RIFF: Rephrasing Inputs for Few-shot Fine-tuning of Language Models

Anonymous ACL submission

Abstract

Pre-trained Language Models (PLMs) have shown remarkable performance when finetuned for downstream text processing tasks. Recently, researchers have effectively guided PLMs to perform specific tasks by optimizing input prompts (prompt optimization) or adjusting a small number of model parameters (efficient tuning). In this study, we explore the impact of altering the input text of the original task while fine-tuning language models with these recent prompt optimization and efficient tuning methods. To most effectively rewrite the input text, we apply a learning objective based on Maximum-Marginal Likelihood estimation in a few-shot setting. Experimenting with seven few-shot text classification datasets, we show that enriching training examples with the input text's paraphrases at train and test time significantly enhances the performance of recent prompt optimization and efficient tuning techniques.

1 Introduction

001

007

017

018

021

024

037

Multiple Pre-trained Language Models (PLMs), such as BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), T5 (Raffel et al., 2019), and GPT2 (Radford et al., 2019), have demonstrated remarkable performance when fine-tuned for downstream text processing tasks. PLM variants with less than 1 billion parameters are easier to train endto-end with commodity hardware. However, very recent PLMs have been trained with a few hundred billion parameters, including PaLM-2 (540B)(Anil et al., 2023), GPT3 (175B)(Brown et al., 2020a), OPT (175B)(Zhang et al., 2022a), or Llama-2 (70B)(Touvron et al., 2023). Training all parameters of these models end-to-end is not straightforward unless done with a dedicated cluster of specialized hardware.

In response, NLP research have developed effective techniques to control or alter the behavior of PLMs by updating the input context through prompt optimization (Liu et al., 2021a) or adapting a few additional parameters within the network itself (Hu et al., 2021). However, current PLM control techniques have not considered altering the *original input text* to improve the performance of the model. Here, we investigate this idea by training a secondary smaller PLM to paraphrase the original input at train and test time, thus augmenting the existing data and improving model performance. 041

042

043

044

045

047

049

052

053

055

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

081

Our inspiration comes from interactions with young children. Determining what a child knows is challenging because they can be sensitive to the wording of the question (Donaldson, 1978). Adults are also influenced by different wordings of a question. For example, opinion polling has been found to be sensitive to the wording of questions (Broughton, 1995). Just like we rephrase questions for humans, we should consider rephrasing input text while querying a PLM. For instance, while classifying the topic of a sentence, phrases related to time may be irrelevant and could be removed to simplify the modeling problem. Slight changes to wording could result in the model producing a correct prediction.

We explore the integration of paraphrased input texts during both the training and testing phases. At *training time*, augmenting data through paraphrase generation has been shown to enhance performance (Wei and Zou, 2019; Feng et al., 2021; Chen et al., 2021). We broaden the scope of previous investigations by using paraphrase augmentation in tandem with recent prompt optimization and efficient tuning methods. At *test time*, recent works have used ensemble predictions with various optimized prompts and tuned weights (Izmailov et al., 2019; Allingham et al., 2023; Li et al., 2023). We further contribute to this line of work by incorporating ensemble predictions based on input paraphrases, again in concert with prompt optimization and efficient tuning techniques.

086

091

094

098

100

101

102

103

104

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

127

We start by pre-training a smaller language model on noisy paraphrases generated by a large language model (i.e., ChatGPT). Subsequently, we explore various training objectives for fine-tuning this paraphrase generator with feedback from the main task's language model. Our experiments on seven text classification datasets demonstrate that incorporating paraphrase augmentation during both training and testing phases significantly enhances the performance of discrete/soft prompt optimization and efficient tuning techniques. In summary, our contributions are as follows:

- We propose an efficient idea for Rephrasing Inputs for Few-shot Fine-tuning of language models (RIFF) with prompt optimization and efficient tuning methods.
- We conduct a comprehensive study on various learning objectives for fine-tuning a paraphrase generator using feedback from the main language model.

• Our augmentation experiments on seven text classification datasets reveal that paraphrase generation, when combined with prompt optimization and adaptation techniques, is a simple yet effective approach for tuning language models on commodity hardware, closing the gap with fine-tuning all the parameters.

2 **Problem Formulation**

We focus on classification problems in Natural Language Understanding (NLU) tasks where we have access to a mini-batch of supervised training examples $B_{\text{supp}} = \{(x_i, y_i)\}_{i=1}^N$. Our goal is to update the parameter set θ_{lm} for a language model by maximizing the probability of the class label y_i given the input x_i : $P_{\theta_{\text{lm}}}(y_i|x_i)$.

Here, we augment B_{supp} with semi-supervised examples. In particular, we generate M paraphrases for each x_i using the paraphrase generator $P_{\theta_{\text{par}}}(z_{i,j}|x_i)$, where $z_{i,j}$ represents the *j*-th paraphrase for the input x_i . In the optimal case, this paraphrase will preserve semantic meaning but vary syntactic/lexical form.

We then incorporate the generated paraphrases to create a new mini-batch of examples $B_{s+p} = B_{supp} \cup B_{para}$. Using this augmented mini-batch, we then optimize the following objective function:

$$J_{\theta_{\mathrm{Im}}} := \sum_{i=1}^{N} \{ \log P_{\theta_{\mathrm{Im}}}(y_i | x_i) + 130 \}$$

$$\frac{1}{M} \sum_{j=1}^{M} \log P_{\theta_{\text{Im}}}(y_i | z_{i,j}) \}$$
 (1) 131

128 129

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

167

168

169

170

171

172

173

174

175

2.1 Baseline LM Tuning Techniques

To train the language model using Equation 1, we need to update the parameter set θ_{lm} . One approach would involve updating every parameter for the language model to optimize the training objective (referred to here as the "All-Finetuning" or AllTune approach). However, this method can be computationally intensive. As a result, we will explore the impact of paraphrase augmentation along with six other efficient baseline tuning techniques (Houlsby et al., 2019a) and prompt optimization (Liu et al., 2021b).

We assume that each input x or its paraphrase z is preceded by the task instruction p, which is often specified in previous works. The task instruction, which we represent using the symbol p to be consist with prompt optimization literature, serves as a parameter-free, gradient-free technique for enhancing the performance of the PLM across various downstream tasks (Brown et al., 2020b; Petroni et al., 2019; Deng et al., 2022). When using only the task instructions, no parameters for the language model are updated ($\theta_{Im} = \emptyset$), and zeroshot predictions are made solely on the evaluation data. We further investigate the following language model tuning techniques while incorporating these task instructions into the input or its paraphrases.

Gradient-Search (GS): The GS technique is based on the recent AUTOPROMPT (Shin et al., 2020) method, which optimizes task instructions without updating any parameters in the model. The search process begins in the vocabulary space, optimizing the change in label log-likelihood when replacing token p_i in the task instruction with another token v from the vocabulary set. In our implementation, each search iteration randomly selects one mini-batch of training examples and then randomly selects a token from the task instruction to update. The top k candidate tokens are determined based on the approximate change in label log-likelihood: $Top_v \{w_v^T \cdot \nabla_{w_{p_i}} \log P_{lm}(y|p, x)\},\$ where w_v is the embedding vector of a candidate token v. The resulting k new task instructions are evaluated again using label log-likelihood on the

265

266

267

268

269

270

271

272

same training examples¹, and the top-performing instruction is retained for the next search iteration. Prompt optimization always uses the original input x when searching new task prompts (Shin et al., 2020; Deng et al., 2022). In our work, we investigate the impact of incorporating paraphrases of xduring search.

176

177

178

179

181

182

185

187

189

190

191

193

194

195

197

198

199

205

210

211

212

213

214

215

216

217

218

219

Input-Finetuning (InTune): As a straightforward and efficient tuning technique, we compare to updating only the input embedding table in the transformer architecture. This method requires gradient computation similar to All-Finetuning (All-Tune) as well as the GS method.

