

# META-WEIGHTED LANGUAGE MODEL TUNING FOR AUGMENTATION-ENHANCED FEW-SHOT LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Recent studies have revealed the intriguing few-shot learning ability of pretrained language models (PLMs): They can quickly adapt to a new task when fine-tuned on a small amount of labeled data formulated as prompts, without requiring abundant task-specific annotations. Despite their promising performance, most existing few-shot approaches that only learn from the small training set still underperform fully supervised training by nontrivial margins. In this work, we study few-shot learning with PLMs from a different perspective: We first tune an autoregressive PLM on the few-shot samples and then use it as a generator to synthesize a large amount of novel training samples which augment the original training set. To encourage the generator to produce label-discriminative samples, we train it via weighted maximum likelihood where the weight of each token is automatically adjusted based on a discriminative meta-learning objective. A classification PLM can then be fine-tuned on both the few-shot and the synthetic samples with regularization for better generalization and stability. Our approach FewGen achieves an overall better result across seven classification tasks of the GLUE benchmark than existing few-shot learning methods, improving no-augmentation methods by 5+ average points, and outperforming augmentation methods by 3+ average points <sup>1</sup>.

## 1 INTRODUCTION

Recent research has demonstrated the appealing few-shot learning potential of pretrained language models (PLMs) (Brown et al., 2020; Clark et al., 2020; Devlin et al., 2019; He et al., 2021; Liu et al., 2019; Meng et al., 2021) on natural language understanding (NLU) tasks (Wang et al., 2019; 2018): Instead of relying on abundant task-specific annotations, PLMs can effectively leverage a small set of training samples to quickly learn a new task. Such training data efficiency is usually achieved by formulating downstream tasks as prompts (Brown et al., 2020; Gao et al., 2021; Scao & Rush, 2021; Schick & Schütze, 2021a;d) which allow the PLM to adapt its language modeling ability acquired through pretraining to new downstream tasks.

The success of prompt-based methods has stimulated numerous explorations along the line of effective few-shot learning with PLMs: The training samples converted to natural language prompts can be used to directly fine-tune PLMs (Gao et al., 2021; Schick & Schütze, 2021a) or as in-context demonstrations to facilitate better inference (Brown et al., 2020; Liu et al., 2022b). More recent approaches aim to automate the design of prompts by gradient-based searching (Shin et al., 2020) or parameterizing prompts as continuous learnable embeddings (Lester et al., 2021; Liu et al., 2021b; Zhang et al., 2022; Zhong et al., 2021). Other studies investigate and address specific issues in prompt-based few-shot learning (Liu et al., 2022a; Tam et al., 2021; Zhao et al., 2021). While remarkable, the model performance still has a nontrivial gap from fully supervised models trained on massive labeled data. Indeed, training deep models is inherently data demanding—model generalization usually benefits from more training samples (Baum & Haussler, 1988).

In this work, we study few-shot learning with PLMs from a different perspective: Instead of proposing new methods for fine-tuning on few-shot samples, we focus on the generation of quality training data based on few-shot samples and using these synthesized training samples to fine-tune the classification models. Motivated by the strong text generation power of autoregressive PLMs (Brown et al., 2020;

<sup>1</sup>Code is shared in the supplementary material.

Keskar et al., 2019; Raffel et al., 2019), previous data augmentation methods enlarge the training set by synthesizing new samples based on the few-shot samples. They either fine-tune the generator on the training set with the standard maximum likelihood objective (Anaby-Tavor et al., 2020; Kumar et al., 2020) or use the training samples as demonstrations (Yoo et al., 2021). However, these methods do not explicitly model the distinction across different labels and may struggle to generate accurate training samples pertaining to the desired labels for challenging NLU tasks.

In this paper, we study how to use few-shot samples to effectively tune PLMs to generate high quality label-discriminative training samples. Our contributions are as follows: (1) We analyze the issues of using standard maximum likelihood for tuning the generator and propose a meta-weighted maximum likelihood objective for generator tuning by automatically learning token weights that emphasize label discriminativeness. (2) We propose a simple and effective training procedure for fine-tuning classification PLMs on generated data by mitigating label noise. (3) Under the same few-shot learning setting, our method FewGen outperforms existing methods by 3+ average points on seven classification tasks of the GLUE benchmark (Wang et al., 2018). Ablation studies demonstrate the effectiveness of our proposed meta-weighted training objective and classifier fine-tuning method.

## 2 RELATED WORK

**Few-Shot Learning with PLMs.** Few-shot learning has gained much attention recently due to its minimal resource assumption—Without requiring massive annotated data but only leveraging a few training samples (*e.g.*, 16 per label), few-shot methods can be widely adopted in many practical scenarios where obtaining large-scale annotations is unaffordable. Standard fine-tuning of PLMs for few-shot learning usually performs poorly because the limited training samples may not be sufficient for optimizing the parameters in the newly introduced classification head. To reuse the language modeling ability of PLMs without introducing randomly initialized parameters, prompt-based approaches (Brown et al., 2020; Gao et al., 2021; Hu et al., 2022; Logan IV et al., 2021; Min et al., 2022; Schick & Schütze, 2021a;b;d; Tam et al., 2021) formulate training samples as natural language prompt templates so that various downstream tasks can be solved as a token prediction problem. They enjoy improved training data efficiency over standard fine-tuning in low-data regimes (Scao & Rush, 2021) and achieve remarkable few-shot learning performance. Later developments in prompt-based methods replace the manual design of prompt templates with automatic search or learning (Cui et al., 2022; Hambardzumyan et al., 2021; Lester et al., 2021; Liu et al., 2021b; Zhang et al., 2022; Zhong et al., 2021). There are also studies focusing on specific issues in prompt-based methods such as densifying the supervision by revising the training objective (Liu et al., 2022a; Tam et al., 2021) and calibrating the biased predictions of PLMs before fine-tuning (Zhao et al., 2021). Instead of focusing on fine-tuning methods for few-shot learning, we study how to effectively generate abundant quality training samples by learning from the few-shot samples and use them to improve the generalization of the classification model.

**Data Augmentation.** Data augmentation methods (Chen et al., 2020; Lee et al., 2021; Miyato et al., 2017; Xie et al., 2020) aim to create similar samples to the existing ones so that the enlarged training set can benefit model generalization. Early approaches simply use manually designed rules (*e.g.*, swapping or inserting tokens) for word-level alterations over the given samples to create new ones (Wei & Zou, 2019). Later methods leverage the strong generation power of PLMs to synthesize novel samples from scratch. Given a training set, the PLMs can be either fine-tuned on the labeled samples to learn label-conditioned generation probability (Kumar et al., 2020; Lee et al., 2021; Yang et al., 2020) or take the labeled data as demonstrations (Wang et al., 2021; Yoo et al., 2021) to generate similar samples pertaining to the same label. In this work, we study how to effectively tune generators on few-shot training data for creating new data—standard fine-tuning of PLMs on a small set of training data is prone to overfitting, and the resulting model may struggle to generate accurate, diverse and novel training data. We address this challenge by leveraging prefix-tuning and proposing a new meta-weighted training objective to emphasize label-discriminative tokens for generator tuning.

**Controlled Text Generation.** Generating training samples for different labels can be viewed as a form of controlled text generation (Hu et al., 2017), whose goal is to generate textual contents of desired semantics, styles or attributes. Such control can be realized through different stages of PLM training and deployment: During pretraining, control codes (Keskar et al., 2019) can be used as

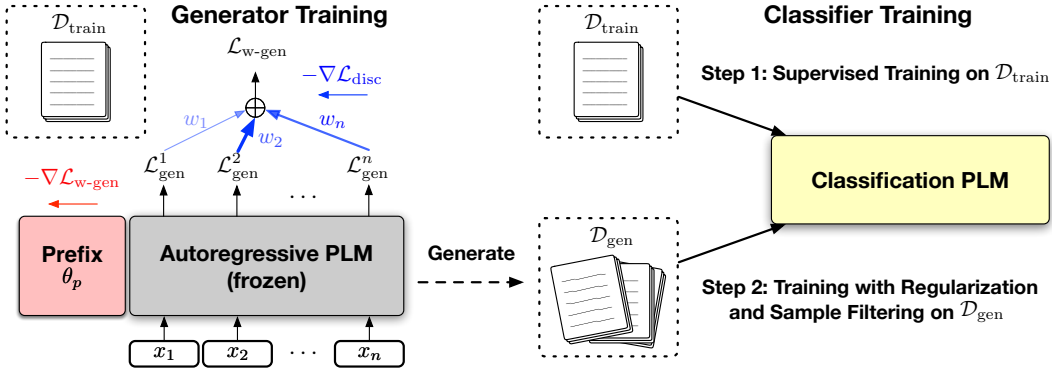


Figure 1: Overview of FewGen. A generator PLM is first tuned on the few-shot samples with our proposed meta-weighted maximum likelihood objective and then used to synthesize new training samples. A classification PLM is finally trained on both the few-shot and the generated samples.

explicit guidance for training the model to generate domain/attribute-specific texts; fine-tuning PLMs with attribute-specific data can also grant high-level control (*e.g.*, certain topics or sentiments (Ziegler et al., 2019)), fine-grained control (*e.g.*, specific words or phrases (Chan et al., 2021)) or both (Khalifa et al., 2021); at inference time, control over desired attributes can also be enforced without updating the PLM parameters (Dathathri et al., 2020; Krause et al., 2021; Kumar et al., 2021; Liu et al., 2021a; Pascual et al., 2021; Yang & Klein, 2021). Recently, a few studies explore fine-tuning autoregressive PLMs (Anaby-Tavor et al., 2020; Yang et al., 2020) with the standard language modeling objective on the training set or using label-specific prompts (Meng et al., 2022; Schick & Schütze, 2021c; Wang et al., 2021; Ye et al., 2022) to steer text generation towards the desired label.

**Meta-Learning for Sample Weighting.** The idea of weighting training samples in the loss calculation originates from the class imbalance (Wang et al., 2017) and noisy label (Hendrycks et al., 2018) learning scenarios—By assigning higher weights to the samples from minority classes or lower weights to the noisy samples, the learning process is less impacted by the imbalance/label noise issues. Meta-learning (Andrychowicz et al., 2016; Finn et al., 2017; Franceschi et al., 2018; Wu et al., 2018) is one way to automatically learn the weight for each sample. Specifically, a meta objective, usually defined as the loss on a clean unbiased validation set (Ren et al., 2018; Shu et al., 2019), can be used to learn the sample weights which become hyperparameters that control the optimization of model parameters. Our work has a different motivation and formulation of the meta objective for token-wise weighted training: Not all tokens in a training sample are equally label-discriminative. We thus design a meta objective to emphasize distinction across different labels (instead of using the validation loss as the meta objective) for learning the token weights.

### 3 METHOD

#### 3.1 PRELIMINARIES

**Overview.** We consider the strict few-shot learning setting (Perez et al., 2021): The training set  $\mathcal{D}_{\text{train}} = \{(\mathbf{x}, y)_i\}$  consists of  $K$  training samples per label where  $\mathbf{x} = [x_1, x_2, \dots, x_n]$  is a text sequence with  $n$  tokens. The development set  $\mathcal{D}_{\text{dev}}$  is of the same size as  $\mathcal{D}_{\text{train}}$ . There is no access to additional task-specific unlabeled data. The number of training samples  $K$  is assumed to be very small (*e.g.*,  $K = 16$ ), making it challenging to train a classification model  $C_\phi$  that generalizes well to unseen data. To mitigate such a training data scarcity issue, we propose to first train an autoregressive PLM on  $\mathcal{D}_{\text{train}}$ , and then use it as a generator  $G_\theta$  to synthesize a large amount of novel samples  $\mathcal{D}_{\text{gen}} = \{(\tilde{\mathbf{x}}, \tilde{y})_i\}$  that augment the original training set. Finally, a classification PLM  $C_\phi$  is fine-tuned on both  $\mathcal{D}_{\text{train}}$  and  $\mathcal{D}_{\text{gen}}$  to perform the task. An overview of our FewGen method is shown in Fig. 1.

