid). Association for Computa-	<i>for Computational Linguistics: Main Volume</i> , pages	773
719	tional Linguistics.	487–503, Online. Association for Computational Lin-	774
720	Brian Lester, Rami Al-Rfou, and Noah Constant. 2021.	guistics.	775
721	The power of scale for parameter-efficient prompt	Sundar Pichai. 2023. An important next step on	776
722	tuning . In <i>Proceedings of the 2021 Conference on</i>	our ai journey. https://blog.google/intl/	777
723	<i>Empirical Methods in Natural Language Processing</i> ,	en-africa/products/explore-get-answers/	778
724	pages 3045–3059, Online and Punta Cana, Domini-	an-important-next-step-on-our-ai-journey/ .	779
725	can Republic. Association for Computational Lin-	Accessed: 2023-08-15.	780
726	guistics.	Reiner Pope, Sholto Douglas, Aakanksha Chowdhery,	781
727	Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning:	Jacob Devlin, James Bradbury, Jonathan Heek, Kefan	782
728	Optimizing continuous prompts for generation . In	Xiao, Shivani Agrawal, and Jeff Dean. 2023. Effi-	783
729	<i>Proceedings of the 59th Annual Meeting of the Asso-</i>	ciently scaling transformer inference. <i>Proceedings</i>	784
730	<i>ciation for Computational Linguistics and the 11th</i>	<i>of Machine Learning and Systems</i> , 5.	785
731	<i>International Joint Conference on Natural Language</i>		
732	<i>Processing (Volume 1: Long Papers)</i> , pages 4582–	Matt Post. 2018. A call for clarity in reporting bleu	786
733	4597, Online. Association for Computational Lin-	scores. <i>WMT 2018</i> , page 186.	787

788 Sharpened Productions. 2023. Slang.net: The slang
789 dictionary. <https://slang.net/>. Accessed: 2023-
790 08-14.

791 Xipeng Qiu, Tianxiang Sun, Yige Xu, Yunfan Shao,
792 Ning Dai, and Xuanjing Huang. 2020. Pre-trained
793 models for natural language processing: A survey.
794 *Science China Technological Sciences*, 63(10):1872–
795 1897.

796 Alec Radford, Jeff Wu, Rewon Child, David Luan,
797 Dario Amodei, and Ilya Sutskever. 2019. Language
798 models are unsupervised multitask learners.

799 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine
800 Lee, Sharan Narang, Michael Matena, Yanqi Zhou,
801 Wei Li, and Peter J Liu. 2020. Exploring the limits
802 of transfer learning with a unified text-to-text trans-
803 former. *The Journal of Machine Learning Research*,
804 21(1):5485–5551.

805 Sudha Rao and Joel Tetreault. 2018. [Dear sir or madam,](#)
806 [may I introduce the GYAFC dataset: Corpus, bench-](#)
807 [marks and metrics for formality style transfer.](#) In
808 *Proceedings of the 2018 Conference of the North*
809 *American Chapter of the Association for Computa-*
810 *tional Linguistics: Human Language Technologies,*
811 *Volume 1 (Long Papers)*, pages 129–140, New Or-
812 leans, Louisiana. Association for Computational Lin-
813 guistics.

814 Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon
815 Lavie. 2020. [COMET: A neural framework for MT](#)
816 [evaluation.](#) In *Proceedings of the 2020 Conference*
817 *on Empirical Methods in Natural Language Process-*
818 *ing (EMNLP)*, pages 2685–2702, Online. Association
819 for Computational Linguistics.

820 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier
821 Martinet, Marie-Anne Lachaux, Timothée Lacroix,
822 Baptiste Rozière, Naman Goyal, Eric Hambro,
823 Faisal Azhar, et al. 2023. Llama: Open and effi-
824 cient foundation language models. *arXiv preprint*
825 *arXiv:2302.13971*.