LM-Head-Finetuning (HTune): The transformer-based pre-trained language models consist of a language modeling head, which maps the hidden vectors to the token logit for each token in the vocabulary. For the HTune technique, we solely update the parameters of the language modeling head.

Classifier-Finetuning (ClsTune): In ClsTune, we first create a feature representation h(x) for the input text x using average pooling of the final hidden vectors in the last layer of the language model. Here, we assume that the language model (feature extractor) remains fixed, and we then construct a two-layer feedforward network with the *gelu* activation function (Hendrycks and Gimpel, 2016) as a classification module on top of the language model.

Softprompt-Tuning (SpTune): In SpTune (Lester et al., 2021), L prompt tokens are prepended to the task instruction. These L tokens are associated with L dedicated prompt embedding vectors, extending the sequence of vectors derived from the task instruction and input text with an additional L trainable feature vectors. During training, the original embedding table of the transformer model remains fixed, while a new prompt embedding table is trained by backpropagating the label log-likelihood into the prompt embedding table. In contrast to InTune, here the prompt vectors do not need to map to vocabulary words.

Low-Rank Adaptation (LoRA): LoRA is one of the latest efficient-tuning techniques specifically designed for PLMs (Hu et al., 2021). It learns lowrank adaptation matrices for the query and value weight matrices within the transformer model. For a pre-trained weight matrix $W_q \in \mathbb{R}^{d \times k}$, LoRA learns the necessary adaptation (i.e., modification) of the weight matrix for a downstream task through a low-rank decomposition, expressed as $W_q + \Delta W_q \approx W_q + BA$. Here, $B \in \mathbb{R}^{d \times r}$, $A \in \mathbb{R}^{r \times k}$, and the rank $r \leq \min(d, k)$. The adaptation matrices A and B are the only parameters subject to training, while the original matrix W_q does not receive any gradient updates. Studies have shown that LoRA performs on par with, or better than, AllTune across various PLMs (Hu et al., 2021).

All language model tuning techniques we have discussed will use the same input format. For example in the sentiment classification task, we use the following format:

"<*s*> {*instruction*} {*text*}. *It was* <*mask*>. <*/s*>".

Except for ClsTune, all of our tuning techniques maximize the probability of the correct label token in place of the <mask> token. In contrast, ClsTune takes the formatted input and classifies it into one of the predefined class labels.

2.2 LM-Friendly Paraphrase Search

Given a training example (x, y), our objective is to assign the gold label y to the input x by maximizing the log likelihood $\log P(y|x)$. We leverage the fact that when x is misclassified, there may exist paraphrases of the input x that lead to the correct class prediction. These paraphrases should retain the semantic meaning of x while exhibiting syntactic differences, akin to the way we rephrase things when we have been misunderstood. Thus, we generate paraphrases z_i based on the input x_i , that enable the downstream language model to predict the correct label y with greater confidence. Consequently, our data log likelihood is factorized into the following marginalization over the space of paraphrases, where θ_{par} and θ_{lm} represent the parameters for the paraphrase generator and the downstream language model, respectively:

$$J_{\theta_{\text{par}}} := \log P(y|x) = \log E_{z \sim \theta_{\text{par}}}[P(y, z|x)]$$
$$= \log \sum_{z} P_{\theta_{\text{par}}}(z|x) \times P_{\theta_{\text{lm}}}(y|z) \quad (2)$$

To train the paraphrase generator and optimize the objective stated in Equation 2, we explore four distinct learning aspects: (a) two methods for gradient approximation, (b) a reward normalization technique, (c) three decoding techniques for sampling paraphrases, and (d) two approaches to ensure

¹The original AUTOPROMPT evaluates the new candidate instructions on another training mini-batch. For fewshot classification, we re-use the drawn training mini-batch to evaluate the complete new candidate instructions.

grammatical integrity during paraphrase generation.
By combining these elements, we examine various
learning approaches to refine the paraphrase generator with the aid of the downstream language model.
In the subsequent paragraphs, we will describe our
suggested options for each aspect.

281

282

290

293

294

299

301

302

305

307

310

Gradient Approximation: Text generation can be reformulated as an episodic reinforcement learning problem where an agent (i.e. a paraphrase generator) generates tokens (i.e. takes actions) one step at a time until reaching the end of the episode (i.e. selecting the end of sequence token). Therefore, for a given training example (x, y) and its paraphrase z, we define the terminal reward (i.e. goodness) for z as $R(z) = \log P_{\theta_{\text{lm}}}(y|z)$. When approximating the gradient vector of objective 2 concerning θ_{par} , we propose two strategies. These include: (i) Maximum Marginal Likelihood (MML) and (ii) approximating the gradient vector of the paraphrase model via the Policy Gradient (PG) theorem. Notably, gradient updates using these two methods exhibit a close relationship, with the main difference lying in the posterior coefficient utilized to score each sample (Guu et al., 2017). We can recast the main objective presented in equation 2 into the following function representing the expected reward:

$$J_{\theta_{\text{par}}} := \log E_{z \sim P_{\theta_{\text{par}}}(.|x|)}[e^{R(z)}]$$

Given each input x, if we extract paraphrase samples from $P_{\theta_{\text{par}}}(.|x)$ and approximate the expectation in $J_{\theta_{\text{par}}}$ via numerical summation, we optimize the objective using MML estimation. This process results in the following gradient update:

$$\nabla J_{\theta_{\text{par}}}^{\text{MML}} := \nabla_{\theta_{\text{par}}} \log E_z[e^{R(z)}]$$

$$= \sum_{j=1}^M \phi^{\text{MML}}(z_j) \times \nabla_{\theta_{\text{par}}} \log P_{\theta_{\text{par}}}(z_j|x)$$

$$\phi^{\text{MML}}(z_j) = \frac{P_{\theta_{\text{par}}}(z_j|x) \times e^{R(z_j)}}{\sum_{j'=1}^M P_{\theta_{\text{par}}}(z_{j'}|x) \times e^{R(z_{j'})}} \quad (4)$$

By introducing the log inside the expectation (applying Jensen's inequality), we can optimize a surrogate lower bound for the objective presented in equation 3, resulting in the following policy gradient approximation (Sutton et al., 1999):

л*л*

$$\nabla J_{\theta_{\text{par}}}^{\text{PG}} := \nabla_{\theta_{\text{par}}} E_z[R(z)]$$
317

$$= \sum_{j=1}^{M} \phi^{\mathrm{PG}}(z_j) \times \nabla_{\theta_{\mathrm{par}}} \log P_{\theta_{\mathrm{par}}}(z_j|x)$$
 318

$$\phi^{\mathrm{PG}}(z_j) = P_{\theta_{\mathrm{par}}}(z_j|x) \times R(z_j) \quad (5)$$

315

316

320

321

322

323

324

325

326

327

328

331

332

333

336

337

338

339

340

341

342

343

344

345

346

347

348

350

351

352

354

355

356

357

358

Reward Normalization: For our secondary learning aspect, we can either utilize the basic reward, denoted as $R(z_j)$, or normalize the rewards among the paraphrases of a given input x. This process of normalization is particularly useful because it prevents the training of the paraphrase generator with rewards of varying magnitudes, as different training examples are not equally challenging for the language model. Prior research suggests that such normalization of rewards can significantly enhance the performance of text generators across a variety of tasks (Guo et al., 2022). The normalized reward R^n is defined as follows:

$$R^{n}(z_{j}) = \frac{R(z_{j}) - \mu_{j}}{\sigma_{j}}, \mu_{j} = \frac{1}{M} \sum_{j=1}^{M} R(z_{j})$$

$$\sigma_{i}^{2} = \frac{1}{M} \sum_{j=1}^{M} (R(z_{j}) - \mu_{j})^{2} \quad (6)$$
334

$$R_j^2 = \frac{1}{M} \sum_{j=1}^{M} (R(z_j) - \mu_j)^2$$
 (6)

Decoding Techniques: To train the paraphrase generator, we use both the MML and PG gradient estimations which necessitates drawing M samples from the paraphrase generator. We implement three decoding techniques for this purpose. Firstly, we utilize diverse beam search decoding (Vijayakumar et al., 2018) to gather these M paraphrases. Secondly, in order to thoroughly explore the paraphrase space, we alternatively collect the M paraphrases using nucleus (top-p) sampling (Holtzman et al., 2020). For the top-p sampling, we establish a sampling threshold of p = 0.99, at which we collect the minimal set of tokens from the vocabulary with a cumulative probability of at least 0.99. We then re-sample tokens from this set. And thirdly, during the training phase we blend diverse beam search and top-p sampling. Here, we initially sample M paraphrases using both methods, then combine the top M/2 samples from each output to construct our final M samples. During the test phase, we only use diverse beam search.