**Text Generation with Autoregressive PLMs.** In standard fine-tuning for text generation, an autoregressive PLM  $G_\theta$  is trained via the maximum likelihood generation loss of each token in a

sequence  $\mathbf{x}$  conditioned on previous tokens:

$$\min_{\theta} -\frac{1}{n} \sum_{j=1}^n \log p_{\theta}(x_j | \mathbf{x}_{<j}), \quad p_{\theta}(x_j | \mathbf{x}_{<j}) = \frac{\exp(\mathbf{e}_j^{\top} \mathbf{h}_j)}{\sum_{j'=1}^{|\mathcal{V}|} \exp(\mathbf{e}_{j'}^{\top} \mathbf{h}_j)}.$$

where the token generation probability  $p_{\theta}(\cdot)$  is usually parameterized using token embeddings  $\mathbf{e}$  and hidden states  $\mathbf{h}$  of a Transformer (Vaswani et al., 2017) model. After training,  $G_{\theta}$  can be used to generate novel texts by iteratively sampling tokens from its generation probability distribution.

**Prefix-Tuning.** Unlike fine-tuning which updates all model parameters  $\theta$  of a PLM, prefix-tuning (Li & Liang, 2021) freezes all pretrained Transformer parameters and only optimizes prefix vectors  $\theta_p$  that are prepended to each Transformer layer. We use prefix-tuning for training  $G_{\theta_p}$  on  $\mathcal{D}_{\text{train}}$  because (1) it offers better effectiveness than fine-tuning for small datasets (Li & Liang, 2021) and (2) the generation models for different labels can share the same backbone Transformer parameters with only the prefix vectors being different, significantly reducing the memory requirement for multi-class classification tasks.

### 3.2 LABEL-DISCRIMINATIVE TEXT GENERATOR TUNING WITH META WEIGHTS

**Motivation.** To model the conditional text generation probability  $p(\mathbf{x}|y_l)$  on different labels, a straightforward way is to parameterize a generation model  $G_{\theta_{p_l}}$  for each label  $y_l$  via a set of prefix vectors  $\theta_p = \{\theta_{p_l}\}_{l=1}^L$  so that  $p(\mathbf{x}|y_l) = p_{\theta_{p_l}}(\mathbf{x})$ , and then tune  $\theta_{p_l}$  on the training samples  $\mathbf{x}$  with label  $y_l$ :

$$\min_{\theta_{p_l}} \mathcal{L}_{\text{gen}}, \quad \mathcal{L}_{\text{gen}}(\theta_{p_l}) = -\frac{1}{n} \sum_{j=1}^n \log p_{\theta_{p_l}}(x_j | \mathbf{x}_{<j}). \quad (1)$$

However, such an approach only optimizes the *generative* likelihood  $p(\mathbf{x}|y_l)$  without accounting for *label discriminativeness*  $p(y_l|\mathbf{x})$  which is essential for generating unambiguous training samples to benefit the final classification task. **Challenging NLU tasks can have largely similar distributions across different labels, with very nuanced differences reflected by a few key tokens.** For example, a negative review text “a movie where the ending feels like a cop-out” may immediately become a positive one by just changing the last word “cop-out” to “revelation”. Indeed, we find that such subtle distinctions over different labels may not be effectively captured by the generators if they are trained with the standard generation objective in Eq. (1). As shown in Fig. 2,  $\mathcal{L}_{\text{disc}}$  (defined in Eq. (2)) can even increase during training—It is possible that the dominating patterns in the training samples are label-indiscriminate (e.g., a movie review dataset may frequently mention “the movie”), making the generators of different labels eventually converge to similar distributions, especially when there are limited training samples per label.

To promote the generation of label-discriminative texts, we encourage each token  $x_j$  to be more likely generated under the corresponding label  $y_l$  instead of other labels (i.e., maximize  $p_{\theta_{p_l}}(x_j | \mathbf{x}_{<j})$  and minimize  $p_{\theta_{p_{l'}}}(x_j | \mathbf{x}_{<j})$  for  $l' \neq l$ ) via a discriminative loss  $\mathcal{L}_{\text{disc}}$ :

$$\mathcal{L}_{\text{disc}}(\theta_p) = -\frac{1}{n} \sum_{j=1}^n \mathcal{L}_{\text{disc}}^j(\theta_p), \quad \mathcal{L}_{\text{disc}}^j(\theta_p) = \frac{p_{\theta_{p_l}}(x_j | \mathbf{x}_{<j})}{\sum_{l'=1}^L p_{\theta_{p_{l'}}}(x_j | \mathbf{x}_{<j})} \quad (2)$$

Although one can directly combine  $\mathcal{L}_{\text{disc}}$  with  $\mathcal{L}_{\text{gen}}$  to train  $G_{\theta_p}$  to enforce distinction across different labels, doing so will result in two undesirable consequences: (1) A hyperparameter needs to be introduced to balance the weights of the two losses, whose optimal value is likely to vary by task; and (2) the generation-irrelevant loss  $\mathcal{L}_{\text{disc}}$  will unavoidably interfere the language modeling process, making the resulting model prone to generating less fluent and coherent texts.

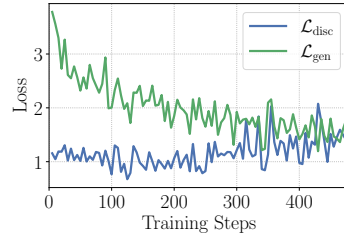


Figure 2: (On MNLI) Training the generator only via  $\mathcal{L}_{\text{gen}}$  does not automatically decrease  $\mathcal{L}_{\text{disc}}$ .

**Algorithm 1** Meta-Weighted Generator Tuning.**Input:**  $\mathcal{D}_{\text{train}}$ : Few-shot training set.**Parameter:**  $T$ : Number of training steps.**Output:**  $\theta_p$ : Prefix parameters for all labels.Initialize  $\theta_p^{(0)}$  (with task-descriptive prompts) and  $\omega^{(0)}$ **for**  $t \in [0, 1, \dots, T - 1]$  **do**     $\mathcal{B} \leftarrow$  Sample a minibatch from  $\mathcal{D}_{\text{train}}$      $\hat{\theta}_p^{(t)}(\omega^{(t)}) \leftarrow$  Take one gradient step to descend  $\mathcal{L}_{\text{w-gen}}(\theta_p^{(t)}; \omega^{(t)})$  on  $\mathcal{B}$      $\omega^{(t+1)} \leftarrow$  Take one gradient step to descend  $\mathcal{L}_{\text{disc}}(\hat{\theta}_p^{(t)}(\omega^{(t)}))$  on  $\mathcal{B}$      $\theta_p^{(t+1)} \leftarrow$  Take one gradient step to descend  $\mathcal{L}_{\text{w-gen}}(\theta_p^{(t)}; \omega^{(t+1)})$  on  $\mathcal{B}$ **end****return**  $\theta_p = \theta_p^{(T)}$ 

**Weighted Maximum Likelihood Generator Tuning.** To preserve the generative learning of  $G_{\theta_p}$  while emphasizing label-discriminative tokens, we assume each token is associated with a weight in the maximum likelihood loss. Intuitively, when our goal is to generate distinctive texts across different labels as training samples, not all tokens should contribute equally to generator training. For example, for sentiment classification tasks, one would expect “good/bad” to be more label-discriminative than “the movie”, and the former should be paid more attention to during training. It is thus natural to generalize  $\mathcal{L}_{\text{gen}}$  in Eq. (1) to  $\mathcal{L}_{\text{w-gen}}$  as follows by assuming a weight  $w_j$  is given for each token.

$$\min_{\theta_{p_l}} \mathcal{L}_{\text{w-gen}}, \quad \mathcal{L}_{\text{w-gen}}(\theta_{p_l}; \mathbf{w}) = - \sum_{j=1}^n w_j \mathcal{L}_{\text{gen}}^j(\theta_{p_l}), \quad \mathcal{L}_{\text{gen}}^j(\theta_{p_l}) = \log p_{\theta_{p_l}}(x_j | \mathbf{x}_{<j}). \quad (3)$$

Note that in  $\mathcal{L}_{\text{w-gen}}$ ,  $\mathbf{w}$  is assumed to be the *hyperparameter* under which  $\theta_{p_l}$  is optimized. When  $w_j$  is the same for every token, Eq. (3) will be equivalent to Eq. (1). While it is possible to manually design weighting rules for setting  $\mathbf{w}$  to promote label-discriminative learning, they will likely necessitate task-specific knowledge and nontrivial tuning. To facilitate the automatic learning of these weights  $\mathbf{w}$ , we propose to parameterize them as learnable hyperparameters using the idea of meta-learning.

**Meta Weight Learning Setup.** To automatically learn token weights using the idea of meta-learning, we formulate a bi-level optimization problem, where the inner objective  $\mathcal{L}_{\text{w-gen}}$  optimizes the generator parameters  $\theta_p$ , and the outer objective  $\mathcal{L}_{\text{disc}}$  optimizes token weights  $\mathbf{w}$  that are used as hyperparameters by the inner objective. We parameterize token weights  $\mathbf{w}$  via a weighting network  $g_{\omega}$  so that  $w_j = w_j(\omega)$ . Details about the implementation of  $g_{\omega}$  are in Appendix E. Overall, the learning objectives are as follows:

$$\begin{aligned} \theta_p^*(\omega) &= \underset{\theta_p}{\operatorname{argmin}} \mathcal{L}_{\text{w-gen}}, & \mathcal{L}_{\text{w-gen}}(\theta_p; \omega) &= - \sum_{j=1}^n w_j(\omega) \mathcal{L}_{\text{gen}}^j(\theta_p) \\ \omega^* &= \underset{\omega}{\operatorname{argmin}} \mathcal{L}_{\text{disc}}, & \mathcal{L}_{\text{disc}}(\theta_p^*(\omega)) &= - \frac{1}{n} \sum_{j=1}^n \mathcal{L}_{\text{disc}}^j(\theta_p^*(\omega)) \end{aligned} \quad (4)$$

Under the above formulation, the token weights  $w_j(\omega)$  are automatically learned such that the resulting generator parameters  $\theta_p^*(\omega)$  capture label-discriminative information (*i.e.*, minimize  $\mathcal{L}_{\text{disc}}$ ). Instead of solving the optimal  $\omega^*$  and  $\theta_p^*$  via nested optimization loops, we use an online optimization strategy (Shu et al., 2019) for training efficiency. It also guarantees convergence to the critical points of both  $\mathcal{L}_{\text{w-gen}}$  and  $\mathcal{L}_{\text{disc}}$  under mild conditions. The initialization prompts can be found in Appendix C. The overall training procedure is shown in Algorithm 1.

**Analysis of Meta Weight Learning.** We analyze the gradient update of meta weights to study its effect in generator tuning. The weighting network parameter  $\omega$  is optimized via Eq. (4), and its

---

**Algorithm 2** Classification model fine-tuning on  $\mathcal{D}_{\text{train}}$  and  $\mathcal{D}_{\text{gen}}$ .

---

**Input:**  $\mathcal{D}_{\text{train}}$ : Few-shot training set;  $\mathcal{D}_{\text{gen}}$ : Synthesized training set.

**Parameter:**  $T$ : Number of training steps.

**Output:**  $\phi$ : Trained classification model parameters.

$\phi^{(0)} \leftarrow$  Train on  $\mathcal{D}_{\text{train}}$  with standard supervised learning

$\bar{z} \leftarrow \mathbf{0}$  // Ensembled prediction initialization

**for**  $t \in [0, 1, \dots, T - 1]$  **do**

$\mathcal{B} \leftarrow$  Sample a minibatch from  $\mathcal{D}_{\text{gen}}$

$\phi^{(t+1)} \leftarrow$  Take one gradient step to descend  $\mathcal{L}_{\text{class}}$  in Eq. (5) on  $\mathcal{B}$

$\bar{z} \leftarrow$  Accumulate the current model prediction

    Update  $\mathcal{D}_{\text{gen}}$  to exclude noisy samples based on  $\bar{z}$

**end**

**return**  $\phi = \phi^{(T)}$

---

gradient is as follows (detailed derivation in Appendix A):

$$-\frac{\partial \mathcal{L}_{\text{disc}}(\hat{\theta}_p^{(t)}(\omega))}{\partial \omega} \bigg|_{\omega=\omega^{(t)}} \propto \sum_{j=1}^n d_j \frac{\partial w_j(\omega)}{\partial \omega} \bigg|_{\omega=\omega^{(t)}}, \quad d_j = \frac{\partial \mathcal{L}_{\text{disc}}(\hat{\theta}_p)}{\partial \hat{\theta}_p} \bigg|_{\hat{\theta}_p=\hat{\theta}_p^{(t)}} \frac{\partial \mathcal{L}_{\text{gen}}^j(\theta_p)}{\partial \theta_p} \bigg|_{\theta_p=\theta_p^{(t)}}^\top.$$

It can be seen that the gradient descent direction of  $\omega$  is determined by a weighted sum of token weight gradient ascent direction (i.e.,  $\frac{\partial w_j(\omega)}{\partial \omega}$ ), where the weight  $d_j$  characterizes the similarity between the gradient of the discriminative objective and the gradient of the generative objective on the  $j$ th token. Therefore, the meta weights will be higher on those tokens where optimizing their generative objective is more beneficial for minimizing the discriminative objective.