826 Ahmet Üstün and Asa Cooper Stickland. 2022. [When](#)
827 [does parameter-efficient transfer learning work for](#)
828 [machine translation?](#) In *Proceedings of the 2022 Con-*
829 *ference on Empirical Methods in Natural Language*
830 *Processing*, pages 7919–7933, Abu Dhabi, United
831 Arab Emirates. Association for Computational Lin-
832 guistics.

833 Guangxuan Xiao, Ji Lin, and Song Han. 2023. Offsite-
834 tuning: Transfer learning without full model. *arXiv*.

835 Yahoo. 2007. L6 - yahoo! answers comprehensive ques-
836 tions and answers version 1.0. [https://webscope.](https://webscope.sandbox.yahoo.com/)
837 [sandbox.yahoo.com/](https://webscope.sandbox.yahoo.com/). Accessed: 2023-07-02.

838 Qingru Zhang, Minshuo Chen, Alexander Bukharin,
839 Pengcheng He, Yu Cheng, Weizhu Chen, and
840 Tuo Zhao. 2023. [Adaptive budget allocation for](#)
841 [parameter-efficient fine-tuning.](#) In *The Eleventh In-*
842 *ternational Conference on Learning Representations*.

Susan Zhang, Stephen Roller, Naman Goyal, Mikel
Artetxe, Moya Chen, Shuohui Chen, Christopher De-
wan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022.
Opt: Open pre-trained transformer language models.
arXiv preprint arXiv:2205.01068.

A Performance on Different Rank Sizes

Model	WMT22 (EN→DE)	GYAFC (F&R)	GYAFC (E&M)
OPT-1.3B (LoRA _{r=8})	3.25	23.13	18.41
OPT-1.3B (LoRA _{r=512})	4.28	20.98	15.57
OPT-1.3B (PEMA _{r=8})	<u>11.75</u>	<u>56.29</u>	<u>54.22</u>
OPT-1.3B (PEMA _{r=512})	12.87	64.82	61.24
LLaMA-7B (LoRA _{r=8})	10.92	14.80	12.69
LLaMA-7B (LoRA _{r=512})	<u>11.46</u>	<u>52.67</u>	<u>52.15</u>
LLaMA-7B (PEMA _{r=8})	3.88	48.88	45.73
LLaMA-7B (PEMA _{r=512})	14.50	63.99	60.88
OPT-30B (LoRA _{r=8})	16.05	61.28	59.48
OPT-30B (LoRA _{r=512})	16.03	61.39	59.76
OPT-30B (PEMA _{r=8})	<u>18.33</u>	<u>62.87</u>	<u>60.12</u>
OPT-30B (PEMA _{r=512})	19.22	70.84	65.43

Table 5: Experiment on LoRA and PEMA on meaning preservation (sBLEU) across rank variations ($r = \{8, 512\}$). The result shows PEMA consistently outperforms LoRA on sBLEU and COMET.

LoRA (Hu et al., 2022) states performance remains comparable with a small rank size. However, AdaLoRA (Zhang et al., 2023) finds a large rank size in the last layer of PLMs is needed for better performance. Performance evaluation on PEMA and baseline PEFT methods is conducted at the last layer of PLMs. For this reason, we set $r = 512$ for LoRA and PEMA to minimize the effect on performance with rank size. However, LoRA uses a rank size between 1 to 64 for their experiment. As PEMA is a LoRA-based PEFT method, we compared the performance on meaning preservation using the rank size employed in LoRA (8) and the rank size used in our experiment (512). As Table 5 shows, a larger rank size generally achieves favorable performance. In the case of LoRA, using a rank size of 512 outperforms 8 in 6 out of 9 cases. PEMA with a rank size of 512 performs better than PEMA with a rank size of 8 at all tasks.