Grammatical Integrity: We describe three distinct modeling techniques for both the MML and

(3)

443

444

445

446

447

448

405

406

PG gradient estimations: On-policy learning, Offpolicy learning and Proximal Policy Optimization (PPO).

359

374

381

384

400

401

402

403

404

As we are sampling paraphrases from $P_{\theta_{par}}(z_j|x)$ and updating θ_{par} using these samples, the paraphrase generator may start generating ungrammatical text during this default on-policy learning setting. Similar instances of degenerate generation have been reported in tasks like question generation (Najafi and Fyshe, 2023) and program synthesis (Liang et al., 2018).

To mitigate degenerate paraphrase generation, we experiment with off-policy sampling. Here, we maintain a fixed sampling module $P_{\text{fixed}}(z_j|x)$ for sample selection, then update the main paraphrase generator $P_{\theta_{\text{par}}}(z_j|x)$ within the frameworks of objectives 4 and 5. Consequently, with these off-policy samples, the posterior coefficients incorporate the importance sampling ratio $s(z_j) = \frac{P_{\theta_{\text{par}}}(z_j|x)}{P_{\text{fixed}}(z_j|x)}$

$$\phi_{\text{off}}^{\text{PG}}(z_j) = s(z_j) \times R(z_j)$$
$$\phi_{\text{off}}^{\text{MML}}(z_j) = \frac{s(z_j) \times e^{R(z_j)}}{\sum_{j'=1}^M s(z_{j'}) \times e^{R(z_{j'})}} \quad (7)$$

Our next solution for degenerate paraphrases involves imposing a penalty in the training objective if the samples drawn from the current paraphrase generator, $P_{\theta_{\text{par}}}(z|x)$, deviate from those of the pre-trained paraphrase generator. We can implement this penalty as a KL-divergence penalty between the distributions of paraphrases produced by the current model and the pre-trained one. This approach resembles the PPO learning with a KL penalty (Schulman et al., 2017) and has been used in fine-tuning InstructGPT with a reward model trained over human feedback (Ouyang et al., 2022). In the case of InstructGPT, researchers prevent the reward fine-tuned model from diverging from a separate language model pre-trained on supervised data (Ouyang et al., 2022). To integrate this penalty we define the following new objective for θ_{par} :

$$J_{\theta_{\text{par}}}^{\text{PPO}} := \log E_z[e^{R(z)}] - \beta E_z[\log s(z)]$$
$$s(z) = \frac{P_{\theta_{\text{par}}}(z|x)}{P_{\text{fixed}}(z|x)} \quad (8)$$

Building upon the previously approximated MML and PG gradients, we can now derive the following regularized gradient vector with respect to θ_{par} . Please note that β is a hyper-parameter in this context:

$$\nabla J_{\theta_{\text{par}}} - \beta E_z[(\log s(z) + 1)\nabla \log P_{\theta_{\text{par}}}(z|x)]$$
$$z \sim P_{\theta_{\text{par}}}(.|x) \quad (9)$$

Note that the KL penalty can be interpreted as the sum of a grammar reward, denoted by $\log P_{\text{fixed}}(z|x)$, and an entropy regularization term over $P_{\theta_{\text{par}}}(z|x)$. The entropy regularization aids in the diverse exploration of the search space (Mnih et al., 2016), while the grammar reward discourages the learning of ungrammatical samples.

2.3 Ensemble Inference

After optimizing Equation 2 and fine-tuning our paraphrase generator, we generate weaklysupervised examples for inclusion in Equation 1 to train our downstream language model.

To predict the class label of a test example, we could either use our fine-tuned language model to predict the class label based on the original input x, or adopt an ensemble approach. For the latter, for a given x, we generate M paraphrases using our fine-tuned paraphrase generator. We then average the prediction scores for a potential class label across the M + 1 values according to Equation 1 to predict the class label for that input example x. This aligns with our earlier assumption that some paraphrases could be easier for the language model to predict the correct class label. During data augmentation for the language model, we select the validation set's best model according to this ensemble prediction.

3 Experiments

3.1 Setup

Pre-trained Models:

For paraphrase generation, we employ a T5-base model (Raffel et al., 2019) which has been trained on paraphrases generated by ChatGPT². These output paraphrases were generated for input texts from various datasets, including Quora paraphrase questions, texts from SQUAD 2.0, and the CNN news dataset³. To create this training data, ChatGPT generated five paraphrases for each input, which were then used as the target for the T5-base model. The

²https://openai.com/blog/chatgpt

³You can find detailed dataset information here: https://huggingface.co/datasets/humarin/ chatgpt-paraphrases

weights for this model are publicly available⁴. In our experiments, this model was able to generate more diverse paraphrases compared to other public pre-trained models.

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

491

492

493

494

495

496

497

For our main language model, we use the RoBERTa-large model pre-trained with the Masked Language Modeling (MLM) objective (Liu et al., 2019), which has demonstrated strong performance on NLU tasks. Our proposed learning framework can be readily extended to other paraphrase generators or backbone language models.

Datasets: Inspired by prior work (Gao et al., 2021; Deng et al., 2022), we experiment on seven classification tasks in the few-shot setting. These include sentiment classification tasks such as the binary sentiment datasets SST2 (Socher et al., 2013), CR (Hu and Liu, 2004), and MR (Pang and Lee, 2005). We also experiment on the 5-label sentiment dataset SST5 (Socher et al., 2013), the binary subjectivity classification dataset Subj (Pang and Lee, 2004), the question type classification dataset TREC (Voorhees and Tice, 2000), and the topic classification dataset AG's News (Zhang et al., 2015). The number of classes per dataset, as well as the used instructions are outlined in Appendix E. Instructions and class verbalizers are based on previous work (Deng et al., 2022) in prompt optimization. Detailed information about the specific learning rates for each LM technique along with other hyper-parameters can be found in Appendix B.

3.2 Few-shot Paraphrase Training

As discussed in Section 2.2, there are four learning aspects to be considered when fine-tuning our paraphrase generator for the downstream language model. We conduct an extensive set of experiments in the 128-shot setting for the SST2 binary sentiment classification task.

We randomly select 128 training examples for each unique label within the dataset. An equal number of examples are gathered to form an internal validation set. We create five train/validation splits using the arbitrarily chosen random seeds {11, 42, 1993, 12321, 2023}. We train the models for 1120 training steps with the batch size of 8 (i.e. 35 epochs). As we are training the models, we evaluate the performance of 140 weight checkpoints per model on the validation splits (i.e one checkpoint per 8 training steps). We examine the mean accuracy, which is averaged over the five validation

⁴https://huggingface.co/humarin/chatgpt_ paraphraser_on_T5_base



Figure 1: Average ensemble accuracy over five validation splits in the 128-shot SST2 classification task. PG gradient estimation is not robust during the training trajectory while doing on-policy learning.

splits. Despite the ensembling approach described in Section 2.3, to accurately capture the quality of the generated paraphrases, we exclude the original input x when computing the ensemble accuracy on the validation splits.