### 3.3 CLASSIFIER FINE-TUNING

With the trained generator  $G_{\theta_p}$ , we can synthesize novel training samples  $\mathcal{D}_{\text{gen}}$  that augment  $\mathcal{D}_{\text{train}}$  for fine-tuning a PLM  $C_\phi$  for classification. The major challenge to effectively leverage  $\mathcal{D}_{\text{gen}}$  is that the label noise (i.e., some generated samples may not accurately pertain to the corresponding label) may deteriorate model performance if standard supervised learning is directly used. We propose a simple noise-robust training procedure to improve the generalization and stability of training: First fine-tune  $C_\phi$  on  $\mathcal{D}_{\text{train}}$  with standard supervised training, and then continue fine-tuning it on  $\mathcal{D}_{\text{gen}}$  by applying *temporal ensembling* (Laine & Aila, 2017) as regularization. Specifically, given a training sample  $(\tilde{x}, \tilde{y}) \in \mathcal{D}_{\text{gen}}$ , we minimize the following classification loss:

$$\min_{\phi} \mathcal{L}_{\text{class}}, \quad \mathcal{L}_{\text{class}}(\phi) = -\log(p_\phi(\tilde{x})_{\tilde{y}}) - \lambda \sum_{l=1}^L \bar{z}_l \log \frac{p_\phi(\tilde{x})_l}{\bar{z}_l}, \quad (5)$$

where  $p_\phi(\tilde{x})$  is the model prediction on  $\tilde{x}$ ;  $\lambda$  is a regularization weight for temporal ensembling; and  $\bar{z}$  is the accumulated moving-average model predictions. We also use the ensembled prediction  $\bar{z}$  to filter out noisy synthesized samples: We only include those samples for training where  $\bar{z}$  strongly agrees with the label  $\tilde{y}$  (i.e.,  $\bar{z}_{\tilde{y}} > \delta$  where  $\delta > 0$  is a threshold parameter). In Eq. (5), the first classification term is the cross-entropy loss; the second regularization term corresponds to temporal ensembling, which requires the current model prediction to be close to its past accumulated predictions. This not only neutralizes the fluctuation in model predictions for better training stability when label noise is present (Nguyen et al., 2020) but also helps prevent catastrophic forgetting (Kirkpatrick et al., 2017) of the information learned previously from the few-shot training set  $\mathcal{D}_{\text{train}}$ . Please refer to Appendix C for details about the temporal ensembling implementation. The overall procedure of classifier fine-tuning is summarized in Algorithm 2.

## 4 EXPERIMENTAL SETUP

**Downstream Tasks and Metrics.** We conduct evaluation on all tasks of the GLUE benchmark (Wang et al., 2018) (more details in Appendix B) except STS-B which is a regression task. We

Table 1: Results on seven classification tasks of the GLUE benchmark. We report average and standard deviation (as subscripts) performance over 5 different  $\mathcal{D}_{\text{train}}/\mathcal{D}_{\text{dev}}$  splits defined in Gao et al. (2021). <sup>†</sup>: Results from Gao et al. (2021). <sup>‡</sup>: Results from Zhang et al. (2022). Methods that use additional models apart from the final classification model are marked.

| Method   | MNLI-(m/mm)<br>(Acc.)                                   | QQP<br>(F1)                | QNLI<br>(Acc.)             | SST-2<br>(Acc.)            | CoLA<br>(Matt.)            | RTE<br>(Acc.)              | MRPC<br>(F1)               | AVG         |
|--|---|----------------------------|----------------------------|----------------------------|----------------------------|----------------------------|----------------------------|-------------|
| <i>Methods without Augmentation: Few-shot samples are directly used for classifier tuning or as demonstrations for inference</i> |   |                            |                            |                            |                            |                            |                            |             |
| Prompting <sup>†</sup>   | 50.8/51.7   | 49.7                       | 50.8                       | 83.6                       | 2.0                        | 51.3                       | 61.9                       | 50.1        |
| Fine-Tuning <sup>†</sup>   | 45.8 <sub>6.4</sub> /47.8 <sub>6.8</sub>                | 60.7 <sub>4.3</sub>        | 60.2 <sub>6.5</sub>        | 81.4 <sub>3.8</sub>        | 33.9 <sub>14.3</sub>       | 54.4 <sub>3.9</sub>        | 76.6 <sub>2.5</sub>        | 59.1        |
| In-Context <sup>†</sup>  | 52.0 <sub>0.7</sub> /53.4 <sub>0.6</sub>                | 36.1 <sub>5.2</sub>        | 53.8 <sub>0.4</sub>        | 84.8 <sub>1.3</sub>        | -1.5 <sub>2.4</sub>        | 60.4 <sub>1.4</sub>        | 45.7 <sub>6.0</sub>        | 47.4        |
| LM-BFF (Man.) <sup>†</sup>   | 68.3 <sub>2.3</sub> /70.5 <sub>1.9</sub>                | 65.5 <sub>5.3</sub>        | 64.5 <sub>4.2</sub>        | 92.7 <sub>0.9</sub>        | 9.3 <sub>7.3</sub>         | 69.1 <sub>3.6</sub>        | 74.5 <sub>5.3</sub>        | 63.6        |
| + demonstration <sup>†</sup>   | 70.7 <sub>1.3</sub> /72.0 <sub>1.2</sub>                | 69.8 <sub>1.8</sub>        | 69.2 <sub>1.9</sub>        | 92.6 <sub>0.5</sub>        | 18.7 <sub>8.8</sub>        | 68.7 <sub>2.3</sub>        | 77.8 <sub>2.0</sub>        | 66.9        |
| LM-BFF (Auto) <sup>†</sup> (w. 2.9B T5)  | 68.3 <sub>2.5</sub> /70.1 <sub>2.6</sub>                | 67.0 <sub>3.0</sub>        | 68.3 <sub>7.4</sub>        | 92.3 <sub>1.0</sub>        | 14.0 <sub>14.1</sub>       | <b>73.9</b> <sub>2.2</sub> | 76.2 <sub>2.3</sub>        | 65.8        |
| + demonstration <sup>†</sup> (w. 2.9B T5)  | 70.0 <sub>3.6</sub> /72.0 <sub>3.1</sub>                | 67.7 <sub>5.8</sub>        | 68.5 <sub>5.4</sub>        | 93.0 <sub>0.6</sub>        | 21.8 <sub>15.9</sub>       | 71.1 <sub>5.3</sub>        | 78.1 <sub>3.4</sub>        | 67.3        |
| P-Tuning <sup>‡</sup>  | 61.5 <sub>2.1</sub> /-                                  | 65.6 <sub>3.0</sub>        | 64.3 <sub>2.8</sub>        | 92.2 <sub>0.4</sub>        | -                          | -                          | 74.5 <sub>7.6</sub>        | -           |
| DART <sup>‡</sup>  | 67.5 <sub>2.6</sub> /-                                  | 67.8 <sub>3.2</sub>        | 66.7 <sub>3.7</sub>        | 93.5 <sub>0.5</sub>        | -                          | -                          | 78.3 <sub>4.5</sub>        | -           |
| <i>Methods with Augmentation: Few-shot samples are used for creating synthesized samples and for classifier tuning</i>           |   |                            |                            |                            |                            |                            |                            |             |
| MixText  | 65.1 <sub>2.6</sub> /66.2 <sub>2.8</sub>                | 60.6 <sub>3.9</sub>        | 68.4 <sub>5.1</sub>        | 89.1 <sub>2.3</sub>        | 12.8 <sub>9.2</sub>        | 66.5 <sub>4.1</sub>        | 64.6 <sub>7.6</sub>        | 61.1        |
| Back Translation (w. trained Marian)   | 66.9 <sub>4.6</sub> /68.3 <sub>3.8</sub>                | 59.8 <sub>4.6</sub>        | 67.8 <sub>4.9</sub>        | 91.1 <sub>1.9</sub>        | 7.5 <sub>3.7</sub>         | 62.4 <sub>5.3</sub>        | 68.0 <sub>11.2</sub>       | 60.6        |
| GPT3Mix (w. 175B GPT3)   | 61.5 <sub>3.2</sub> /62.6 <sub>2.2</sub>                | 70.4 <sub>1.9</sub>        | 69.2 <sub>0.3</sub>        | <b>93.6</b> <sub>0.6</sub> | <b>48.9</b> <sub>1.9</sub> | 70.4 <sub>10.0</sub>       | 69.9 <sub>12.4</sub>       | 69.2        |
| Generator Fine-Tuning (w. 1.6B CTRL)   | 68.9 <sub>5.1</sub> /70.8 <sub>5.3</sub>                | 60.4 <sub>8.7</sub>        | 70.9 <sub>4.1</sub>        | 91.2 <sub>1.2</sub>        | 18.8 <sub>10.0</sub>       | 66.1 <sub>4.4</sub>        | 60.8 <sub>15.4</sub>       | 62.6        |
| FewGen (w. 1.6B CTRL)  | <b>75.7</b> <sub>1.6</sub> / <b>77.1</b> <sub>1.0</sub> | <b>71.5</b> <sub>1.7</sub> | <b>76.3</b> <sub>4.4</sub> | 93.1 <sub>0.8</sub>        | 40.0 <sub>7.5</sub>        | 71.2 <sub>2.4</sub>        | <b>81.1</b> <sub>2.5</sub> | <b>72.8</b> |
| Fully Supervised Fine-Tuning <sup>†</sup>  | 89.8/89.5   | 81.7                       | 93.3                       | 95.0                       | 62.6                       | 80.9                       | 91.4                       | 84.9        |

follow the same data split and evaluation protocol as Gao et al. (2021): Both  $\mathcal{D}_{\text{train}}$  and  $\mathcal{D}_{\text{dev}}$  contain 16 samples per label and are sampled from the original training set with 5 different random seeds. The original development sets are used for testing. For all reported results, we include the average and standard deviation over the 5 different  $\mathcal{D}_{\text{train}}/\mathcal{D}_{\text{dev}}$  splits. F1 score is used as the metric for QQP and MRPC, Matthews correlation for CoLA, and accuracy for the remaining tasks.

**Models, Training Settings and Hyperparameters.** FewGen is a training data generation method and can be used with any fine-tuning method on any classification model. We use moderate-sized PLMs to ensure our results are reproducible on typical research hardware: CTRL (1.6B parameters) (Keskar et al., 2019) as the generator  $G_\theta$  and RoBERTa<sub>Large</sub> (356M parameters) (Liu et al., 2019) as the classifier  $C_\phi$ . We use prefix-tuning for training  $G_\theta$  and prompt-based fine-tuning for training  $C_\phi$ . For simplicity, we use the most basic manual prompt version of LM-BFF (Gao et al., 2021). The only exception is CoLA for which we use the standard fine-tuning since the input data might be out of the distribution of  $C_\phi$  (Gao et al., 2021). The hyperparameter tuning is performed on  $\mathcal{D}_{\text{dev}}$ . More details are in Appendix C.

**Compared Methods.** No-augmentation baselines include zero-shot prompting, standard fine-tuning, in-context learning, and the following strong few-shot learning methods: Four versions of LM-BFF (Gao et al., 2021), P-Tuning (Liu et al., 2021b) and DART (Zhang et al., 2022). We also compare FewGen with data augmentation methods for few-shot learning: MixText (Chen et al., 2020), using back translation systems to generate paraphrases (UDA-style (Xie et al., 2020) augmentation), GPT3Mix (Yoo et al., 2021) and standard fine-tuning of generator on the few-shot samples with prompts. All augmentation methods use LM-BFF (Man.) for fine-tuning the RoBERTa<sub>Large</sub> classifier. More details about data augmentation baselines can be found in Appendix D.

