# PERFECT: Prompt-free and Efficient Language Model Fine-Tuning

Anonymous ACL submission

### Abstract

Current methods for few-shot fine-tuning of pretrained masked language model (PLM) require carefully engineered prompts and verbalizers for each new task, to convert examples into a 004 cloze-format that the PLM can score. In this work, we propose PERFECT, a simple and efficient method for few-shot fine-tuning of PLMs without relving on any such handcrafting, which is highly effective given as few as 32 data points. PERFECT makes two key design choices: First, we show that manually engineered task prompts can be replaced with task-specific adapters that enable 012 sample-efficient fine-tuning and reduce memory and storage costs by roughly factors of 5 and 100, respectively. Second, instead of using handcrafted 016 verbalizers, we learn a new multi-token label embedding during fine-tuning which are not tied to 017 the model vocabulary and which allow us to avoid complex auto-regressive decoding. These embeddings are not only learnable from limited data but also enable nearly 100x faster training and inference. Experiments on a wide range of few shot NLP tasks demonstrate that PERFECT, while being simple and efficient, also outperforms existing state-of-the-art few-shot learning methods.<sup>1</sup>

## 1 Introduction

026

027

033

040

Recent methods for few-shot language model tuning obtain impressive performance but require careful engineering of prompts and verbalizers to convert inputs to a cloze-format (Taylor, 1953) that can be scored with pre-trained language models (PLMs) (Radford et al., 2018; Radford et al.; Brown et al., 2020; Schick and Schütze, 2021a,b). For example, as Figure 1 shows, a sentiment classifier can be designed by inserting the input text *x* in a *prompt template* "*x* It was [MASK]" where *verbalizers* (e.g., 'great' and 'terrible') are substituted for the [MASK] to score target task labels ('positive' or 'negative'). In this paper, we show that such engineering is not needed for few-shot learning and instead can







Figure 1: Existing few-shot fine-tuning methods require manual engineering to reduce new tasks to masked language modeling. PERFECT does not rely on any handcrafting, removing both patterns and verbalizers (see Figure 3).

be replaced with simple methods for data-efficient fine-tuning with as few as 32 end-task examples.

041

042

043

044

047

048

050

051

052

053

054

059

060

061

062

063

064

065

067

068

069

More specifically, we propose PERFECT, a Prompt-free and Efficient paRadigm for FEw-shot Cloze-based fine-Tuning. To remove handcrafted patterns, PERFECT uses *task-specific adapter layers* (Houlsby et al., 2019) (§3.1). Freezing the underlying PLM with millions or billions of parameters (Liu et al., 2019; Raffel et al., 2020), and only tuning adapters with very few new parameters saves on memory and storage costs (§4.2), while allowing very sample-efficient tuning (§4). It also stabilizes the training by increasing the worst-case performance and decreasing variance across the choice of examples in the few shot training sets (§4.3).

To remove handcrafted verbalizers (with variable token lengths), we introduce a new *multi-token fixed-length classifier scheme* that learns task label embeddings which are independent from the language model vocabulary during fine-tuning (§3.2). We show (§4) that this approach is sample efficient and outperforms carefully engineered verbalizers from *random initialization* (§4). It also allows us to avoid previously used expensive auto-regressive decoding schemes (Schick and Schütze, 2021b), by leveraging prototypical networks (Snell et al., 2017) over multiple tokens. Overall, these changes enable up to 100x faster learning and inference (§4.2).

Overall, PERFECT has several advantages: It avoids engineering patterns and verbalizers for each

new task, which can be cumbersome. Recent work 071 has shown that even some intentionally irrelevant 072 or misleading prompts can perform as well as more 073 interpretable ones (Webson and Pavlick, 2021). Unlike the zero-shot or extreme few-shot case, where prompting might be essential, we argue in this paper that all you need is tens of training examples to avoid 077 these challenges by adopting PERFECT or a similar data-efficient learning method. Experiments on a 079 wide variety of NLP tasks demonstrate that PERFECT outperforms state-of-the-art prompt-based methods while being significantly more efficient in inference and training time, storage, and memory usage (§4.2). To the best of our knowledge, we are the first to 084 propose a few-shot learning method with PLMs that successfully removes all per-task manual engineering.

## 2 Background

087

089

098

100

101

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

Problem formulation: We consider a general problem of fine-tuning language models in a few-shot setting, on a small training set with K unique classes and N examples per class, such that the total number of examples is  $|\mathcal{D}| = N \times K$ . Let  $\mathcal{D} = \bigcup_{k=1}^{K} \mathcal{D}_k$  be the given training set, where  $\mathcal{D}_k = \{(\boldsymbol{x_k^i}, y_k^i)\}_{i=1}^N$  shows the set of examples labeled with class k and  $y_k^i \in \mathcal{Y}$ is the corresponding label, where  $|\mathcal{Y}| = K$ . We additionally assume access to a development set with the same size as the training data. Note that larger validation sets can grant a substantial advantage (Perez et al., 2021), and thus it is important to use a limited validation size to be in line with the goal of few-shot learning. Unless specified otherwise, in this work, we use 16 training examples (N = 16) and a validation set with 16 examples, for a total of 32-shot learning.