B Measuring Informal Language Patterns

The GYAFC dataset for style transfer includes common informal input patterns that are frequently occur. To analyze the amount of mitigation, we categorize these patterns into four types. The four informal patterns are as follows. **Slang abbreviations**

	Informal Input	Formal Reference	Naive OPT-30B	kNN-LM	LoRA	UniPELT	Offsite-Tuning	PEMA
Family & Relationships								
Slang Abbreviation	525	307.75	346	339	356	322	361	289
All Capital	68	0	61	60	8	5	65	3
Redundant Word	39	2	1	1	2	0	17	3
Non-Capital Start	636	1.5	16	2	1	1	2	0
Entertainment & Music								
Slang Abbreviation	651	485.75	541	538	530	534	529	463
All Capital	36	0	31	34	9	9	37	0
Redundant Word	49	17.75	5	5	7	3	16	32
Non-Capital Start	655	7	24	2	0	1	3	0

Table 6: Count of informal patterns for each generated formal sentence. The result shows that PEMA performs better in mitigating informal patterns than baseline approaches. Lower is better.

are informal short forms of words or phrases (e.g., "LOL"- "laughing out loud"). To identify the presence of slang words, we check how many words from the predicted target sentence are present in the slang dictionary from (Productions, 2023). **All capital** is a pattern in which all characters in a generated word are capitalized (e.g., "FUNNY"). We calculate how many generated words are all capitalized. **Redundant word** occurs when two consecutive words are the same. For example, "I lie lie lie and then I lie some more." has two redundant words. **Non-capital start** is counted when a sentence does not start with a capital letter (e.g., "i only want points").

Table 6 shows the count of each informal pattern in generated sentences for both the baseline and PEMA. We also show an informal pattern count on informal input and formal reference. There are four reference sentences for each example in the test set. We show the average count for each pattern using the formal reference. It shows PEMA is good at mitigating slang abbreviation, all capital, and non-capital start compared to other baseline approaches. Interestingly, PEMA outperforms formal references in mitigating slang abbreviations and non-capital start.

C Dataset

C.1 Data Statistic

Table 7 shows data statistics of GYAFC and WMT22. For WMT22, we use a news-commentary v16 (EN→DE) for training. The test set for GYAFC has four references, while WMT22 has one reference for each test input.

Dataset	Train	Valid	Test	Length of \mathcal{E}
GYAFC (F&R)	51,967	2,788	1,332	691,531
GYAFC (E&M)	52,595	2,877	1,416	695,465
WMT22	388,482	2,203	1,984	20,983,482

Table 7: Data statistic of GYAFC and WMT22 with length of external memory \mathcal{E} .

Task	Example
WMT22	English: In better shape, but not alone. German: In besserer Verfassung, aber nicht allein.
GYAFC	Informal: I'd say it is punk though. Formal: However, I do believe it to be punk.

Table 8: Example of parallel dataset GYAFC and WMT22.

C.2 Dataset Examples

Table 8 demonstrates examples of parallel datasets of GYAFC and WMT22.

C.3 Prompts

Table 9 presents prompt input used for evaluation. WMT22 and GYAFC have two placeholders. This includes [English Input] and [Informal Input]. [Generated Output] is a predicted output sentence generated by PLMs.

[English Input] represents the English input sentence in WMT22. [Informal Input] is the informal input sentence in GYAFC. An example of the parallel data input can be found in Table 8.

C.4 Post-processing

We use three decoder-based pre-trained language models for evaluation: OPT-IML-MAX-1.3B, LLaMA-7B, and OPT-IML-MAX-30B. These

Task	Prompt
WMT22	Translate this from English to German: [English Input] German: [Generated Output]
GYAFC	Convert the following informal sentence into a formal sentence: Informal: [Informal Input] Formal: [Generated Output]

Table 9: Prompt used for evaluation. [] represents the placeholder.

Model	Common hallucination patterns
OPT	I'm not sure ... I 50% ... Convert the following informal sentence ... Translate this from English to German: ... I
LLaMA	Informal: ... ### ... Comment: ... \\ ... \\begin ... Answer: ...