## 5 EVALUATION

### 5.1 MAIN RESULTS

We present the results of FewGen and baselines in Table 1. FewGen achieves overall better performance across the GLUE tasks, on average 5+ points higher than the previous best few-shot method without augmentation, and 3+ points better than GPT3Mix<sup>2</sup> (Yoo et al., 2021) which uses a 100 times larger generator model (175B) than FewGen. The promising results confirm the effectiveness of our

<sup>2</sup>The CoLA results reported in the original GPT3Mix paper use accuracy as the metric instead of Matthews correlation; our reimplemented GPT3Mix achieves 79.4<sub>0.6</sub> on CoLA if measured by accuracy.



Table 2: Ablation studies by removing (–) or switching (w.) one component of FewGen.

| Method   | MNLI-(m/mm)                              | QQP                 | QNLI                | SST-2               | CoLA                | RTE                 | MRPC                |
|--|--|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|
| FewGen   | 75.7 <sub>1.6</sub> /77.1 <sub>1.0</sub> | 71.5 <sub>1.7</sub> | 76.3 <sub>4.4</sub> | 93.1 <sub>0.8</sub> | 40.0 <sub>7.5</sub> | 71.2 <sub>2.4</sub> | 81.1 <sub>2.5</sub> |
| w. $\mathcal{L}_{\text{gen}}$  | 74.9 <sub>1.0</sub> /76.2 <sub>1.0</sub> | 70.7 <sub>1.9</sub> | 75.0 <sub>4.8</sub> | 92.5 <sub>0.7</sub> | 37.8 <sub>8.2</sub> | 69.5 <sub>2.2</sub> | 80.8 <sub>3.0</sub> |
| w. $\mathcal{L}_{\text{gen}} + \mathcal{L}_{\text{disc}}$                  | 74.6 <sub>1.6</sub> /76.0 <sub>1.5</sub> | 68.8 <sub>2.1</sub> | 76.1 <sub>4.3</sub> | 92.4 <sub>0.8</sub> | 41.2 <sub>9.0</sub> | 70.1 <sub>2.2</sub> | 79.6 <sub>2.4</sub> |
| – temporal ensemble  | 72.2 <sub>2.5</sub> /74.0 <sub>2.2</sub> | 65.8 <sub>2.1</sub> | 75.1 <sub>2.7</sub> | 92.1 <sub>1.7</sub> | 33.9 <sub>4.4</sub> | 66.6 <sub>2.4</sub> | 80.4 <sub>3.2</sub> |
| w. fine-tune on $\mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{gen}}$ | 68.9 <sub>1.8</sub> /70.6 <sub>1.9</sub> | 64.3 <sub>1.5</sub> | 71.1 <sub>4.1</sub> | 91.8 <sub>1.3</sub> | 34.0 <sub>3.2</sub> | 59.6 <sub>1.0</sub> | 80.4 <sub>3.5</sub> |

proposed FewGen method in generating quality training data and leveraging them in combination with the few-shot training set for fine-tuning the classification model.

**Comparison with Back Translation.** Using back translation to paraphrase the few-shot samples does not improve the results, even with prompt-based fine-tuning to train the classifier – this is probably because it does not produce samples that are sufficiently different from the few-shot training set. The success of UDA (Xie et al., 2020) is grounded in the augmentations from abundant unlabeled data that improve the classifier generalization. However, under the strict few-shot learning setup, there is no access to additional task-specific unlabeled data (Gao et al., 2021), making it challenging for paraphrase-based methods to create sufficiently diverse training samples only based on the small few-shot set. The new training samples produced by our FewGen method are not limited to the paraphrases of the few-shot samples, as the generator is trained via prefix-tuning to preserve the PLM’s pretraining knowledge, based on which novel training samples can be synthesized.

**Comparison with GPT3Mix.** The gigantic size of GPT3 makes it challenging for tuning on few-shot samples. Therefore, GPT3Mix (Yoo et al., 2021) uses few-shot samples as demonstrations for creating the augmentations. Such an approach suffers from two limitations: (1) Without any parameter update to the PLM, its learning ability is not fully leveraged to adapt to the few-shot training set. (2) The PLM can only use a small subset of the few-shot samples at a time for creating each augmentation, as the number of demonstrations received by the model is bounded by its maximum input sequence length. This makes the quality of the created augmentations more sensitive to the randomly drawn training samples. Our FewGen method, on the other hand, can use the entire few-shot set for tuning the PLM and achieves overall even better classification results with a much smaller PLM (< 1% the size of the GPT3 model) which can be deployed much more easily in practice.

## 5.2 ABLATION STUDIES

We further analyze the effectiveness of each important component in FewGen. Specifically, we compare FewGen with the following ablations: (1) Using the standard  $\mathcal{L}_{\text{gen}}$  in Eq. (1) instead of our proposed  $\mathcal{L}_{\text{w-gen}}$  in Eq. (3) for generator tuning (w.  $\mathcal{L}_{\text{gen}}$ ); (2) using the directly combined  $\mathcal{L}_{\text{gen}}$  and  $\mathcal{L}_{\text{disc}}$  for generator tuning (w.  $\mathcal{L}_{\text{gen}} + \mathcal{L}_{\text{disc}}$ ); (3) without applying temporal ensembling in Eq. (5) (– temporal ensemble); (4) directly fine-tuning the classification model on the combination of  $\mathcal{D}_{\text{gen}}$  and  $\mathcal{D}_{\text{train}}$  (w. fine-tune on  $\mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{gen}}$ )<sup>3</sup>. As shown in Table 2, (1) & (2) using the standard maximum likelihood loss or the combination of generation and discrimination losses to tune the generator both yield lower-quality training data and lead to degraded classification performance; (3) not applying temporal ensembling for fine-tuning the classifier is more prone to label noise in the generated samples; (4) fine-tuning the classifier on the combination of  $\mathcal{D}_{\text{gen}}$  and  $\mathcal{D}_{\text{train}}$  significantly underperforms our two-step fine-tuning method. To study the impact of the amount of generated training samples on the model performance, we plot the MNLI-m accuracy (mean and standard deviation) with different sizes of  $\mathcal{D}_{\text{gen}}$  in Fig. 3. Both the average model performance and stability improve with more generated samples.

## 5.3 ANALYSES OF LOSS FUNCTIONS FOR GENERATOR TUNING

As shown in Table 2, the choice of generator loss has a significant impact on the synthesized data quality and thus the final model performance. We conduct further analyses to compare the training processes of the generator under the following three loss functions and the resulting generated

<sup>3</sup>For this ablation, we upsample  $\mathcal{D}_{\text{train}}$  by  $\times 100$  so that its size is comparable with  $\mathcal{D}_{\text{gen}}$ ; without upsampling, the result is much worse.



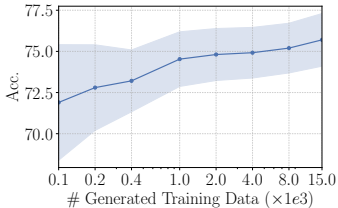


Figure 3: MNLi-m accuracy with different amounts of generated training data.

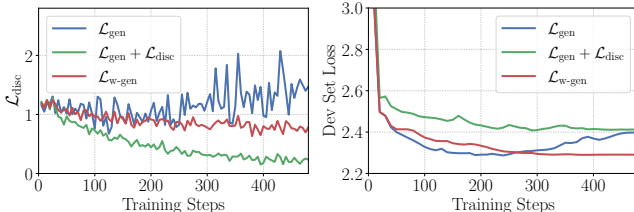


Figure 4: With different loss functions used for generator tuning, (Left)  $\mathcal{L}_{\text{disc}}$  and (Right) standard language modeling loss on the development set. Best viewed in color.

samples: (1)  $\mathcal{L}_{\text{gen}}$  which is the standard language modeling loss; (2)  $\mathcal{L}_{\text{gen}} + \mathcal{L}_{\text{disc}}$  which directly adds the discriminative loss to generator training; and (3)  $\mathcal{L}_{\text{w-gen}}$  which is our meta-weighted objective. Fig. 4 shows the discriminative loss  $\mathcal{L}_{\text{disc}}$  and the standard language modeling loss on the held-out development set throughout training. Although using  $\mathcal{L}_{\text{gen}} + \mathcal{L}_{\text{disc}}$  helps reduce the discriminative loss, it comes at the cost of hindering language modeling—the generator loss on the development set is high. Using our meta-weighted objective  $\mathcal{L}_{\text{w-gen}}$  for tuning the generator not only encourages discriminativeness but also mitigates overfitting, yielding the lowest validation set loss. This is probably because the model receives contrastive information from other labels which facilitates more accurate modeling of the texts with the target label. We present more quantitative analyses of different generator training objectives in Appendix G. We visualize the token weights  $w$  automatically learned and used in  $\mathcal{L}_{\text{w-gen}}$  in Appendix F.

## 6 DISCUSSIONS AND CONCLUSIONS

**Ethical Considerations.** Despite the impressive text generation and representation power of PLMs, they can also come with the risk (Bender et al., 2021; Bender & Koller, 2020; Brown et al., 2020) of generating disinformation (Pagnoni et al., 2021) or exacerbating biases (Prabhumoye et al., 2018). Instead of improving upon PLM architectures or generation techniques, our work focuses on using existing PLMs to create training data for NLU tasks. Therefore, our method can be combined with any bias reduction and correction strategies (Gehman et al., 2020; Ma et al., 2020) in practice to reduce the adverse effects of PLMs.

**Limitations.** Compared to few-shot learning methods that directly train classification models on the small training set, FewGen requires tuning a generator PLM and using it to synthesize novel training samples, resulting in higher computation costs and longer running time. Still, we believe that our method may bring more good than harm—when the small training data size becomes the performance bottleneck for NLU tasks, a simple yet costly solution is to obtain more human annotations. Our method may replace or reduce the human efforts in such training data creation processes.

**Conclusions.** In this work, we propose FewGen, which leverages few-shot training samples to tune a generator PLM for synthesizing novel training data. The generated data can be then used in combination with few-shot samples to fine-tune a classification model for better generalization. To emphasize label-discriminative information during generator tuning, we propose a weighted maximum likelihood objective where the token weights are automatically learned via a discriminative meta objective. Since the generated samples may contain label noise, we propose a simple training procedure that first trains classifiers on the few-shot training set and then on the generated set by applying temporal ensembling for noise-robustness. Across seven classification tasks from the GLUE benchmark, FewGen significantly outperforms existing approaches under the same few-shot learning setting. The effectiveness of each important component in FewGen is validated via ablation studies. Future work directions may include: Using larger PLMs as the generator and the classifier, jointly training both models with each other’s high-confident predictions, and developing systematic metrics for evaluating the quality of generated training samples.

## REFERENCES

- Ateret Anaby-Tavor, Boaz Carmeli, Esther Goldbraich, Amir Kantor, George Kour, Segev Shlomov, N. Tepper, and Naama Zwerdling. Do not have enough data? deep learning to the rescue! In *AAAI*, 2020.
- Marcin Andrychowicz, Misha Denil, Sergio Gomez Colmenarejo, Matthew W. Hoffman, David Pfau, Tom Schaul, and Nando de Freitas. Learning to learn by gradient descent by gradient descent. In *NIPS*, 2016.
- Eric Baum and David Haussler. What size net gives valid generalization? In *NIPS*, 1988.
- Emily M. Bender and Alexander Koller. Climbing towards NLU: On meaning, form, and understanding in the age of data. In *ACL*, 2020.
- Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big? In *ACM Conference on Fairness, Accountability, and Transparency*, 2021.
- Luisa Bentivogli, Peter Clark, Ido Dagan, and Danilo Giampiccolo. The fifth pascal recognizing textual entailment challenge. In *TAC*, 2009.
- 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, T. J. Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeff 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. Language models are few-shot learners. In *NeurIPS*, 2020.
- Alvin Chan, Y. Ong, Bill Tuck Weng Pung, Aston Zhang, and Jie Fu. CoCon: A self-supervised approach for controlled text generation. In *ICLR*, 2021.
- Jiaao Chen, Zichao Yang, and Diyi Yang. MixText: Linguistically-informed interpolation of hidden space for semi-supervised text classification. In *ACL*, 2020.
- Kevin Clark, Minh-Thang Luong, Quoc V Le, and Christopher D Manning. ELECTRA: Pre-training text encoders as discriminators rather than generators. In *ICLR*, 2020.
- Ganqu Cui, Shengding Hu, Ning Ding, Longtao Huang, and Zhiyuan Liu. Prototypical verbalizer for prompt-based few-shot tuning. In *ACL*, 2022.
- Ido Dagan, Oren Glickman, and Bernardo Magnini. The pascal recognising textual entailment challenge. In *Machine Learning Challenges Workshop*, 2005.
- Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. Plug and play language models: A simple approach to controlled text generation. In *ICLR*, 2020.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2019.
- William B Dolan and Chris Brockett. Automatically constructing a corpus of sentential paraphrases. In *International Workshop on Paraphrasing (IWP)*, 2005.
- Chelsea Finn, P. Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, 2017.
- Luca Franceschi, Paolo Frasconi, Saverio Salzo, Riccardo Grazi, and Massimiliano Pontil. Bilevel programming for hyperparameter optimization and meta-learning. In *ICML*, 2018.
- Tianyu Gao, Adam Fisch, and Danqi Chen. Making pre-trained language models better few-shot learners. In *ACL*, 2021.

- Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. RealToxicityPrompts: Evaluating neural toxic degeneration in language models. In *EMNLP Findings*, 2020.
- Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and Bill Dolan. The third pascal recognizing textual entailment challenge. In *ACL-PASCAL workshop on textual entailment and paraphrasing*, 2007.
- R Bar Haim, Ido Dagan, Bill Dolan, Lisa Ferro, Danilo Giampiccolo, Bernardo Magnini, and Idan Szpektor. The second pascal recognising textual entailment challenge. In *PASCAL Challenges Workshop on Recognising Textual Entailment*, 2006.
- Karen Hambardzumyan, H. Khachatrian, and Jonathan May. WARP: Word-level adversarial reprogramming. In *ACL*, 2021.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. DeBERTa: Decoding-enhanced BERT with disentangled attention. In *ICLR*, 2021.
- Dan Hendrycks, Mantas Mazeika, Duncan Wilson, and Kevin Gimpel. Using trusted data to train deep networks on labels corrupted by severe noise. In *NeurIPS*, 2018.
- Shengding Hu, Ning Ding, Huadong Wang, Zhiyuan Liu, Juan-Zi Li, and Maosong Sun. Knowledgeable prompt-tuning: Incorporating knowledge into prompt verbalizer for text classification. In *ACL*, 2022.
- Zhiting Hu, Zichao Yang, Xiaodan Liang, Ruslan Salakhutdinov, and Eric P. Xing. Toward controlled generation of text. In *ICML*, 2017.
- Marcin Junczys-Dowmunt, Roman Grundkiewicz, Tomasz Dwojak, Hieu T. Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F. T. Martins, and Alexandra Birch. Marian: Fast neural machine translation in C++. In *ACL System Demo*, 2018.
- Nitish Shirish Keskar, Bryan McCann, Lav R. Varshney, Caiming Xiong, and Richard Socher. CTRL: A conditional transformer language model for controllable generation. *ArXiv*, abs/1909.05858, 2019.
- Muhammad Khalifa, Hady ElSahar, and Marc Dymetman. A distributional approach to controlled text generation. In *ICLR*, 2021.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 2017.
- Ben Krause, Akhilesh Deepak Gotmare, Bryan McCann, Nitish Shirish Keskar, Shafiq R. Joty, Richard Socher, and Nazneen Rajani. GeDi: Generative discriminator guided sequence generation. In *EMNLP*, 2021.
- Sachin Kumar, Eric Malmi, Aliaksei Severyn, and Yulia Tsvetkov. Controlled text generation as continuous optimization with multiple constraints. In *NeurIPS*, 2021.
- Varun Kumar, Ashutosh Choudhary, and Eunah Cho. Data augmentation using pre-trained transformer models. In *Workshop on Life-long Learning for Spoken Language Systems*, 2020.
- Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. In *ICLR*, 2017.
- Kenton Lee, Kelvin Guu, Luheng He, Tim Dozat, and Hyung Won Chung. Neural data augmentation via example extrapolation. *arXiv preprint arXiv:2102.01335*, 2021.
- Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. In *EMNLP*, 2021.
- Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In *ACL*, 2021.

- Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith, and Yejin Choi. DExperts: Decoding-time controlled text generation with experts and anti-experts. In *ACL*, 2021a.
- Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and Colin Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. In *NeurIPS*, 2022a.
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. What makes good in-context examples for GPT-3? In *Proceedings of Deep Learning Inside Out*, 2022b.
- Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. GPT understands, too. *ArXiv*, abs/2103.10385, 2021b.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A robustly optimized BERT pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- Robert L Logan IV, Ivana Balažević, Eric Wallace, Fabio Petroni, Sameer Singh, and Sebastian Riedel. Cutting down on prompts and parameters: Simple few-shot learning with language models. *arXiv preprint arXiv:2106.13353*, 2021.
- Xinyao Ma, Maarten Sap, Hannah Rashkin, and Yejin Choi. PowerTransformer: Unsupervised controllable revision for biased language correction. In *EMNLP*, 2020.
- Yu Meng, Chenyan Xiong, Payal Bajaj, Saurabh Tiwary, Paul Bennett, Jiawei Han, and Xia Song. COCO-LM: Correcting and contrasting text sequences for language model pretraining. In *NeurIPS*, 2021.
- Yu Meng, Jiaxin Huang, Yu Zhang, and Jiawei Han. Generating training data with language models: Towards zero-shot language understanding. In *NeurIPS*, 2022.
- Sewon Min, Michael Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. Noisy channel language model prompting for few-shot text classification. In *ACL*, 2022.
- Takeru Miyato, Andrew M. Dai, and Ian J. Goodfellow. Adversarial training methods for semi-supervised text classification. In *ICLR*, 2017.
- Duc Tam Nguyen, Chaithanya Kumar Mummadi, Thi-Phuong-Nhung Ngo, Thi Hoai Phuong Nguyen, Laura Beggel, and Thomas Brox. SELF: Learning to filter noisy labels with self-ensembling. In *ICLR*, 2020.
- Artidoro Pagnoni, Vidhisha Balachandran, and Yulia Tsvetkov. Understanding factuality in abstractive summarization with FRANK: A benchmark for factuality metrics. In *NAACL*, 2021.
- Damian Pascual, Béni Egressy, Clara Meister, Ryan Cotterell, and Roger Wattenhofer. A plug-and-play method for controlled text generation. In *EMNLP Findings*, 2021.
- Ethan Perez, Douwe Kiela, and Kyunghyun Cho. True few-shot learning with language models. In *NeurIPS*, 2021.
- Shrimai Prabhumoye, Yulia Tsvetkov, Ruslan Salakhutdinov, and Alan W. Black. Style transfer through back-translation. In *ACL*, 2018.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 2019.
- Mengye Ren, Wenyan Zeng, Binh Yang, and Raquel Urtasun. Learning to reweight examples for robust deep learning. In *ICML*, 2018.
- Teven Le Scao and Alexander M. Rush. How many data points is a prompt worth? In *NAACL*, 2021.
- Timo Schick and Hinrich Schütze. Exploiting cloze-questions for few-shot text classification and natural language inference. In *EACL*, 2021a.

- Timo Schick and Hinrich Schütze. Few-shot text generation with natural language instructions. In *EMNLP*, 2021b.
- Timo Schick and Hinrich Schütze. Generating datasets with pretrained language models. In *EMNLP*, 2021c.
- Timo Schick and Hinrich Schütze. It’s not just size that matters: Small language models are also few-shot learners. In *NAACL*, 2021d.
- Iyer Shankar, Dandekar Nikhil, and Csernai Kornél. First Quora dataset release: Question pairs, 2017. URL <https://www.quora.com/q/quoradata/First-Quora-Dataset-Release-Question-Pairs>.
- Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. Eliciting knowledge from language models using automatically generated prompts. In *EMNLP*, 2020.
- Jun Shu, Qi Xie, Lixuan Yi, Qian Zhao, Sanping Zhou, Zongben Xu, and Deyu Meng. Meta-weightnet: Learning an explicit mapping for sample weighting. In *NeurIPS*, 2019.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *EMNLP*, 2013.
- Trevor Scott Standley, Amir Roshan Zamir, Dawn Chen, Leonidas J. Guibas, Jitendra Malik, and Silvio Savarese. Which tasks should be learned together in multi-task learning? In *ICML*, 2019.
- Derek Tam, Rakesh R Menon, Mohit Bansal, Shashank Srivastava, and Colin Raffel. Improving and simplifying pattern exploiting training. In *EMNLP*, 2021.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *EMNLP Workshop BlackboxNLP*, 2018.
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. SuperGLUE: A stickier benchmark for general-purpose language understanding systems. In *NeurIPS*, 2019.
- Yu-Xiong Wang, Deva Ramanan, and Martial Hebert. Learning to model the tail. In *NIPS*, 2017.
- Zirui Wang, Adams Wei Yu, Orhan Firat, and Yuan Cao. Towards zero-label language learning. *ArXiv*, abs/2109.09193, 2021.
- Alex Warstadt, Amanpreet Singh, and Samuel R Bowman. Neural network acceptability judgments. In *TACL*, 2019.
- Jason Wei and Kai Zou. EDA: Easy data augmentation techniques for boosting performance on text classification tasks. In *EMNLP*, 2019.
- Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus for sentence understanding through inference. In *NAACL-HLT*, 2018.
- Lijun Wu, Fei Tian, Yingce Xia, Yang Fan, Tao Qin, Jianhuang Lai, and Tie-Yan Liu. Learning to teach with dynamic loss functions. In *NeurIPS*, 2018.
- Qizhe Xie, Zihang Dai, Eduard H. Hovy, Minh-Thang Luong, and Quoc V. Le. Unsupervised data augmentation for consistency training. In *NeurIPS*, 2020.
- Kevin Yang and Dan Klein. FUDGE: Controlled text generation with future discriminators. In *NAACL*, 2021.

- Yiben Yang, Chaitanya Malaviya, Jared Fernandez, Swabha Swayamdipta, Ronan Le Bras, Ji ping Wang, Chandra Bhagavatula, Yejin Choi, and Doug Downey. G-daug: Generative data augmentation for commonsense reasoning. In *EMNLP Findings*, 2020.
- Jiacheng Ye, Jiahui Gao, Qintong Li, Hang Xu, Jiangtao Feng, Zhiyong Wu, Tao Yu, and Lingpeng Kong. Zerogen: Efficient zero-shot learning via dataset generation. *ArXiv*, abs/2202.07922, 2022.
- Kang Min Yoo, Do-Hyoung Park, Jaewoo Kang, Sang-Woo Lee, and Woomyeong Park. GPT3Mix: Leveraging large-scale language models for text augmentation. In *EMNLP Findings*, 2021.
- Ningyu Zhang, Luoqiu Li, Xiang Chen, Shumin Deng, Zhen Bi, Chuanqi Tan, Fei Huang, and Huajun Chen. Differentiable prompt makes pre-trained language models better few-shot learners. In *ICLR*, 2022.
- Tony Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. Calibrate before use: Improving few-shot performance of language models. In *ICML*, 2021.
- Zexuan Zhong, Dan Friedman, and Danqi Chen. Factual probing is [mask]: Learning vs. learning to recall. In *NAACL*, 2021.
- Daniel M. Ziegler, Nisan Stiennon, Jeff Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *ArXiv*, abs/1909.08593, 2019.

## A DERIVATION OF META WEIGHT GRADIENT UPDATE

We first write out the gradient update of  $\hat{\theta}_p^{(t)}$  ( $\omega^{(t)}$ ) and  $\omega^{(t+1)}$  according to Algorithm 1 as follows:

$$\hat{\theta}_p^{(t)}(\omega^{(t)}) = \theta_p^{(t)} - \alpha \left. \frac{\partial \mathcal{L}_{w\text{-gen}}(\theta_p; \omega^{(t)})}{\partial \theta_p} \right|_{\theta_p = \theta_p^{(t)}} = \theta_p^{(t)} - \alpha \sum_{j=1}^n w_j(\omega^{(t)}) \left. \frac{\partial \mathcal{L}_{\text{gen}}^j(\theta_p)}{\partial \theta_p} \right|_{\theta_p = \theta_p^{(t)}} \quad (6)$$

$$\omega^{(t+1)} = \omega^{(t)} - \beta \left. \frac{\partial \mathcal{L}_{\text{disc}}(\hat{\theta}_p^{(t)}(\omega))}{\partial \omega} \right|_{\omega = \omega^{(t)}}. \quad (7)$$

where  $\alpha$  and  $\beta$  are step sizes.