#### 2.1 Adapters

Recent work has shown that fine-tuning *all* parameters of PLMs with a large number of parameters in low-resource datasets can lead to a sub-optimal solution (Peters et al., 2019; Dodge et al., 2020). As shown in Figure 2, Rebuffi et al. (2018) and Houlsby et al. (2019) suggest an efficient alternative, by inserting small task-specific modules called *adapters* within layers of a PLMs. They then only train the newly added adapters and layer normalization, while fixing the remaining parameters of a PLM.

Each layer of a transformer model is composed of two primary modules: a) an attention block, and b) a feed-forward block, where both modules are followed by a skip connection. As depicted in Figure 2, adapters are normally inserted after each of these blocks before the skip connection.



Figure 2: Left: Adapter integration in a PLM. Right: An adapter architecture. Adapters are usually inserted after the feed-forward and self-attention modules. During training, we only optimize the green components

Adapters are bottleneck architectures. By keeping input and output dimensions the same, they introduce no additional architectural changes. Each adapter,  $A(.) \in \mathbb{R}^{H}$ , consists of a down-projection,  $D(.) \in \mathbb{R}^{H \times B}$ , a non-linearity, such as GeLU (Hendrycks and Gimpel, 2016), and an up-projection  $U(.) \in \mathbb{R}^{B \times H}$ , where H is the dimension of input hidden states x, and B is the bottleneck size. Formally defined as:

$$A(\boldsymbol{x}) = U(\text{GeLU}(D(\boldsymbol{x}))) + \boldsymbol{x}, \quad (1)$$

121

122

123

124

125

127

128

139

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

## 2.2 Prompt-based Fine-tuning

**Standard Fine-tuning:** In standard fine-tuning with PLMs (Devlin et al., 2019), first a special [CLS] token is appended to the input x, and then the PLM maps it to a sequence of hidden representations  $h = (h_1, ..., h_S)$  with  $h_i \in \mathbb{R}^H$ , where H is the hidden dimension, and S is the maximum sequence length. Then, a classifier,  $softmax(W^Th_{[CLS]})$ , using the embedding of the classification token  $(h_{[CLS]})$ , is trained end-to-end for each downstream task. The main drawback of this approach is the discrepancy between the pre-training and fine-tuning phases since PLMs have been trained to *predict mask tokens* in a masked language modeling task (Devlin et al., 2019).

**Prompt-based tuning:** To address this discrepancy, *prompt-based fine-tuning* (Schick and Schütze, 2021a,b; Gao et al., 2021) formulates tasks in a cloze-format (Taylor, 1953). This way, the model can predict targets with a *masked language modeling (MLM) objective*. For example, as shown in Figure 1, for a sentiment classification task, inputs are converted to:

$$x_{\text{prompt}} = [\text{CLS}] x \cdot \underbrace{\text{It was}}_{\text{pattern}} [\text{MASK}] \cdot [\text{SEP}]$$
 152

Then, the PLM determines which verbalizer (e.g., 153 'great' and 'terrible') is the most likely substitute for 154 the mask in the  $x_{\text{prompt}}$ . This subsequently determines 155 the score of targets ('positive' or 'negative'). In detail: 156

157

161

163

196

197

198

199

**Training strategy:** Let  $\mathcal{M}: \mathcal{Y} \to \mathcal{V}$  be a mapping from target labels to individual words in a PLM's vocabulary. We refer to this mapping as verbalizers. 159 Then the input is converted to  $x_{\text{prompt}} = \mathcal{T}(x)$  by 160 appending a *pattern* and a *mask token* to x so that it has the format of a masked language modeling input. 162 Then, the classification task is converted to a MLM objective (Tam et al., 2021; Schick and Schütze, 164 2021a), and the PLM computes the probability of the 165 label y as:

167 
$$p(y|\boldsymbol{x}) = p([MASK] = \mathcal{M}(y)|\boldsymbol{x}_{prompt})$$
168 
$$= \frac{\exp(\boldsymbol{W}_{\mathcal{M}(y)}^{T}\boldsymbol{h}_{[MASK]})}{\sum_{v' \in \mathcal{V}} \exp(\boldsymbol{W}_{v'}^{T}\boldsymbol{h}_{[MASK]})}, \quad (2)$$

where  $h_{[{
m MASK}]}$  is the last hidden representation of the mask, and  $W_v$  shows the output embedding of the 170 PLM for each verbalizer  $v \in \mathcal{V}$ . For many tasks, ver-171 balizers have multiple tokens. Schick and Schütze 172 (2021b) extended (2) to multiple mask tokens by 173 adding the maximum number of mask tokens M174 needed to express the outputs (verbalizers) for a task. 175 In that case, Schick and Schütze (2021b) computes the 176 probability of each class as the summation of the log probabilities of each token in the corresponding verbal-178 izer, and then they add a hinge loss to ensure a margin 179 between the correct verbalizer and the incorrect ones.