We assess the impact of reward normalization in the context of on-policy, off-policy, and PPO learning, considering both PG and MML gradient estimations. Table 1 lists the best performance out of all the checkpoints evaluated on the validation splits, which is further averaged over five validation splits. With both PG and MML gradient estimations, reward normalization is boosting the performance across the three text decoding techniques for both on-policy and PPO learning techniques (see 'AVG' column in Table 1). Conversely, reward normalization is not improving performance with off-policy learning (follow discussion in Appendix A and see Table 3).

Table 1 verifies that MML gradient estimation outperforms PG gradient estimation on average across three decoding techniques for both on-policy and PPO learning techniques. The highest accuracy is acheived by 'PG + zscoring' with on-policy learning and top-p decoding, however it is not robust during the entire training trajectory. Figure 1 shows that PG gradient estimation is not robust throughout the training trajectory, which causes the paraphrase generator to produce nonsensical paraphrases. This results in downstream classification performance on par with random guessing. In contrast, off-policy and PPO learning circumvent this divergence. MML gradient estimation maintains robustness throughout the training phase (further

531

498

499

Table 1: The accuracy of the best performing validation checkpoints in the 128-shot SST2 classification task for both the on-policy and PPO learning techniques. Highest performance per column bolded. Last column reports the macro-average among each row.

Learning	0	On-Polic	су	PPO			AVG
Technique	Top-P	Beam	Mixed	Top-P	Beam	Mixed	
No Tuning	67.5	67.5	67.5	67.5	67.5	67.5	67.5
PG	67.9	68.0	67.9	68.5	68.3	69.1	68.3
PG +zscoring	71.3	70.2	71.2	68.9	68.8	69.8	70.0
MML	69.6	69.1	69.8	69.5	69.9	70.5	69.7
MML +zscoring	70.3	70.2	70.2	68.9	70.3	70.6	70.1

discussion in Table 4 of Appendix A).

532

533

534

536

538

540

541

542

543

544

545

546

549

552

553

554

557

558

559

560

566

570

Upon investigating various elements of our learning objectives for fine-tuning the paraphrase generator, the combination that delivers the best performance across the validation splits, which is also robust during the entire training trajectory, includes: MML gradient approximation, PPO learning, mixed decoding for sample generation, and finally applying reward normalization. We name this combined approach our proposed RIFF algorithm. Table 5 lists an example sentence with its generated paraphrases on the SST2 dataset. In the subsequent experiments, we applied RIFF to generate paraphrases that augment the training mini-batches while tuning the downstream language model in a few-shot setting.

3.3 Paraphrases for Few-shot LM Tuning

Our primary hypothesis is that various LM tuning techniques could benefit from diverse views of the original input text. To test this hypothesis, we fine-tuned our paraphrase generators in a 16shot classification setup using the RIFF algorithm. Subsequently, we fine-tuned the downstream classification model in the same 16-shot setting, while introducing M = 8 paraphrases as per the objective outlined in Equation 1. For evaluation, we used the best model from the validation set to make predictions on standard evaluation splits, following the ensemble approach described in Section 2.3. For consistency with prior research, we used the random dataset splits provided by RLPrompt (Deng et al., 2022), aligning with the random seeds used by LM-BFF (Gao et al., 2021).

Table 2 illustrates the average accuracy on standard test sets across seven text classification datasets. The reported scores correspond to four distinct LM tuning techniques: GS, SpTune, All-Tune, and LoRA. Results for ClsTune, HTune, and InTune are presented in Table 8 of Appendix C.

Notably, recent prompt optimization techniques (i.e. GS and SpTune) exhibit significant benefits from paraphrase augmentation during training. Particularly, SpTune demonstrates the most substantial improvement, with a 2.7% increase in accuracy. Furthermore, LoRA consistently outperforms prompt optimization techniques, aligning with findings from prior studies (Hu et al., 2021). Interestingly, paraphrase augmentation still aids LoRA in efficiently learning adaptation matrices, resulting in a 0.5% accuracy increase on SST2, a 0.4% accuracy boost on Subj, and a 0.3% accuracy improvement on AGN datasets. Coupled with ensemble predictions, denoted in rows labeled '+RIFF (train+test)' all LM tuning techniques benefit from paraphrase augmentation.

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

4 Related Works

To improve prompt optimization and efficient tuning techniques for LMs, we incorporate the generated paraphrases into the training mini-batches. It is important to note that data augmentation in NLP has been explored through various methods. Paraphrase generation represents just one technique of data augmentation. For a comprehensive overview of diverse data augmentation techniques for NLP tasks, we direct interested readers to a recent survey by Chen et al. (2021).

Prompt Optimization & Efficient Tuning: Recent research proposes various techniques for prompt optimization and efficient tuning of language models. In our experiments, we have used successful techniques from each of these areas. Appendix D provides our detailed description of these recent techniques. All of the recent techniques for prompt optimization and efficient tuning of the language models use the original input task (or the original input context) provided within the dataset. **Paraphrase Generation**: Our objective is not to

Table 2: Average accuracy on the standard evaluation sets for the 16-shot single-sentence text classification. The last column is the micro averaged performance across the datasets. Numbers in parentheses are the deltas compared to non-RIFF baseline per LM Tuning technique. Highest performance per dataset bolded, second highest underlined. †: the average 16-shot fine-tuning results with automatically searched templates. *****: reported zero-shot results for GPT3 with in-context learning (Gao et al., 2021).

Tuning Method	SST2	SST5	CR	MR	TREC	Subj	AGN	AVG
(LM = RoBERTa-large)								
GPT-3 In-context Learning*	84.8	30.6	87.4	80.5	26.2	53.6	-	64.4
(Deng et al., 2022) (RLPrompt)	92.5	41.4	89.5	<u>87.1</u>	60.5	81.9	80.2	78.1
(Gao et al., 2021)†	92.3	49.2	89.0	85.5	<u>88.2</u>	91.2	-	80.9
No Tuning + Instructions	84.6	31.0	77.8	81.3	27.6	57.7	51.5	58.5
+GS	85.5	37.0	80.2	83.0	45.3	74.5	82.0	74.9
+RIFF (train)	86.4	37.8	82.7	84.7	51.0	74.4	81.0	75.3 (+0.4)
+RIFF (train+test)	87.3	38.2	85.1	84.7	52.4	77.2	83.3	77.0 (+2.1)
+SpTune	89.7	39.4	82.4	86.1	35.2	72.4	82.0	75.7
+RIFF (train)	91.2	44.5	84.6	86.1	38.4	79.7	84.0	78.5 (+2.7)
+RIFF (train+test)	91.6	45.1	86.2	86.6	38.4	81.6	86.0	79.9 (+4.1)
+AllTune	93.1	48.0	89.2	87.3	87.2	85.8	87.7	<u>83.3</u>
+RIFF (train)	<u>93.6</u>	<u>50.6</u>	<u>90.2</u>	85.8	84.2	85.3	87.2	<u>83.3</u> (+0.0)
+RIFF (train+test)	93.8	51.2	91.0	85.5	84.4	<u>86.6</u>	87.2	83.6 (+0.3)
+LoRA	92.5	48.1	88.6	86.0	89.3	81.5	87.3	82.5
+RIFF (train)	92.7	48.0	87.5	85.1	84.8	81.9	<u>87.6</u>	82.3 (-0.2)
+RIFF (train+test)	93.1	49.2	89.0	85.4	85.9	84.4	87.9	83.1 (+0.6)

609

610

611

619 620

621

623

62

62

628

629

630 631

632 633 present a state-of-the-art paraphrase generator, but rather to examine the impact of incorporating input paraphrases on the efficient tuning of LMs. Recent advancements in generating diverse paraphrases (Zhou and Bhat, 2021) could provide better pre-trained models, thereby these techniques can enhance performance in all our experiments as our proposed RIFF technique can be seen as an extra fine-tuning step for the paraphrase model. These recent techniques encompass various approaches, including the use of copy mechanisms, Variational Autoencoders, Generative Adversarial Networks, and Reinforcement Learning techniques. For a comprehensive overview of neural paraphrase generation, please refer to a recent survey by Zhou and Bhat (Zhou and Bhat, 2021). While previous studies have applied RL techniques for paraphrase generation, we propose the use of MML gradients instead of policy gradients to fine-tune our paraphrase model.