The gradient in Equation (7) is calculated as:

$$\begin{aligned} & \left. \frac{\partial \mathcal{L}_{\text{disc}}(\hat{\theta}_p^{(t)}(\omega))}{\partial \omega} \right|_{\omega = \omega^{(t)}} \\ &= \left. \frac{\partial \mathcal{L}_{\text{disc}}(\hat{\theta}_p)}{\partial \hat{\theta}_p} \right|_{\hat{\theta}_p = \hat{\theta}_p^{(t)}} \left. \frac{\partial \hat{\theta}_p(\omega)}{\partial \omega} \right|_{\omega = \omega^{(t)}} \\ &= \left. \frac{\partial \mathcal{L}_{\text{disc}}(\hat{\theta}_p)}{\partial \hat{\theta}_p} \right|_{\hat{\theta}_p = \hat{\theta}_p^{(t)}} \left( -\alpha \sum_{j=1}^n \left. \frac{\partial \mathcal{L}_{\text{gen}}^j(\theta_p)}{\partial \theta_p} \right|_{\theta_p = \theta_p^{(t)}}^\top \left. \frac{\partial w_j(\omega)}{\partial \omega} \right|_{\omega = \omega^{(t)}} \right) \quad \text{Plugging in Eq. (6)} \\ &= -\alpha \sum_{j=1}^n \left( \underbrace{\left. \frac{\partial \mathcal{L}_{\text{disc}}(\hat{\theta}_p)}{\partial \hat{\theta}_p} \right|_{\hat{\theta}_p = \hat{\theta}_p^{(t)}} \left. \frac{\partial \mathcal{L}_{\text{gen}}^j(\theta_p)}{\partial \theta_p} \right|_{\theta_p = \theta_p^{(t)}}^\top}_{\triangleq d_j} \left. \frac{\partial w_j(\omega)}{\partial \omega} \right|_{\omega = \omega^{(t)}} \right) \end{aligned}$$

Therefore,

$$-\left. \frac{\partial \mathcal{L}_{\text{disc}}(\hat{\theta}_p^{(t)}(\omega))}{\partial \omega} \right|_{\omega = \omega^{(t)}} \propto \sum_{j=1}^n d_j \left. \frac{\partial w_j(\omega)}{\partial \omega} \right|_{\omega = \omega^{(t)}}, \quad d_j = \left. \frac{\partial \mathcal{L}_{\text{disc}}(\hat{\theta}_p)}{\partial \hat{\theta}_p} \right|_{\hat{\theta}_p = \hat{\theta}_p^{(t)}} \left. \frac{\partial \mathcal{L}_{\text{gen}}^j(\theta_p)}{\partial \theta_p} \right|_{\theta_p = \theta_p^{(t)}}^\top$$

## B GLUE TASKS

We provide the details of the seven classification tasks included in the GLUE benchmark.

**MNLI:** Multi-genre Natural Language Inference (Williams et al., 2018) requires predicting whether a given premise sentence entails, contradicts or neutral with respect to a given hypothesis sentence.

**QQP:** Quora Question Pairs (Shankar et al., 2017) requires judging whether a pair of questions asked are semantically equivalent.

**QNLI:** Question Natural Language Inference requires predicting whether a given sentence contains the answer to a given question sentence.

**SST-2:** Stanford Sentiment Treebank (Socher et al., 2013) requires determining if a movie review has positive or negative sentiment.

**CoLA:** Corpus of Linguistic Acceptability (Warstadt et al., 2019) requires determining whether a given sentence is linguistically acceptable or not.

**RTE:** Recognizing Textual Entailment (Bentivogli et al., 2009; Dagan et al., 2005; Giampiccolo et al., 2007; Haim et al., 2006) requires predicting whether a given premise sentence entails a given hypothesis sentence or not.



**MRPC:** Microsoft Research Paraphrase Corpus (Dolan & Brockett, 2005) requires predicting whether two sentences are semantically equivalent or not.

## C IMPLEMENTATION DETAILS

Table 3: Prompts used for initializing the prefix vectors and control codes (required by CTRL (Keskar et al., 2019)) used in generator training. The control codes are selected to approximate the domain of the task. For single-sequence tasks,  $x$  denotes the training sample; for sequence-pair tasks,  $x_1$  and  $x_2$  denote the first and second sequence in the training sample, respectively.

| Task  | Task Type       | Control Code | Label                                  | Initialization Prompt  |
|-------|-----------------|--------------|--|--|
| SST-2 | single-sequence | Reviews      | positive<br>negative                   | Rating: 5.0 positive movie review: $x$<br>Rating: 1.0 negative movie review: $x$   |
| CoLA  | single-sequence | Links        | grammatical<br>not grammatical         | Linguistically correct sentence: $x$<br>Linguistically incorrect sentence: $x$   |
| MNLI  | sequence-pair   | Wikipedia    | entailment<br>neutral<br>contradiction | Sentence 1 implies Sentence 2. Sentence 1: $x_1$ Sentence 2: $x_2$<br>Sentence 2 supplements Sentence 1. Sentence 1: $x_1$ Sentence 2: $x_2$<br>Sentence 2 contradicts Sentence 1. Sentence 1: $x_1$ Sentence 2: $x_2$ |
| QNLI  | sequence-pair   | Links        | entailment<br>not entailment           | Paragraph is relevant to Question. Question: $x_1$ Paragraph: $x_2$<br>Paragraph is irrelevant to Question. Question: $x_1$ Paragraph: $x_2$   |
| RTE   | sequence-pair   | Wikipedia    | entailment<br>not entailment           | Sentence 1 implies Sentence 2. Sentence 1: $x_1$ Sentence 2: $x_2$<br>Sentence 2 supplements Sentence 1. Sentence 1: $x_1$ Sentence 2: $x_2$   |
| MRPC  | sequence-pair   | Wikipedia    | equivalent<br>not equivalent           | Sentence 1 is equivalent to Sentence 2. Sentence 1: $x_1$ Sentence 2: $x_2$<br>Sentence 1 is different from Sentence 2. Sentence 1: $x_1$ Sentence 2: $x_2$  |
| QQP   | sequence-pair   | Links        | equivalent<br>not equivalent           | Question 1 is equivalent to Question 2. Question 1: $x_1$ Question 2: $x_2$<br>Question 1 is different from Question 2. Question 1: $x_1$ Question 2: $x_2$  |

**Details of Initialization Prompts Used for Generator Tuning on Different Tasks.** For generator tuning, we find it beneficial to initialize the prefix vectors with task-descriptive prompts, similar to the observations in Li & Liang (2021). The prefix lengths (*i.e.*, number of trained prefix token positions) are equal to the number of tokens in the prompts. We present details about the prompts used for initializing the prefix vectors for different tasks in Table 3. For sequence-pair tasks, an additional infix prompt is used between the two sequences, and we also tune the embeddings of the infix (*i.e.*, prompt-tuning (Lester et al., 2021)) for generator training.

**Details of Generator Tuning.** The meta-weighted generator tuning procedure (Algorithm 1) involves three forward and backward passes, and thus its time complexity is approximately 3 times of standard generator training without meta learning. However, since the few-shot training sets have a small amount of training data, the extra time cost is usually affordable. In practice, our generator tuning with meta weight learning takes 10 minutes to train on each task (the standard generator training time without meta-learning is 3.5 minutes). We use a fixed set of hyperparameters for all tasks without task-specific hyperparameter tuning: We set batch size to be 2, the learning rate for optimizing  $\theta_p$  to be  $5e - 3$ , the learning rate for optimizing  $\omega$  to be  $1e - 2$ , and training epoch to be 20.

**Details of Generating Training Data.** Following Meng et al. (2022), for sequence-pair tasks (MNLI, QQP, QNLI, RTE and MRPC), we randomly sample the first sequence from the pretraining corpus (*e.g.*, Wikipedia) and use greedy sampling for generating the second sequence. For single-sequence tasks (SST-2 and CoLA), we use top- $k$  sampling with temperature to generate training data from scratch where  $k = 10$ . For all tasks, we generate 5,000 samples per label.

For SST-2, we use one of the following tokens to start generation: “a”, “one”, “the”, “this”, “that”, “i”, “you”, “it”, “what”. For CoLA, we use a random stop word to start generation.

**Hyperparameters for Fine-Tuning Classifier PLMs.** For fine-tuning on the few-shot training samples  $\mathcal{D}_{\text{train}}$ , we search among the following hyperparameter ranges based on development set ( $\mathcal{D}_{\text{dev}}$ ) model performance and pick the best performing model for further fine-tuning on synthesized data: Learning rate in  $[1e - 5, 2e - 5]$  and batch size in  $[2, 4, 8]$ . The number of training steps is fixed to be 1000. We also find it beneficial to apply label smoothing (smoothing weight set to 0.15) for fine-tuning on the few-shot training set.

For fine-tuning on the synthesized training samples  $\mathcal{D}_{\text{gen}}$ , we use the following hyperparameters:  $5e-6$  as the learning rate; 16 as the batch size; label smoothing weight  $\epsilon = 0.15$ ; temporal ensemble momentum  $\gamma = 0.9$ ; temporal ensemble loss weight  $\lambda = 20$ ; training steps  $T = 6,000$ .

**Details of Temporal Ensembling for Fine-Tuning Classifier PLMs on Synthetic Data.** We update ensembled predictions  $\bar{z}$  as follows where  $p_\phi$  is the current model prediction,  $\gamma$  is the momentum parameter,  $\hat{z}$  is the accumulated model prediction before bias correction,  $\bar{z}$  is the accumulated model prediction after bias correction, and  $t$  is the number of updates  $\bar{z}$  has received:

$$\hat{z} \leftarrow \gamma \hat{z} + (1 - \gamma)p_\phi, \bar{z} \leftarrow \hat{z}/(1 - \gamma^t).$$

The accumulated model prediction  $\hat{z}$  has a zero initialization; the division  $(1 - \gamma^t)$  is for bias correction (Laine & Aila, 2017). After each update of  $\hat{z}$ , it will be compared to a threshold value  $\delta$ ; each synthesized sample  $(\tilde{x}, \tilde{y})$  will be included in training only if  $\bar{z}_{\tilde{y}} > \delta$ .

We update the ensembled predictions  $\bar{z}$  on all samples in  $\mathcal{D}_{\text{gen}}$  every 200 steps, and set the threshold value for sample filtering  $\delta = 0.8$ .

**Computation Environment.** The experiments are conducted on NVIDIA A100 GPUs.

## D DATA AUGMENTATION BASELINE DETAILS

**Details About MixText (Chen et al., 2020).** We use the TMix version of MixText to perform data interpolation on the few-shot labeled dataset (since there is no access to unlabeled task-specific data under the strict few-shot learning setting Gao et al. (2021)). We adapt the label mix-up operation to fit prompt-based fine-tuning by interpolating the label words instead of categorical labels; we observe that this results in better few-shot performance than the original TMix, probably analogous to why prompt-based fine-tuning outperforms standard fine-tuning for few-shot learning. We train the classifier with supervised loss combined with consistency loss over the interpolated samples as in the original paper. We follow the default hyperparameters in MixText.

**Details About Back Translation.** We use two trained Marian (Junczys-Dowmunt et al., 2018) models to perform data augmentation via back translation. We translate our labeled examples from English to French, and then back to English. As in UDA (Xie et al., 2020), we employ random sampling with a tunable temperature to generate a diverse set of derivative examples. We generate 32 examples from each few-shot training example and let the synthesized samples share the same label with the original few-shot training sample. After combining with the original examples, we fine-tune the classifier and observe performance.

**Details About GPT3Mix (Yoo et al., 2021).** We use the 175B GPT3 model for generating the augmentations. For creating each augmentation, we randomly sample  $k = 4$  (the optimal setting according to GPT3Mix) examples from the few-shot training set as demonstrations. The prompts follow the suggested format proposed in the original paper (Yoo et al., 2021) and are shown in Table 4. We create 5,000 augmented samples per label to make the resulting training set size equal to that of FewGen. After obtaining the augmented examples and their pseudo labels (the probability predictions over all labels by GPT3), we use them along with the real few-shot samples for fine-tuning the classifier, following the setting in GPT3Mix (Yoo et al., 2021).

**Details About Standard Generator Fine-Tuning.** We fine-tune the same 1.6B CTRL (Keskar et al., 2019) model as used in FewGen with the standard maximum likelihood objective. Different from previous studies (Anaby-Tavor et al., 2020; Kumar et al., 2020) that prepend categorical labels to the training samples, we enhance the generator fine-tuning with label-descriptive prompts (shown in Table 3) used in FewGen. We create 5,000 augmented samples per label to make the resulting training set size equal to that of FewGen.