Table 10: Common hallucination patterns after generating a predicted sentence.

models are capable of generating tokens continuously. This characteristic makes decoder-based language models generate beyond the predicted sentences, typically called hallucinations. We find common hallucination patterns in each pre-trained language model. We post-process hallucinations generated after the predicted sentence for evaluation. Table 10 shows common hallucination patterns that are removed.

D Implementation Details

We use three RTX 8000 GPUs with 48GB GDDR6 memory for our experiment. For OPT-IML-MAX-1.3B, we use full precision (FP32) for training and inference. For LLaMA-7B and OPT-IML-MAX-30B, we use half-precision (FP16) and distribute the model across three GPUs using the HuggingFace Accelerate library. The hyperparameters for PEMA and the baselines are in Table 11. The best hyperparameter is selected using a grid search.

E Examples of Generated Outputs

The generated formal outputs of GYAFC are shown in Table 13 and Table 12. In WMT22, the German output generated is presented in Table 14. It shows PEMA understands the meaning of abbreviated

PEMA	
Random seed	123
Batch size	40,960
Adam lr	1e-03
Adam (β_1, β_2)	(0.9, 0.999)
Adam eps	1e-08
Number of rank	512
Optimal λ_{max}	0.7 to 0.9
Offsite-Tuning	
Random seed	42
Batch size	18
Emulator size	$\frac{1}{3}$ of PLM
Adam lr	1e-04
Adam (β_1, β_2)	(0.9, 0.999)
Adam eps	1e-08
LoRA	
Random seed	123
Batch size	10 to 30
Adam lr	1e-03
Adam (β_1, β_2)	(0.9, 0.999)
Adam eps	1e-08
Number of rank	512
LoRA α	1
Merge weight	FALSE
k NN-LM	
Random seed	1
Number of centroids learn	4,096
Quantized vector size	64
Number of clusters to query	32
Distance function	L2 Distance
UniPELT	
Random seed	123
Batch size	10 to 30
Adam lr	1e-03
Adam (β_1, β_2)	(0.9, 0.999)
Adam eps	1e-08
Prefix gate	True
Prefix length	10
Prefix mid dimension	512
LoRA gate	True
Number of rank	10
LoRA α	16
Adapter gate	True
Adapter down sample	$D_{hid}/2$ Adapter
Used PEFT methods	Prefix tuning LoRA

Table 11: Hyper-parameter setup of each baseline method. We select the batch size between 10 to 30. D_{hid} represent hidden size of a model.

format (e.g., translating "5'4" to "5 feet 4 inches"), or removing the informal word (e.g., "flirt" which typically refers to playful or teasing behavior). Mit-

948
949
950

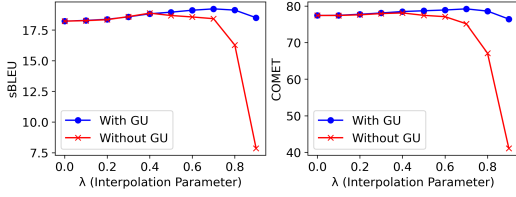


Figure 5: Performance variation for each interpolation value λ_{max} in the WMT22 task. With both Gradual Unrolling (*GU*) (blue) and without *GU* (red), there is a decline in performance at a specific point of λ_{max} . However, when utilizing *GU*, the model is not only robust to performance degradation but also gains performance improvement.

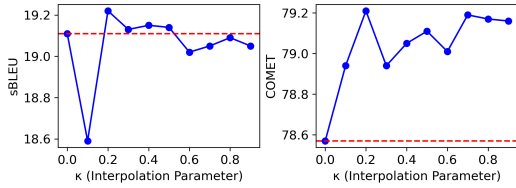


Figure 6: Impact of mixing ratio values between reconstruction loss and predicting the next-token loss in the WMT22 task. When κ is 0, it means excluding reconstruction loss (red dashed line). We fix the λ_{max} value as 0.7. The graphs show that combining reconstruction loss and predicting the next-token loss is superior to excluding reconstruction loss.

igating common informal patterns such as all capital words (e.g., "PINK FLOYD" to "Pink Floyd") while preserving the meaning of input (e.g., "Wir" means "We" in German).