5 Conclusion

We investigated the impact of incorporating input paraphrases while fine-tuning PLMs with recent prompt optimization and efficient tuning techniques. We also provided extensive experiments for reducing noise in a distantly supervised paraphrase generator.

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

Recent study suggests that large PLMs face constraints related to the quantity of unique data points (Muennighoff et al., 2023). A potential avenue for future research could explore whether the introduction of paraphrased inputs can mitigate these challenges for large PLMs during the pretraining stage.

Limitations

Our paraphrase generator is pre-trained on semisupervised paraphrases given by a truly large language model (i.e. ChatGPT). Although these large models are capable of generating high quality paraphrases for the English language. It is not clear if these semi-supervised paraphrases are available for other languages.

Our investigation into the effects of paraphrases on efficient LM tuning techniques relies on language models that can be trained on commodity hardware. It would be interesting to explore how recent LLMs respond to input paraphrases. Nevertheless, it is important to note that tuning larger models, those with more than 30 billion parameters, requires access to the gradients, however, the gradients for proprietary models are not available.

Throughout this study, our primary focus has been on classification tasks. A potential avenue for future research could involve assessing the sensitivity of LM tuning techniques to input paraphrases in generative tasks.

To enhance language model tuning with paraphrases, we augment the training mini-batches. However, this approach does result in increased training time. It would be ideal to devise a form of regularization for efficient language model tuning, one that exposes the model's parameters to various paraphrases of the original input text.

Ethics Statement

659

660

662

672

673

675

676

677

679

684

692

694

701

702

703

704

707

710

711

Many language models show biases in their output due to the data used to train them (Liang et al., 2021). It is possible that even with few-shot language model tuning, we might continue to detect analogous biases in the downstream classification task, for instance, resulting in diminished classification accuracy for specific minority groups.

References

- James Urquhart Allingham, Jie Ren, Michael W Dusenberry, Xiuye Gu, Yin Cui, Dustin Tran, Jeremiah Zhe Liu, and Balaji Lakshminarayanan. 2023. A simple zero-shot prompt weighting technique to improve prompt ensembling in text-image models.
- Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, Eric Chu, Jonathan H. Clark, Laurent El Shafey, Yanping Huang, Kathy Meier-Hellstern, Gaurav Mishra, Erica Moreira, Mark Omernick, Kevin Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao, Yuanzhong Xu, Yujing Zhang, Gustavo Hernandez Abrego, Junwhan Ahn, Jacob Austin, Paul Barham, Jan Botha, James Bradbury, Siddhartha Brahma, Kevin Brooks, Michele Catasta, Yong Cheng, Colin Cherry, Christopher A. Choquette-Choo, Aakanksha Chowdhery, Clément Crepy, Shachi Dave, Mostafa Dehghani, Sunipa Dev, Jacob Devlin, Mark Díaz, Nan Du, Ethan Dyer, Vlad Feinberg, Fangxiaoyu Feng, Vlad Fienber, Markus Freitag, Xavier Garcia, Sebastian Gehrmann, Lucas Gonzalez, Guy Gur-Ari, Steven Hand, Hadi Hashemi, Le Hou, Joshua Howland, Andrea Hu, Jeffrey Hui, Jeremy Hurwitz, Michael Isard, Abe Ittycheriah, Matthew Jagielski, Wenhao Jia, Kathleen Kenealy, Maxim Krikun, Sneha Kudugunta, Chang Lan, Katherine Lee, Benjamin Lee, Eric Li, Music Li, Wei Li, YaGuang Li, Jian Li, Hyeontaek Lim, Hanzhao Lin, Zhongtao Liu, Frederick Liu, Marcello Maggioni, Aroma Mahendru, Joshua Maynez, Vedant Misra, Maysam Moussalem,

Zachary Nado, John Nham, Eric Ni, Andrew Nystrom, Alicia Parrish, Marie Pellat, Martin Polacek, Alex Polozov, Reiner Pope, Siyuan Qiao, Emily Reif, Bryan Richter, Parker Riley, Alex Castro Ros, Aurko Roy, Brennan Saeta, Rajkumar Samuel, Renee Shelby, Ambrose Slone, Daniel Smilkov, David R. So, Daniel Sohn, Simon Tokumine, Dasha Valter, Vijay Vasudevan, Kiran Vodrahalli, Xuezhi Wang, Pidong Wang, Zirui Wang, Tao Wang, John Wieting, Yuhuai Wu, Kelvin Xu, Yunhan Xu, Linting Xue, Pengcheng Yin, Jiahui Yu, Qiao Zhang, Steven Zheng, Ce Zheng, Weikang Zhou, Denny Zhou, Slav Petrov, and Yonghui Wu. 2023. Palm 2 technical report. 712

713

714

715

716

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

- David Broughton. 1995. *The assumptions and theory of public opinion polling*, pages 15–33. Macmillan Education UK, London.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020a. Language models are few-shot learners.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam Mc-Candlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020b. Language models are few-shot learners. *CoRR*, abs/2005.14165.
- Jiaao Chen, Derek Tam, Colin Raffel, Mohit Bansal, and Diyi Yang. 2021. An empirical survey of data augmentation for limited data learning in NLP. *CoRR*, abs/2106.07499.
- Ganqu Cui, Wentao Li, Ning Ding, Longtao Huang, Zhiyuan Liu, and Maosong Sun. 2023. Decoder tuning: Efficient language understanding as decoding. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15072–15087, Toronto, Canada. Association for Computational Linguistics.
- Mingkai Deng, Jianyu Wang, Cheng-Ping Hsieh, Yihan Wang, Han Guo, Tianmin Shu, Meng Song, Eric Xing, and Zhiting Hu. 2022. RLPrompt: Optimizing discrete text prompts with reinforcement learning. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3369–3391, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

880

881

882

883

- 771 772
- 77 77
- 77
- 77
- 77
- 7
- 781 782 783
- 784
- 7
- _
- 7
- 7
- 79 79

795

- 796
- 7 7
- 79 80
- 801 802 803
- 804 805
- 8
- 809 810
- 811 812
- 813
- 814
- 816
- 817 818
- 819 820 821

- 82
- 825 826

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Margaret C. Donaldson. 1978. Children's Minds.
 - Wanyu Du and Yangfeng Ji. 2019. An empirical comparison on imitation learning and reinforcement learning for paraphrase generation. 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), pages 6012–6018, Hong Kong, China. Association for Computational Linguistics.
 - Steven Y. Feng, Varun Gangal, Jason Wei, Sarath Chandar, Soroush Vosoughi, Teruko Mitamura, and Eduard Hovy. 2021. A survey of data augmentation approaches for NLP. In *Findings of the Association* for Computational Linguistics: ACL-IJCNLP 2021, pages 968–988, Online. Association for Computational Linguistics.
 - Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. Making pre-trained language models better few-shot learners. 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), pages 3816–3830, Online. Association for Computational Linguistics.
 - Sonal Garg, Sumanth Prabhu, Hemant Misra, and G. Srinivasaraghavan. 2021. Unsupervised contextual paraphrase generation using lexical control and reinforcement learning. *CoRR*, abs/2103.12777.
 - Han Guo, Bowen Tan, Zhengzhong Liu, Eric Xing, and Zhiting Hu. 2022. Efficient (soft) Q-learning for text generation with limited good data. In *Findings of the Association for Computational Linguistics: EMNLP* 2022, pages 6969–6991, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
 - Han Guo, Bowen Tan, Zhengzhong Liu, Eric P. Xing, and Zhiting Hu. 2021. Efficient (soft) q-learning for text generation with limited good data.
 - Kelvin Guu, Panupong Pasupat, Evan Liu, and Percy Liang. 2017. From language to programs: Bridging reinforcement learning and maximum marginal likelihood. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1051–1062, Vancouver, Canada. Association for Computational Linguistics.
 - Dan Hendrycks and Kevin Gimpel. 2016. Bridging nonlinearities and stochastic regularizers with gaussian error linear units. *CoRR*, abs/1606.08415.

- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019a. Parameter-efficient transfer learning for NLP. In Proceedings of the 36th International Conference on Machine Learning, volume 97 of Proceedings of Machine Learning Research, pages 2790–2799. PMLR.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019b. Parameter-efficient transfer learning for NLP. *CoRR*, abs/1902.00751.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *CoRR*, abs/2106.09685.
- Minqing Hu and Bing Liu. 2004. Mining and summarizing customer reviews. In *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '04, page 168–177, New York, NY, USA. Association for Computing Machinery.
- Pavel Izmailov, Dmitrii Podoprikhin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. 2019. Averaging weights leads to wider optima and better generalization.
- Takeshi Kojima, Shixiang (Shane) Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. In *Advances in Neural Information Processing Systems*, volume 35, pages 22199–22213. Curran Associates, Inc.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Xiang Lisa Li and Percy Liang. 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), pages 4582– 4597, Online. Association for Computational Linguistics.
- Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. 2023. Making language models better reasoners with step-aware verifier. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pages 5315–5333, Toronto, Canada. Association for Computational Linguistics.

995

996

997

941

- Zichao Li, Xin Jiang, Lifeng Shang, and Hang Li. 2018. Paraphrase generation with deep reinforcement learning. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 3865–3878, Brussels, Belgium. Association for Computational Linguistics.
 - Chen Liang, Mohammad Norouzi, Jonathan Berant, Quoc V Le, and Ni Lao. 2018. Memory augmented policy optimization for program synthesis and semantic parsing. In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.

895

896

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

917

918

919

921

923

925

926

927

928

929

930

931

932

933

934

935

937

- Paul Pu Liang, Chiyu Wu, Louis-Philippe Morency, and Ruslan Salakhutdinov. 2021. Towards Understanding and Mitigating Social Biases in Language Models. In *International Conference on Machine Learning*, pages 6565–6576. PMLR.
- Zhaojiang Lin, Andrea Madotto, and Pascale Fung. 2020. Exploring versatile generative language model via parameter-efficient transfer learning. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 441–459, Online. Association for Computational Linguistics.
- Mingtong Liu, Erguang Yang, Deyi Xiong, Yujie Zhang, Yao Meng, Changjian Hu, Jinan Xu, and Yufeng Chen. 2020. A learning-exploring method to generate diverse paraphrases with multi-objective deep reinforcement learning. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 2310–2321, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021a. Pretrain, prompt, and predict: A systematic survey of prompting methods in natural language processing. *CoRR*, abs/2107.13586.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021b. Pretrain, prompt, and predict: A systematic survey of prompting methods in natural language processing.
- Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. 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)*, pages 61–68, Dublin, Ireland. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019.
 Roberta: A robustly optimized BERT pretraining approach. *CoRR*, abs/1907.11692.
- Ilya Loshchilov and Frank Hutter. 2017. Fixing weight decay regularization in adam. *CoRR*, abs/1711.05101.

- Volodymyr Mnih, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. 2016. Asynchronous methods for deep reinforcement learning. *CoRR*, abs/1602.01783.
- Niklas Muennighoff, Alexander M. Rush, Boaz Barak, Teven Le Scao, Aleksandra Piktus, Nouamane Tazi, Sampo Pyysalo, Thomas Wolf, and Colin Raffel. 2023. Scaling data-constrained language models.
- Saeed Najafi and Alona Fyshe. 2023. Weaklysupervised questions for zero-shot relation extraction. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 3075–3087, Dubrovnik, Croatia. Association for Computational Linguistics.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback.
- Bo Pang and Lillian Lee. 2004. A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. In *Proceedings* of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL-04), pages 271–278, Barcelona, Spain.
- Bo Pang and Lillian Lee. 2005. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *Proceedings of the* 43rd Annual Meeting of the Association for Computational Linguistics (ACL'05), pages 115–124, Ann Arbor, Michigan. Association for Computational Linguistics.
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? 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), pages 2463–2473, Hong Kong, China. Association for Computational Linguistics.
- Lihua Qian, Lin Qiu, Weinan Zhang, Xin Jiang, and Yong Yu. 2019. Exploring diverse expressions for paraphrase generation. 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), pages 3173–3182, Hong Kong, China. Association for Computational Linguistics.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

- 999 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1018 1019 1020 1021 1022 1023 1024 1025 1026 1027 1028 1031 1032 1033 1034

- 1040 1041 1042 1043 1044 1046

1048 1049

1054

- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yangi Zhou, Wei Li, and Peter J. Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. CoRR, abs/1910.10683.
- Sashank J. Reddi, Satyen Kale, and Sanjiv Kumar. 2019. On the convergence of adam and beyond. CoRR, abs/1904.09237.
- Steven J. Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. 2016. Self-critical sequence training for image captioning. CoRR, abs/1612.00563.
- Abhilasha Sancheti, Balaji Vasan Srinivasan, and Rachel Rudinger. 2022. Entailment relation aware paraphrase generation. Proceedings of the AAAI Conference on Artificial Intelligence, 36(10):11258-11266.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. CoRR, abs/1707.06347.
- Weijia Shi, Xiaochuang Han, Hila Gonen, Ari Holtzman, Yulia Tsvetkov, and Luke Zettlemoyer. 2022. Toward human readable prompt tuning: Kubrick's the shining is a good movie, and a good prompt too?
- Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. 2020. AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 4222-4235, Online. Association for Computational Linguistics.
- A. B. Siddique, Samet Oymak, and Vagelis Hristidis. 2020. Unsupervised paraphrasing via deep reinforcement learning. In Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD '20, page 1800–1809, New York, NY, USA. Association for Computing Machinery.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing, pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.
- Tianxiang Sun, Yunfan Shao, Hong Qian, Xuanjing Huang, and Xipeng Qiu. 2022. Black-box tuning for language-model-as-a-service. CoRR. abs/2201.03514.
- Richard S. Sutton, David McAllester, Satinder Singh, and Yishay Mansour. 1999. Policy gradient methods for reinforcement learning with function approximation. In Proceedings of the 12th International Conference on Neural Information Processing Systems, NIPS'99, page 1057–1063, Cambridge, MA, USA. MIT Press.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-1062 thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, 1063 Isabel Kloumann, Artem Korenev, Punit Singh Koura, 1064 Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-1065 ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-1066 tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-1067 bog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subrama-1070 nian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, 1073 Melanie Kambadur, Sharan Narang, Aurelien Ro-1074 driguez, Robert Stojnic, Sergey Edunov, and Thomas 1075 Scialom. 2023. Llama 2: Open foundation and fine-1076 tuned chat models. 1077

1055

1056

1058

1078

1079

1080

1081

1082

1083

1084

1086

1089

1090

1091

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

- Mojtaba Valipour, Mehdi Rezagholizadeh, Ivan Kobyzev, and Ali Ghodsi. 2023. DyLoRA: Parameter-efficient tuning of pre-trained models using dynamic search-free low-rank adaptation. In Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics, pages 3274–3287, Dubrovnik, Croatia. Association for Computational Linguistics.
- Ashwin Vijayakumar, Michael Cogswell, Ramprasaath Selvaraju, Qing Sun, Stefan Lee, David Crandall, and Dhruv Batra. 2018. Diverse beam search for improved description of complex scenes. Proceedings of the AAAI Conference on Artificial Intelligence, 32(1).
- Ellen M. Voorhees and Dawn M. Tice. 2000. Building a question answering test collection. In Proceedings of the 23rd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '00, page 200-207, New York, NY, USA. Association for Computing Machinery.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H. Chi, Quoc Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. CoRR, abs/2201.11903.
- Jason Wei and Kai Zou. 2019. EDA: Easy data augmentation techniques for boosting performance on text classification tasks. 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), pages 6382-6388, Hong Kong, China. Association for Computational Linguistics.
- Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. 2023. Adaptive budget allocation for parameter-efficient fine-tuning.