## E DETAILS OF WEIGHTING NETWORK IMPLEMENTATION

Table 4: Prompts used for GPT3Mix augmentation. For sequence-pair tasks,  $x_1$  and  $x_2$  denote the first and second input sequence, respectively. For single-sequence tasks,  $x$  denotes the input sequence.  $y$  denotes the label name. Only one example is shown in the template for clarity; in practice, we concatenate  $k = 4$  samples according to the optimal setting in GPT3Mix (Yoo et al., 2021).

| Task  | Template  | Label name   |
|-------|---|--|
| SST-2 | Each item in the following list contains a movie review and the respective sentiment.<br>The sentiment is one of 'positive' or 'negative'.<br>Movie review: $x$ (Sentiment: $y$ ) . . .   | positive: positive<br>negative: negative                                   |
| CoLA  | Each item in the following list contains a text and the respective grammar.<br>The grammar is one of 'correct' or 'incorrect'.<br>Text: $x$ (Grammar: $y$ ) . . .   | grammatical: correct<br>not grammatical: incorrect                         |
| MNLI  | Each item in the following list contains a premise, a hypothesis and their logical relation.<br>The logical relation is one of 'entailment', 'neutral' or 'contradiction'.<br>Premise: $x_1$ Hypothesis: $x_2$ (Logical relation: $y$ ) . . . | entailment: entailment<br>neutral: neutral<br>contradiction: contradiction |
| QNLI  | Each item in the following list contains a question, an answer and their logical relation.<br>The logical relation is one of 'entailment' or 'neutral'.<br>Question: $x_1$ Answer: $x_2$ (Logical relation: $y$ ) . . .                       | entailment: entailment<br>not entailment: neutral                          |
| RTE   | Each item in the following list contains a premise, a hypothesis and their logical relation.<br>The logical relation is one of 'entailment' or 'neutral'.<br>Premise: $x_1$ Hypothesis: $x_2$ (Logical relation: $y$ ) . . .                  | entailment: entailment<br>not entailment: neutral                          |
| MRPC  | Each item in the following list contains two sentences and their semantic relation.<br>The semantic relation is one of 'equivalent' or 'different'.<br>Sentence 1: $x_1$ Sentence 2: $x_2$ (Semantic relation: $y$ ) . . .                    | equivalent: equivalent<br>not equivalent: different                        |
| QQP   | Each item in the following list contains two questions and their semantic relation.<br>The semantic relation is one of 'equivalent' or 'different'.<br>Question 1: $x_1$ Question 2: $x_2$ (Semantic relation: $y$ ) . . .                    | equivalent: equivalent<br>not equivalent: different                        |

Since the token weights  $w$  used in Eq. (4) need to characterize the discriminativeness of each token, we use the value of discriminative objective at each token  $\mathcal{L}_{\text{disc}}^j$  as the input to the weighting network, and we use softmax to normalize the weights:

$$w_j(\omega) = \frac{\exp\left(g_\omega(\mathcal{L}_{\text{disc}}^j)\right)}{\sum_{j'=1}^n \exp\left(g_\omega(\mathcal{L}_{\text{disc}}^{j'})\right)}.$$

Following Shu et al. (2019), we instantiate  $g_\omega$  to be a feedforward network (FFN) with only one 100-dimension hidden layer by default. We explore an alternative instantiation that adds one self-attention layer on top of the generator PLM’s output hidden states. The meta weights are finally obtained by projecting the outputs of the self-attention layer using another linear layer. We evaluate the resulting generator quality via the same two metrics as in Appendix G. Table 5 shows that using more complicated architectures (*e.g.*, adding another self-attention layer) does not result in a better generator compared to using a simple FFN for meta weight learning. This is probably because the generator PLM’s output representations are sufficiently contextualized and contain the information necessary for learning the token weights, thus a simple FFN as the weighting network will be enough. Using more complicated networks, on the other hand, will introduce more randomly initialized new parameters which may not be learned well using the limited amount of few-shot training data.

## F VISUALIZATION OF TOKEN WEIGHT LEARNING

To gain intuitive understanding of what tokens are assigned more weight during generator tuning, we visualize the learned weights in Fig. 5. The tokens with higher weights (*e.g.*, “weak” in the first example and “hates” in the second example) are learned to be important tokens that decide the relation of the second sentence to the first sentence (*i.e.*, the label of the training sample). With such tokens emphasized during training, the generator is encouraged to capture label-discriminative information that facilitates the generation of unambiguous training samples.

Table 5: Study of weighting network instantiation. The default architecture is a feedforward network (FFN) with one hidden layer. We also explore adding a self-attention layer on top of the generator PLM’s output hidden states (Self-attention). We use the same two metrics with Table 6 to evaluate the resulting generators.

| Architecture   | MNLI        |             | SST-2       |             |
|----------------|-------------|-------------|-------------|-------------|
|                | Acc. (↑)    | PPL (↓)     | Acc. (↑)    | PPL (↓)     |
| FFN            | <b>72.3</b> | <b>11.9</b> | <b>93.2</b> | <b>43.5</b> |
| Self-attention | 70.3        | 12.9        | 92.3        | 44.2        |

|  |   |
|--|---|
| <b>Sentence 1:</b> But prophecy is always strongest when based on coincidence--that is a prime rule.   | <b>Label:</b> Contradiction                                 |
| <b>Sentence 2:</b> Prophecies based on coincidences are widely known to be weak and unreliable.  |   |
| <b>weights</b>   | 0.03 0.02 0.03 0.11 0.06 0.06 0.03 0.03 0.05 0.33 0.06 0.21 |
| <b>Sentence 1:</b> But Rodgers did tell Lewis that he despises Amelio because Amelio supported Clinton, so it is Rodgers' mistake, not our author's, that we are correcting. | <b>Label:</b> Entailment                                    |
| <b>Sentence 2:</b> Rodgers told Lewis he hates Amelio.   |   |
| <b>weights</b>   | 0.14 0.08 0.07 0.08 0.47 0.17                               |

Figure 5: Visualization of learned token weights on two samples from MNLI’s few-shot training set. The generator is trained given the first sentence to generate the second. The tokens associated with higher weights are more label indicative.

Table 6: Evaluation of generator training objectives. We use two metrics: Generated data accuracy (Acc; higher is better) and generator’s perplexity on the test set (PPL; lower is better). The results are averaged over 5  $\mathcal{D}_{\text{train}}/\mathcal{D}_{\text{dev}}$  splits.

| Objective  | MNLI        |             | QQP         |             | QNLI        |             | SST-2       |             | CoLA        |             | RTE         |            | MRPC        |            |
|--|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|------------|-------------|------------|
|  | Acc. (↑)    | PPL (↓)     | Acc. (↑)    | PPL (↓)     | Acc. (↑)    | PPL (↓)     | Acc. (↑)    | PPL (↓)     | Acc. (↑)    | PPL (↓)     | Acc. (↑)    | PPL (↓)    | Acc. (↑)    | PPL (↓)    |
| $\mathcal{L}_{\text{gen}}$                             | 69.4        | 13.1        | 87.5        | 10.9        | 57.0        | 23.4        | 91.5        | 43.8        | 59.1        | 85.6        | 82.9        | 9.3        | 87.6        | 5.0        |
| $\mathcal{L}_{\text{gen}} + \mathcal{L}_{\text{disc}}$ | 70.2        | 13.5        | 87.3        | 11.2        | 57.2        | 24.8        | 92.0        | 49.5        | 59.2        | 87.0        | 82.8        | 9.6        | 86.3        | 5.3        |
| $\mathcal{L}_{\text{w-gen}}$                           | <b>72.3</b> | <b>11.9</b> | <b>89.5</b> | <b>10.7</b> | <b>60.1</b> | <b>23.2</b> | <b>93.2</b> | <b>43.5</b> | <b>60.7</b> | <b>83.8</b> | <b>83.4</b> | <b>8.9</b> | <b>90.5</b> | <b>4.9</b> |

Table 7: (For MNLI) Examples of generated second sequence (hypothesis) by generators tuned with three different objectives conditioned on a given first sequence (premise) “*In 2009, hair samples from 1,137 Taiji residents were tested for mercury by the National Institute for Minamata Disease (NIMD)*”. The true label of the generated sequence is marked at the end of the sequence (if the generated sequence correctly pertains to the target label, it is marked in blue; otherwise, it is in red).

| Objective  | Label         | Generated Second Sequence   |
|--|---------------|---|
| $\mathcal{L}_{\text{gen}}$                             | entailment    | The NIMD found that hair samples from 1,137 Taiji residents had mercury in their hair. ( <i>neutral</i> )                               |
|  | neutral       | The NIMD found that there was no evidence of a link between exposure to high levels of mercury and thyroid cancer. ( <i>neutral</i> )   |
|  | contradiction | There was no evidence of mercury in hair samples from Taiji. ( <i>neutral</i> )   |
| $\mathcal{L}_{\text{gen}} + \mathcal{L}_{\text{disc}}$ | entailment    | The number of hairs in a sample is equal to the number of people who lived in Taiji. ( <i>neutral</i> )                                 |
|  | neutral       | The results showed that there was no significant difference in levels of mercury. ( <i>neutral</i> )                                    |
|  | contradiction | Hair samples from 1,137 Taiji residents were not tested. ( <i>contradiction</i> )   |
| $\mathcal{L}_{\text{w-gen}}$                           | entailment    | The NIMD tested hair samples from 1,137 residents of Taiji. ( <i>entailment</i> )   |
|  | neutral       | There was no significant difference in levels between people who lived near a nickel mine and those living far away. ( <i>neutral</i> ) |
|  | contradiction | The NIMD did not test any of the hair samples. ( <i>contradiction</i> )   |

## G QUANTITATIVE ANALYSES OF GENERATOR TRAINING OBJECTIVES

Apart from the final classification model performance which indirectly reflects the synthetic data quality, we additionally conduct more direct quantitative analyses of different generator training objectives. We use two metrics: (1) The accuracy of generated texts, which is judged by fully-supervised RoBERTa<sub>Large</sub> models fine-tuned on the original training sets of each task. We choose to adopt such an automatic evaluation instead of human evaluation because it is efficient and reliable—fully-supervised RoBERTa<sub>Large</sub> models have comparable or better accuracy than human baselines according to the GLUE benchmark<sup>4</sup>. (2) The generator’s perplexity on the test sets, which reflects how well the generator models the task distribution. As shown in Table 6, using  $\mathcal{L}_{\text{w-gen}}$  for generator training consistently outperforms using  $\mathcal{L}_{\text{gen}}$  or  $\mathcal{L}_{\text{gen}} + \mathcal{L}_{\text{disc}}$ , both in generated text accuracy and in language modeling ability.

Comparing  $\mathcal{L}_{\text{w-gen}}$  with  $\mathcal{L}_{\text{gen}}$ , the meta weights automatically learned emphasize discriminative tokens in generator training and help the generator capture the subtle semantic differences across different labels, resulting in better language modeling quality and more distinctive generated data.

Comparing  $\mathcal{L}_{\text{w-gen}}$  with  $\mathcal{L}_{\text{gen}} + \mathcal{L}_{\text{disc}}$ , the generator training objective is not directly impacted by the discriminative objective, thus avoiding the gradient interference issue in multi-task learning (Standley et al., 2019)—the gradient for optimizing the generative probability  $p(\mathbf{x}|y_l)$  will be interfered by

<sup>4</sup><https://gluebenchmark.com/leaderboard>

that optimizing the discriminative probability  $p(y_l|\mathbf{x})$  if  $\mathcal{L}_{\text{gen}} + \mathcal{L}_{\text{disc}}$  is used. Therefore, using  $\mathcal{L}_{\text{w-gen}}$  results in better language modeling quality and more fluent and coherent generation results.

We also showcase concrete generation results for the three labels of MNLI by models trained with the three different loss functions in Table 7. The model trained with  $\mathcal{L}_{\text{gen}}$  produces fluent and coherent sentences, but the generated sentences do not accurately pertain to the desired label (*i.e.*, the “entailment” and “contradiction” generation results are in fact neutral with respect to the given sentence), lacking label discriminativeness. When  $\mathcal{L}_{\text{gen}} + \mathcal{L}_{\text{disc}}$  is used, the generated samples of different labels are more distinctive, but also become less natural and coherent due to the model’s language modeling ability being hampered. The generator tuned with  $\mathcal{L}_{\text{w-gen}}$  produces both coherent and label-discriminative samples which can serve as quality training data.