F Difference Between PEMA and LoRA at W_{hd}

Applying LoRA to $W_{hd} \in \mathbb{R}^{v \times d}$, a larger set of parameters is required due to the difference in input and output sizes (d and v). Conversely, PEMA operates more efficiently, utilizing computation resources by receiving an input of size d and yielding an output of the same size. For instance, OPT-1.3B has $d = 2,048$ and $v = 50,272$.

G Impact on Interpolation λ and κ

In the WMT22 task, we observe performance variation with different interpolation values, λ_{max} in Figure 5. Additionally, we investigate the impact of the mixing ratio values between reconstruction loss and predicting the next-token loss in Figure 6.

H Licensing Information

Models OPT is licensed under the MIT License. The LLaMA is licensed under the GNU General Public License (GPL) version 3.

Fine-tuning Methods k NN-LM, LoRA, and Offsite-Tuning are licensed under the MIT License. UniPELT is licensed under the Creative Commons Attribution-NonCommercial (CC-BY-NC) license. **Dataset** GYAFC is based on the Yahoo Answers corpus (L6 - Yahoo! Answers Comprehensive Questions and Answers version 1.0) (Yahoo, 2007), and is designated for research purposes. Access to the GYAFC dataset requires prior access to Yahoo Answers corpus. WMT22 is freely available for research purposes for academic and educational activities.

Input		he is probably wondering if your interested in him at all....flirt back!!	sBLEU
Reference	1	He is likely wondering if you are interested in him at all; Flirt back with him.	
	2	He probably wants to know if you're interested in him.	
	3	He is probably wondering if you are interested in him at all, so flirt back.	
	4	He is probably wondering if you are interested in him at all. Flirt back.	
Output	PEMA	He is probably wondering if you are interested in him at all.	100.0
	LoRA	He is probably wondering if you are interested in him at all. If you are interested, flirt back.	66.78
	kNN-LM	It is most likely that he is wondering if you are interested in him at all....flirt back!!	42.60
	UniPELT	He is probably wondering if your interested in him at all....flirt back!	50.82
	Offsite-Tuning	He probably is wondering if you are interested in him at all. Flirt back!!	72.98
	Naïve OPT-30B	In informal situations he is probably wondering if your interested in him at all.	46.03
Input		I don't know!...I just want the points...lol	
Reference	1	I only want points.	
	2	I do not know. I merely want the points.	
	3	I do not know; I just want the points.	
	4	I do not know, I only want the points.	
Output	PEMA	I do not know, but I just want the points.	73.49
	LoRA	I don't know!... I just want the points. I am not sure what I am doing.	25.31
	kNN-LM	I don't know!...I just want the points...lol	34.90
	UniPELT	I don't know!...I just want the points...lol	34.90
	Offsite-Tuning	-	0.00
	Naïve OPT-30B	I don't know!...I just want the points...lol	34.90
Input		No way im 5'4 and he's 6'2	
Reference	1	No, I am 5ft 4inches and he is 6ft and 2inches.	
	2	No way, I am only 5'4" and he is 6'2".	
	3	Not at all. I am five feet four inches tall and he is 6 feet 2 inches tall. '	
	4	No chance, I am five feet four inches tall and he is six feet two inches tall.	
Output	PEMA	No way, I am 5 feet 4 inches tall and he is 6 feet 2 inches tall.	74.44
	LoRA	No way, I am 5'4 and he is 6'2.	51.52
	kNN-LM	No way, I am 5'4 and he is 6'2	50.05
	UniPELT	No way, I am 5'4 and he is 6'2	50.05
	Offsite-Tuning	No way im 5'4 and he's 6'2.	7.78
	Naïve OPT-30B	No way, I am 5'4 and he is 6'2	45.72

Table 12: Examples of generated formal output of GYAFC (Family & Relationships) for given informal input. One interesting example is PEMA can understand the meaning of abbreviated height descriptions like "5'4" and "6'2". And rewrite them into more formal forms "5 feet 4 inches" and "6 feet 2 inches".