Table 3: The accuracy of the best performing validation checkpoints in the 128-shot SST2 classification task trained with the off-policy learning technique. Highest performance per column bolded. Last column reports the macro-average among each row.

	C	AVG		
Learn Tech	Top-P	Beam	Mixed	
No Tuning	67.5	67.5	67.5	67.5
PG	68.6	68.4	69.1	68.7
PG +zscoring	68.8	68.7	68.0	68.5
MML	69.2	70.1	70.1	69.8
MML +zscoring	69.2	69.7	70.1	69.7

Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022a. Opt: Open pre-trained transformer language models.

1114

1115

1116

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

- Tianjun Zhang, Xuezhi Wang, Denny Zhou, Dale Schuurmans, and Joseph E. Gonzalez. 2022b. Tempera: Test-time prompting via reinforcement learning.
- Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. *CoRR*, abs/1509.01626.
- Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. 2022c. Automatic chain of thought prompting in large language models.
- Jianing Zhou and Suma Bhat. 2021. Paraphrase generation: A survey of the state of the art. In *Proceedings* of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 5075–5086, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. 2023. Large language models are human-level prompt engineers.

A Few-shot Paraphrase Training (Further Results)

This section provides additional results that compare our training objectives for fine-tuning the paraphrase generator using the feedback from the downstream language model.

The off-policy learning technique improves performance when using basic rewards (i.e., 69.1% compared to 67.9% with mixed decoding). However, the combined effect of off-policy learning and reward normalization decreases performance. With mixed decoding, 'PG + zscoring' yields an accuracy of 71.2% in on-policy learning compared to an accuracy of 68.0% with off-policy learning. The 'AVG' column in Table 3 further verifies this conclusion that reward normalization is not improving the final performance while training the model with off-policy learning. We hypothesize that with the off-policy learning technique, the normalized rewards should be re-weighted properly if the sampled paraphrases are from the fixed paraphrase model. 1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

In Table 4, we also report the average accuracy of all the checkpoints as we are training the models. The learning technique 'MML + zscoring' is more robust during the training trajectory compared to 'PG + zscoring'.

B Further Training Details

The learning rate for each LM tuning technique was separately fine-tuned from the set $\{0.5, 0.3,$ 0.1, 0.01, 0.001, 0.0001, 0.00001 using the train/validation split created for the seed 11 on the SST2 dataset. The tuned learning rates were then applied globally across other datasets and experiments. For paraphrase fine-tuning, we train all the parameters in T5-base with the learning rate of 0.00001. In Tables 6 and 7, we list the hyper-parameters and learning rates used across all datasets. For optimization, we utilized the AdamW (Loshchilov and Hutter, 2017)⁵ optimizer with the AMSGrad variant set to True (Reddi et al., 2019). We implemented the methods using the HuggingFace⁶ library and the PyTorch⁷ machine learning framework. We report the accuracy metric on these classification datasets. The experiments were conducted using multiple NVIDIA's A40 GPU cards.

C Paraphrases for Few-shot LM Tuning (Further Results)

Due to space limitations, we present the results for ClsTune, HTune, and InTune in Table 8.

D Extended Related Works

Prompt Optimization & Efficient Tuning: Recent research proposes various techniques for prompt optimization and efficient tuning of language models. In our experiments, we have used successful techniques from each of these areas.

⁵https://pytorch.org/docs/stable/generated/ torch.optim.AdamW.html

⁶https://huggingface.co/

⁷https://pytorch.org/

Learn Tech	On-Policy				PPO		AVG
	Top-P	Beam	Mixed	Top-P	Beam	Mixed	
No Tuning + Instructions	67.5	67.5	67.5	67.5	67.5	67.5	67.5
+PG	53.2	52.0	52.4	67.4	67.4	68.2	60.1
+PG +zscoring	63.8	68.6	66.6	67.6	67.9	68.5	67.2
+MML	68.5	67.6	68.6	68.2	69.0	68.9	68.5
+MML +zscoring	69.1	68.9	69.1	67.9	69.0	68.9	68.8

Table 4: The average accuracy computed over all the validation checkpoints in the 128-shot SST2 classification task. Highest performance per column bolded. Last column reports the macro-average among each row.

Table 5: Generated paraphrases after fine-tuning the paraphrase model with the RIFF method for the following input sentence: "It's difficult to get beyond the overall blandness of American chai, despite its likable performances and refreshingly naive point of view."

Paraphrase	Text
P1	"Despite its enjoyable performances and refreshingly naive perspective, american chai is
	difficult to surpass due to its overall blandness."
P2	"American chai's enjoyable performances and refreshing perspective make it difficult to
	elevate beyond its blandness."
P3	"The overall blandness of American chai is difficult to overcome, despite its enjoyable
	performances and refreshingly naive perspective."
P4	"It's difficult to surpass the blandness of American chai, even with its enjoyable performances
	and refreshingly naive perspective. "
P5	"Although american chai has decent performances and a refreshingly naive viewpoint, it's
	difficult to elevate its overall blandness. "
P6	"Americans chai is often bland, but it's difficult to surpass its enjoyable performances and
	refreshingly naive perspective."
P7	"American Chai's lack of quality is difficult to overcome, even with enjoyable performances
	and a refreshing perspective. "
P8	"Even with enjoyable performances and a refreshingly simplistic viewpoint, american
	chapin's blandness is difficult to shake off."

Hyper-parameter	Value
Top- k candidates in GS	<i>k</i> =4
batch size (RoBERTa-large)	8
batch size in GS (RoBERTa-large)	2
Weight decay	0.0001
Max epochs	100
length cutoff	128 tokens
Paraphrase sample size	M=8
Checkpointing steps	8
D' in ClsTune	128
Prompt len in SpTune	L=25
β in MML	0.1
β in PG	0.6
Lora α	32
Lora r	8
Lora dropout	0.1
Diversity penalty for Div beam	3.0
Repetition penalty for Div beam	10.0
Temperature in Div beam	0.7
P value for top-p	0.99

Table 6: Shared hyper-parameters used across all experiments and datasets.

Table 7: Learning rates used per Language Model (LM) tuning technique.

LM Tuning Technique	Learning Rate
GS	No rate
AllTune	0.00001
InTune	0.001
HTune	0.001
ClsTune	0.001
SpTune	0.001
LoRA	0.0001

FluentPrompt (Shi et al., 2022) is a recent discrete prompting technique based on the projected gradient-descent and Langevin dynamics. FluentPrompt introduces a fluency constraint within Langevin dynamics to generate a sample of highperforming prompts for more interpretable analysis of these discrete prompts. The optimized prompts by FluentPrompt performs on-par to the Auto-Prompt, however they have lower perplexity (Shi et al., 2022).

1197

1198

1199

1200

1201

1202

1204

1205

1206

1207 1208

1209

1210

1211

1212

Building upon SpTune (Lester et al., 2021) and P-tuning (Li and Liang, 2021), P-tuning V2 (Liu et al., 2022) introduced the concept of deep prompt tuning. This method involves injecting prompt vectors into the deeper layers of the transformer model to close the performance gap with AllTuning in medium-sized language models. We have experimented with LoRA (Hu et al., 2021), a recent lowrank adaptation technique for tuning language models. Other potential methods include training bottleneck adapter modules (Houlsby et al., 2019b; Lin et al., 2020) added per sub-layer of the transformer model. LoRA outperforms adapter tuning and P-Tuning V2 techniques (Hu et al., 2021). The successors of LoRA include DyLoRA (Valipour et al., 2023) which dynamically learns a range of adaptation ranks, thus eliminating the need to search the rank of the adaptation matrices as a hyperparameter. Similarly, AdaLoRA dynamically allocates the parameter budget among the weight matrices during adaptation, with matrices of higher priority (i.e., those with greater importance to the downstream task) receiving higher adaptation ranks than less important matrices (Zhang et al., 2023).