Table 8: 16-shot training samples of SST-2.

| Label    | Example | Review Text  |
|----------|---------|--|
| positive | #1      | (ramsay) visually transforms the dreary expanse of dead-end distaste the characters inhabit into a poem of art , music and metaphor .  |
|          | #2      | the film jolts the laughs from the audience – as if by cattle prod .   |
|          | #3      | the film presents visceral and dangerously honest revelations about the men and machines behind the curtains of our planet .   |
|          | #4      | a film that will enthrall the whole family .   |
|          | #5      | serious movie-goers embarking upon this journey will find that the road to perdition leads to a satisfying destination .   |
|          | #6      | sweet and memorable film .   |
|          | #7      | shyamalan takes a potentially trite and overused concept (aliens come to earth) and infuses it into a rustic , realistic , and altogether creepy tale of hidden invasion .   |
|          | #8      | a crisp psychological drama (and) a fascinating little thriller that would have been perfect for an old “ twilight zone ” episode .  |
|          | #9      | my big fat greek wedding is not only the best date movie of the year , it ’s also a – dare i say it twice – delightfully charming – and totally american , i might add – slice of comedic bliss .  |
|          | #10     | a comedy-drama of nearly epic proportions rooted in a sincere performance by the title character undergoing midlife crisis .   |
|          | #11     | diggs and lathan are among the chief reasons brown sugar is such a sweet and sexy film .   |
|          | #12     | you ’re not merely watching history , you ’re engulfed by it .   |
|          | #13     | the concept is a hoot .  |
|          | #14     | the filmmakers ’ eye for detail and the high standards of performance convey a strong sense of the girls ’ environment .   |
|          | #15     | a haunting tale of murder and mayhem .   |
|          | #16     | neil burger here succeeded in ... making the mystery of four decades back the springboard for a more immediate mystery in the present .  |
| negative | #1      | nothing happens , and it happens to flat characters .  |
|          | #2      | as lively an account as seinfeld is deadpan .  |
|          | #3      | so we got ten little indians meets friday the 13th by way of clean and sober , filmed on the set of carpenter ’s the thing and loaded with actors you ’re most likely to find on the next inevitable incarnation of the love boat .                                |
|          | #4      | the plot is nothing but boilerplate cliches from start to finish , and the script assumes that not only would subtlety be lost on the target audience , but that it ’s also too stupid to realize that they ’ve already seen this exact same movie a hundred times |
|          | #5      | ultimately , sarah ’s dedication to finding her husband seems more psychotic than romantic , and nothing in the movie makes a convincing case that one woman ’s broken heart outweighs all the loss we witness .   |
|          | #6      | the big finish is a bit like getting all excited about a chocolate eclair and then biting into it and finding the filling missing .  |
|          | #7      | this picture is mostly a lump of run-of-the-mill profanity sprinkled with a few remarks so geared toward engendering audience sympathy that you might think he was running for office – or trying to win over a probation officer .                                |
|          | #8      | just because a walk to remember is shrewd enough to activate girlish tear ducts does n’t mean it ’s good enough for our girls .  |
|          | #9      | often lingers just as long on the irrelevant as on the engaging , which gradually turns what time is it there ?  |
|          | #10     | this movie , a certain scene in particular , brought me uncomfortably close to losing my lunch .   |
|          | #11     | but it would be better to wait for the video .   |
|          | #12     | a rude black comedy about the catalytic effect a holy fool has upon those around him in the cutthroat world of children ’s television .  |
|          | #13     | just a collection of this and that – whatever fills time – with no unified whole .   |
|          | #14     | although god is great addresses interesting matters of identity and heritage , it ’s hard to shake the feeling that it was intended to be a different kind of film .   |
|          | #15     | the chocolate factory without charlie .  |
|          | #16     | in that setting , their struggle is simply too ludicrous and borderline insulting .  |

## H CONCRETE GENERATION RESULTS

We present few-shot training samples ( $\mathcal{D}_{\text{train}}$ ) of SST-2 in Table 8, and some concrete generation results (from  $\mathcal{D}_{\text{gen}}$ ) for all tasks in Tables 9, 10, 11, 12, 13, 14, and 15.

As shown in Tables 8 and 9, the synthetic samples are accurate and quite different from the given training samples to serve as effective augmentations. For sequence-pair tasks, because we randomly sample the first sequence from the pretraining corpus and let the generator create the second sequence given certain labels, the resulting generated samples will be certainly different from the given training samples.

Table 9: Examples of concrete generations for SST-2.

| Example | Label    | Generated Text  |
|---------|----------|---|
| #1      | positive | i was a bit skeptical when i first saw the trailer for this movie, but after watching it, it was quite enjoyable.   |
|         | negative | i was expecting a fun and entertaining movie with some nice special effects, but instead it was just plain boring.  |
| #2      | positive | that is, a film that has the potential to move an audience to tears, or at least to a state of mind where you feel something.   |
|         | negative | a film that has the potential to be one of the most moving and powerful experiences ever made is instead marred by an unfocused script, bad acting and a lack of emotional depth. |
| #3      | positive | this is a good movie – it’s not just for people who like to see things blow up or have explosions or shoot guns.  |
|         | negative | i am a big fan of the first two movies but this one is just plain bad.  |

Table 10: Examples of concrete generations for CoLA.

| Example | Label           | Generated Text   |
|---------|-----------------|--|
| #1      | grammatical     | Between the two of us, I think you have a good chance to make it to the finals.  |
|         | not grammatical | A few weeks later we saw them on television. We thought there were a hundred more. It wasn’t that they came out one at time. |
| #2      | grammatical     | Doing the dishes is a task that I have to do because I am not allowed to go out and play with my friends.                    |
|         | not grammatical | Only those whose lives depend on it, but only when you are at its mercy do it harm.  |
| #3      | grammatical     | To be able to speak the language of the other person is a prerequisite for communication.                                    |
|         | not grammatical | The only good way to avoid getting into trouble for being too friendly would to let me go on a trip.                         |

Table 11: Examples of concrete generations for QQP. The first question (*italicized*) is randomly sampled from the pretraining corpus; the second question (underlined) is generated by  $G_{\theta_p}$ .

| Example | Label          | Generated Text   |
|---------|----------------|--|
| #1      | equivalent     | Question 1: <i>How long does it take for a project to be completed?</i><br>Question 2: <u>How long does it take to complete a project?</u>   |
|         | not equivalent | Question 1: <i>How long does it take for a project to be completed?</i><br>Question 2: <u>What is the total cost of completing this project?</u>   |
| #2      | equivalent     | Question 1: <i>What mascots or characters would you like to see included as Super Mario Maker DLC in the future?</i><br>Question 2: <u>What would you like to see in Super Mario Maker DLC that you did not see in the game?</u> |
|         | not equivalent | Question 1: <i>What mascots or characters would you like to see included as Super Mario Maker DLC in the future?</i><br>Question 2: <u>How do I get a copy of this game?</u>   |

Table 12: Examples of concrete generations for MNL. The first sentence (*italicized*) is randomly sampled from the pretraining corpus; the second sentence (underlined) is generated by  $G_{\theta_p}$ .

| Example | Label         | Generated Text  |
|---------|---------------|---|
| #1      | entailment    | Sentence 1: <i>Air is provided for the combustion by an electric blower.</i><br>Sentence 2: <u>The blower provides air to a combustion chamber.</u>   |
|         | neutral       | Sentence 1: <i>Air is provided for the combustion by an electric blower.</i><br>Sentence 2: <u>Electric blowers are available in most gas stations.</u>   |
|         | contradiction | Sentence 1: <i>Air is provided for the combustion by an electric blower.</i><br>Sentence 2: <u>The blower does not work.</u>  |
| #2      | entailment    | Sentence 1: <i>Since its base is almost at sea level, it is only the 15th highest light in the United States, the first 14 being built on higher ground.</i><br>Sentence 2: <u>It is the 15th highest light in the United States.</u>   |
|         | neutral       | Sentence 1: <i>Since its base is almost at sea level, it is only the 15th highest light in the United States, the first 14 being built on higher ground.</i><br>Sentence 2: <u>The lighthouse was originally constructed to be a beacon for ships passing by and as such has been used since before World War II.</u> |
|         | contradiction | Sentence 1: <i>Since its base is almost at sea level, it is only the 15th highest light in the United States, the first 14 being built on higher ground.</i><br>Sentence 2: <u>It is located on a mountain top.</u>   |

Table 13: Examples of concrete generations for QNLI. The question (*italicized*) is randomly sampled from the pretraining corpus; the answer (underlined) is generated by  $G_{\theta_p}$ .

| Example | Label          | Generated Text   |
|---------|----------------|--|
| #1      | entailment     | Question: <i>What makes you want to step up to the next level?</i><br>Answer: <u>I want to be the best player I can be.</u>  |
|         | not entailment | Question: <i>What makes you want to step up to the next level?</i><br>Answer: <u>The new program will be called "Project 10" and it is expected that a total of \$450 million in federal funding would go toward it.</u>                             |
| #2      | entailment     | Question: <i>How do all those shops know what you would like to buy?</i><br>Answer: <u>The stores are able to track your preferences and provide you with a list of products that are best for you.</u>  |
|         | not entailment | Question: <i>How do all those shops know what you would like to buy?</i><br>Answer: <u>The stores are not required by law or regulation in the United States and Canada but they have been known for years as a source of illegal sales on eBay.</u> |

Table 14: Examples of concrete generations for RTE. The first sentence (*italicized*) is randomly sampled from the pretraining corpus; the second sentence (underlined) is generated by  $G_{\theta_p}$ .

| Example | Label          | Generated Text  |
|---------|----------------|---|
| #1      | entailment     | Sentence 1: <i>In 1970, the Science Council of Canada recommended that the government of Canada immediately invest in industrial development of the aviation industry, including construction of aircraft, navigation aids, and regulation of air traffic.</i><br>Sentence 2: <u>The government of Canada has invested in the aviation industry.</u>          |
|         | not entailment | Sentence 1: <i>In 1970, the Science Council of Canada recommended that the government of Canada immediately invest in industrial development of the aviation industry, including construction of aircraft, navigation aids, and regulation of air traffic.</i><br>Sentence 2: <u>The Aviation Industry was established by a Royal Decree on June 1, 1970.</u> |
| #2      | entailment     | Sentence 1: <i>All of the Centre's staff are fluently bilingual in both English and Chinese and are familiar with the traditions of the Chinese culture.</i><br>Sentence 2: <u>The Centre is a bilingual institution.</u>   |
|         | not entailment | Sentence 1: <i>All of the Centre's staff are fluently bilingual in both English and Chinese and are familiar with the traditions of the Chinese culture.</i><br>Sentence 2: <u>The Centre is a cultural centre for learning about China.</u>  |



Table 15: Examples of concrete generations for MRPC. The first sentence (*italicized*) is randomly sampled from the pretraining corpus; the second sentence (underlined) is generated by  $G_{\theta_p}$ .

| Example | Label          | Generated Text   |
|---------|----------------|--|
| #1      | entailment     | Sentence 1: <i>It embodies the aesthetic ideals of the Queen Anne style through its design and detailing and shows the type of house that attracted prosperous Napans at the end of the Victorian era.</i>   |
|         |                | Sentence 2: <u>It is a fine example of the Queen Anne style of architecture through its design and detailing and shows the type of house that attracted prosperous Napans at the end of the Victorian era.</u>   |
|         | not entailment | Sentence 1: <i>It embodies the aesthetic ideals of the Queen Anne style through its design and detailing and shows the type of house that attracted prosperous Napans at the end of the Victorian era.</i>   |
|         |                | Sentence 2: <u>The building is a fine example in this style, with an elegant facade reminiscent to those found on many grand mansions built by wealthy merchants during America’s Gilded Age.</u>  |
| #2      | entailment     | Sentence 1: <i>Crosbie ran unsuccessfully for the leadership of the Liberal Party of Newfoundland and Labrador in 1969, losing to Smallwood, and was also a candidate in the Progressive Conservative Party of Canada’s 1983 leadership election, placing third.</i> |
|         |                | Sentence 2: <u>Crosbie was a candidate in the Progressive Conservative Party of Canada’s 1983 leadership election, placing third.</u>  |
|         | not entailment | Sentence 1: <i>Crosbie ran unsuccessfully for the leadership of the Liberal Party of Newfoundland and Labrador in 1969, losing to Smallwood, and was also a candidate in the Progressive Conservative Party of Canada’s 1983 leadership election, placing third.</i> |
|         |                | Sentence 2: <u>He lost his bid as leader after he failed twice at running against John Diefenbaker.</u>  |