Input		I respect that tyra banks is trying show hosting..but the big O [aka OPRAH] won...	sBLEU
Reference	1	While I respect that Tyra Banks is branching out into television hosting, Oprah Winfrey is the better example.	
	2	I respect the hostess Tyra Banks, but I prefer Oprah.	
	3	I admire Tyra Banks for trying to become a tv show host. However, Oprah is clearly the best hostess.	
	4	I respect that Tyra Banks is trying out show hosting, but I think that Oprah won.	
Output	PEMA	I respect that Tyra Banks is trying to show hosting..but the big O [aka Oprah Winfrey] won.	35.56
	LoRA	I respect that tyra banks is trying show hosting, but the big O won.	24.32
	kNN-LM	I respect that tyra banks is trying show hosting..but the big O [aka OPRAH] won.	12.02
	UniPELT	I respect that tyra banks is trying show hosting..but the big O won...	13.40
	Offsite-Tuning	I respect that tyra banks is trying show hosting..but the big O [aka OPRAH] won...	10.90
	Naive OPT-30B	In informal: I respect that tyra banks is trying show hosting..but the big O [aka OPRAH] won...	9.56
Input		Not sure what's gonna happen, sure it will be good and OH YEAH!!!!!!	
Reference	1	I am not sure what is going to happen but I am sure it will be good.	
	2	I'm not sure what is going to happen, I'm sure it will be good.	
	3	I am not sure what will happen, but I am sure it will be good. Oh, yes.	
	4	I am not sure what will happen, however, I am sure that it will be good.	
Output	PEMA	I am not sure what is going to happen, but I am sure it will be good and OH YEAH!!!!!!	61.11
	LoRA	Not sure what's going to happen, but I am sure it will be good and OH YEAH!!!!!!	43.05
	kNN-LM	Not sure what's gonna happen, sure it will be good and OH YEAH!!!!!!	20.22
	UniPELT	Not sure what's going to happen, but	16.74
	Offsite-Tuning	Not sure what's gonna happen, sure it will be good and OH YEAH!!!!!!	20.22
	Naive OPT-30B	Not sure what's gonna happen, sure it will be good and OH YEAH!!!!!!	20.22
Input		Journey, PINK FLOYD, The POLICE, The EAGLES & RUSH... omg!	
Reference	1	Journey, Pink Floyd, The Police, The Eagles, and Rush - oh my!	
	2	I like Journey, Pink Floyd, The Police, The Eagles, and Rush.	
	3	Oh goodness, Journey, Pink Floyd, The Police, the Eagles, and Rush!	
	4	Journey, Pink Floyd, The Police, The Eagles, and Rush are all great classic bands.	
Output	PEMA	I love Journey, Pink Floyd, The Police, The Eagles and Rush.	69.01
	LoRA	Journey, PINK FLOYD, The Police, The Eagles & Rush. I love it!	36.45
	kNN-LM	Journey, PINK FLOYD, The Police, The Eagles & Rush... omg!	35.66
	UniPELT	Journey, PINK FLOYD, The Police, The Eagles & Rush... omg!	35.66
	Offsite-Tuning	Journey, Pink Floyd, The Eagles, Rush, and The Police. Oh my god!	47.29
	Naive OPT-30B	Journey, PINK FLOYD, The Police, The Eagles & Rush... omg!	35.66

Table 13: Examples of generated formal output of GYAFC (Entertainment & Music) for given informal input. It shows that PEMA is capable of restoring All Capital patterns to their formal format. For example, PEMA successfully restore "Oprah Winfrey" given "OPRAH" as an input.