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1262

1263

1264

In scenarios where gradients are absent, Black-Box Tuning (Sun et al., 2022) applies derivativefree algorithms for optimizing continuous prompts. For discrete prompt optimization, RLPrompt (Deng et al., 2022) employs the on-policy version of soft Q-learning (Guo et al., 2021) to find the optimal prompt tokens in a gradient-free setting. Decoder Tuning (Cui et al., 2023) learns a decoder network over the language model, thus circumventing the need for gradient computation and input-side prompt tuning in few-shot classification. In a recent study, TEMPERA (Zhang et al., 2022b) introduced a novel approach that involves test-time discrete prompt editing using a trained RL agent. This agent is capable of modifying the instruction, in-context examples, or the verbalizers based on the given task input.

The use of Language Models (LLMs) in generating instructions for downstream tasks has involved a two-step process. Initially, LLMs generate a set of candidate instructions, and subsequently, the highest-scoring instruction is utilized to prompt another LLM to perform the downstream task. This approach, known as promptbased generation-then-filtering, has been investigated in the recent APE method (Zhou et al., 2023). APE demonstrates the ability to generate prompts that achieve performance comparable to humandesigned prompts (Zhou et al., 2023).

To prompt language models for reasoning tasks, another line of research augment the input context with demonstration examples outlining the intermediate reasoning steps to form the answer. Providing manually or automatically generated chain-of-

Table 8: Average accuracy on the standard evaluation sets for the 16-shot single-sentence text classification using the ClsTune, HTune, and InTune fine-tuning techniques. The last column is the micro averaged performance across the datasets.

Tuning Method	SST2	SST5	CR	MR	TREC	Subj	AGN	AVG
(LM = RoBERTa-large)								
ClsTune	72.6	34.4	71.4	67.3	74.8	87.7	81.7	72.8
+RIFF (train)	72.5	33.9	68.3	70.3	75.8	88.4	84.0	73.8 (+1.0)
+RIFF (train+test)	74.0	35.0	71.1	72.0	76.8	89.4	84.9	75.1 (+2.3)
HTune	87.4	37.4	84.0	83.1	62.4	83.8	81.4	76.9
+RIFF (train)	88.1	40.3	84.5	83.4	70.7	87.3	83.4	78.8 (+2.0)
+RIFF (train+test)	89.1	40.4	86.4	83.2	71.6	87.8	85.2	80.0 (+3.1)
InTune	91.5	42.3	87.3	84.0	67.7	82.7	83.8	79.4
+RIFF (train)	92.6	43.2	87.5	85.9	63.8	82.7	85.6	80.5 (+1.1)
+RIFF (train+test)	93.1	43.9	89.0	85.9	69.6	84.6	86.9	81.7 (+2.3)

thoughts within these demonstrations strikingly improve LLMs performance in reasoning tasks (Wei et al., 2022; Zhang et al., 2022c; Kojima et al., 2022).

All of the aforementioned techniques for prompt optimization and efficient tuning of the language model use the original task's input text (or the original input context) provided within the dataset.

RL for Paraphrase Generation: In the following paragraphs, we provide a brief overview of similar reinforcement learning objectives employed for paraphrase generation. Li et al. (Li et al., 2018) used a deep RL technique, training a pointergenerator network as the paraphrase generator and a decomposable attention model as the evaluator which assigns a paraphrase score to pairs of sentences. The generator was trained using the policy gradient objective, with reward shaping and scaling to stabilize the training process (Li et al., 2018). Another approach by Qian et al. (Qian et al., 2019) focused on generating diverse paraphrases by training multiple generators, accompanied by a paraphrase discriminator and a generator discriminator. Policy gradient objective and self-critical learning (Rennie et al., 2016) were employed for training the generators, with the baseline reward used in the policy gradient objective being the reward obtained from the greedy-decoded sequence. Liu et al. (Liu et al., 2020) also applied the policy gradient objective with self-critical learning, incorporating multiple reward functions such as Rouge score with the reference paraphrase, negative Rouge score with the input sentence to encourage lexical variations, and semantic similarity score between the paraphrase and the input sentence to ensure semantic fidelity.

Another study by Du and Ji (Du and Ji, 2019) compared the use of imitation learning algorithm DAGGER with policy gradient REINFORCE for paraphrase generation. The policy gradient objective has also been applied in generating paraphrases while considering multiple objectives for entailment relation-aware paraphrase generation (Sancheti et al., 2022). In the context of chatbot responses, a recent work studies unsupervised paraphrase generation with proximal policy optimization, aiming to maximize a combination of rewards such as textual entailment, semantic similarity, language fluency, and lexical dissimilarity (Garg et al., 2021). Similarly, the policy gradient objective has been employed to optimize multiple rewards, similar to previous work, for unsupervised paraphrase generation (Siddique et al., 2020).

1301

1302

1303

1304

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

While previous studies have applied RL techniques for paraphrase generation, we propose the use of MML gradients instead of policy gradients to train our paraphrase model.

E Task Instructions & Input Format

Table 9 provides a summary of the task instructions that we append before the inputs, as well as the class verbalizers for classifying the input text. The instructions and input templates are derived from prior work in prompt optimization (Deng et al., 2022).

1298

1300

Table 9: Number of classes C, the input format	, and the instruction used pe	r dataset. The label	words are provided
within the instructions.			

Dataset	C	Input Format	Instruction
SST2	2	" <s> {Instruction}</s>	"In this task, you are given sentences from movie
		{Text} . It was	reviews. The task is to classify a sentence as 'great' if
		<mask> . "</mask>	the sentiment of the sentence is positive or as 'terrible'
			if the sentiment of the sentence is negative."
SST5	5	" <s> {Instruction}</s>	"In this task, you are given sentences from movie
		{Text} . It was	reviews. Based on the given review, classify it to one
		<mask> . "</mask>	of the five classes: (1) terrible, (2) bad, (3) okay, (4)
			good, and (5) great."
CR	2	" <s> {Instruction}</s>	"In this task, you are given sentences from customer
		{Text} . It was	reviews. The task is to classify a sentence as 'great' if
		<mask> . "</mask>	the sentiment of the sentence is positive or as 'terrible'
			if the sentiment of the sentence is negative."
MR	2	" <s> {Instruction}</s>	"In this task, you are given sentences from movie
		{Text} . It was	reviews. The task is to classify a sentence as 'great' if
		<mask> . "</mask>	the sentiment of the sentence is positive or as 'terrible'
			if the sentiment of the sentence is negative."
TREC	6	" <s> {Instruction}</s>	"You are given a question. You need to detect which
		<mask>: {Text} .</mask>	category better describes the question. Answer with
		"	'Description', 'Entity', 'Expression', 'Human', 'Lo-
			cation', and 'Number'."
Subj	2	" <s> {Instruction}</s>	"In this task, you are given sentences from reviews.
		{Text} . This is	The task is to classify a sentence as 'subjective' if the
		<mask> . "</mask>	opinion of the sentence is subjective or as 'objective'
			if the opinion of the sentence is objective."
AG's	4	" <s> {Instruction}</s>	"In this task, you are given a news article. Your task
News		<mask> News:</mask>	is to classify the article to one out of the four topics
		{Text} . "	'World', 'Sports', 'Business', 'Tech' if the article's
			main topic is relevant to the world, sports, business,
			and technology, correspondingly. If you are not sure
			about the topic, choose the closest option."