Input		In Berlin he trained as an opera singer and in 1986 he founded the Palast Orchester together with his classmates.	sBLEU
Reference		In Berlin studierte er Operngesang, 1986 gründete er zusammen mit Kommilitonen das Palast Orchester.	
Output	PEMA	In Berlin studierte er als Opernsänger und gründete 1986 mit seinen Schülern das Palast Orchester.	33.16
	LoRA	In Berlin wurde er als Opernsänger trainiert und 1986 gründete er mit seinen Schülern das Palast Orchester.	22.79
	kNN-LM	In Berlin schloss er 1986 mit seinen Schülern die Palast-Orchester zusammen.	6.64
	UniPELT	In Berlin trainierte er als Opernsänger und gründete 1986 mit seinen Schülern das Palast Orchester. When	19.35
	Offsite-Tuning	ich war in Berlin als Opernstätte und 1986 gründete er die Palast Orchester gemeinsam mit seinen Kollegen.	11.31
	Naive OPT-30B	In Berlin wurde er als Opernsänger trainiert und 1986 gründete er mit seinen Schülern das Palast Orchester.	22.79
Input		Since then the baritone has established himself with his repertoire of hits and chansons in the style of the Comedian Harmonists.	
Reference		Seitdem hat sich der Bariton mit seinem Repertoire aus Schlagern und Chansons im Stil der Comedian Harmonists etabliert.	
Output	PEMA	Seitdem hat er sich mit seinem Repertoire aus Hits und Chansons im Stil der Comedian Harmonists etabliert.	60.26
	LoRA	Seitdem ist er mit seinem Repertoire aus Hitparaden und Chansons im Stil der Comedian Harmonists aufgestiegen.	46.00
	kNN-LM	Seitdem ist er mit seinem Repertoire aus Hitparaden und Chansons im Stil der Comedian Harmonistslipperlipper dieserenedenaut festgelegt.	39.81
	UniPELT	Seitdem ist der Bariton mit seinem Repertoire aus Hitparaden und Chansons im Stil der Comedian Harmonists auf dem Gebiet der Kulturlandschaft in der Schweiz etabliert.	42.85
	Offsite-Tuning	ich habe seitdem die Baritone seine Rekorde aus der Stil der Comedian Harmonisten errichtet.	8.91
	Naive OPT-30B	Seitdem ist er mit seinem Repertoire aus Hitparaden und Chansons im Stil der Comedian Harmonists auf dem Gebiet der Baritonen geschrieben.	39.81
Input		We have the book Terra, but I am generally not very satisfied with the assignments in the book.	
Reference		Wir haben das Buch Terra und mit den meisten Aufgabenstellungen im Buch bin ich eher unzufrieden.	
Output	PEMA	Wir haben das Buch Terra, aber ich bin im Allgemeinen nicht sehr zufrieden mit den Aufgaben in dem Buch.	22.37
	LoRA	ich habe das Buch Terra, aber ich bin im Allgemeinen nicht sehr zufrieden mit den Aufgaben in dem Buch.	10.11
	kNN-LM	ich habe das Buch Terra, aber ich bin im Allgemeinen nicht sehr zufrieden mit denenteilen in dem Buch.	9.38
	UniPELT	ich habe das Buch Terra, aber in der Regel bin ich nicht sehr zufrieden mit den Aufgaben in dem Buch.	10.06
	Offsite-Tuning	ich habe die Buch Terra, aber ich bin allgemein nicht sehr begeistert mit den Schreibungen in der Buch.	6.44
	Naive OPT-30B	ich habe das Buch Terra, aber ich bin im Allgemeinen nicht sehr zufrieden mit den Aufgaben in dem Buch.	10.11

Table 14: Examples of generated German output in WMT22 test set. The result shows that PEMA is capable of generating German output that preserves its meaning.