**Inference strategy:** During inference, the model 181 needs to select which verbalizer to use in the given context. Schick and Schütze (2021b) predicts the verbalizer tokens in an autoregressive fashion. They first trim the number of mask tokens from M to each candidate verbalizer's token length, and compute the probability of each mask token. They then choose the predicted token with the highest probability and 188 replace the corresponding mask token. Conditioning 189 on this new token, the probabilities of the remaining 190 mask positions are recomputed. They repeat this autoregressive decoding until they fill all mask 192 positions. This inference strategy is very slow, as the 193 number of forward passes increases with the number 194 of classes and the number of verbalizer's tokens.

> This formulation obtained impressive few-shot performance with PLMs. However, the success of this approach heavily relies on engineering handcrafted patterns and verbalizers. Coming up with suitable verbalizers and patterns can be difficult (Mishra et al.,



Figure 3: We remove handcrafted patterns and verbalizers. We replace patterns using task-specific adapters and design label embeddings for the classes. We only train the green blocks (the label embeddings, adapters, and layer norms).

2021). Additionally, the performance is sensitive to the wording of patterns (Zhao et al., 2021; Perez et al., 2021; Schick and Schütze, 2021a; Jiang et al., 2020) or to the chosen verbalizers (Webson and Pavlick, 2021). 201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

221

In addition, handcrafted verbalizers cause problems for efficient training: a) they require updating the PLM embedding layer, causing large memory overhead; b) fine-tuning PLMs also requires a very small learning rate (usually  $10^{-5}$ ), which slows down tuning the parameters of the verbalizers; c) modeling verbalizers as one of the tokens of the PLM vocabulary (perhaps unintentionally) impacts the input representation during tuning; d) verbalizers have variable token lengths, complicating the implementation in a vectorized format, thereby making it challenging to efficiently fine-tune PLMs.

#### 3 Method

We propose PERFECT, a verbalizer and patterns-free few-shot learning method. We design PERFECT to be close to the pre-training phase, similar to the PET family of models (Schick and Schütze, 2021b; Gao et al., 2021), while replacing handcrafted patterns and

259

260

261

263

264

265

224

verbalizers with new components that are designed to describe the task and learn the labels. As shown in Figure 3, we first convert each input  $x_{input}$  to its masked language modeling (MLM) input containing M mask tokens [MASK]<sup>2</sup> with no added patterns, denoted as  $x_{masked} = \mathcal{T}'(x_{input})$ .<sup>3</sup> PERFECT then trains a classifier per-token and optimizes the average multi-class hinge loss over each mask position.

Three main components play a role in the success of PERFECT: a) a pattern-free task description, where we use task-specific adapters to efficiently tell the model about the given task, replacing previously manually engineered patterns (§3.1), b) multi-token label-embedding as an efficient mechanism to learn the label representations, removing manually designed verbalizers (§3.2). c) an efficient inference strategy building on top of the idea of prototypical networks (Snell et al., 2017) (§??), which replaces prior iterative autoregressive decoding methods (Schick and Schütze, 2021b).

As shown in Figure 3, we fix the underlying PLM model and only optimize the new parameters that we add (green boxes). This includes the task-specific adapters to adapt the representations for a given task and the multi-token label representations. We detail each of these components below.

#### 3.1 Pattern-Free Task Description

We use task-specific adapter layers to provide the model with learned, implicit task descriptions. Adapters additionally bring multiple other benefits: a) fine-tuning all weights of PLMs with millions or billions of parameters is sample-inefficient, and can be unstable in low-resource settings (Dodge et al., 2020); adapters allow sample-efficient fine-tuning, by keeping the underlying PLM fixed, b) adapters reduce the storage and memory footprints (§4.2), c) they also increase stability and performance (§4), making them an excellent choice for few-shot fine-tuning. To our knowledge, this is the first approach for using task-specific adapters to effectively and efficiently remove patterns in few-shot learning. Experimental results in §4 show its effectiveness compared to handcrafted patterns and soft prompts (Li and Liang, 2021; Lester et al., 2021).

#### 3.2 Multi-Token Label Embeddings

We freeze the weights of the PLM's embedding layer and introduce a separate label embedding  $L \in \mathbb{R}^{K \times M \times H}$ , which is a multi-token label representation where M is the number of tokens representing each label, K indicates the number of classes, H is the input hidden dimension. Using a fixed number of tokens M for each label, versus variable-token length verbalizers used in prior work (Schick and Schütze, 2021a,b) substantially simplifies the implementation and accelerates the training (§4.2). 267

268

269

270

271

272

274

275

276

277

279

281

282

283

285

286

287

288

289

291

293

294

296

299

300

301

302

303

306

307

308

310

311 **312** 

314

315

### 3.3 Training PERFECT

As shown in Figure 3, we optimize label embeddings so that the PLM predicts the correct label, and optimize adapters to adapt the PLM for the given task. For label embeddings, PERFECT trains a classifier per token and optimizes the average multi-class hinge loss over all mask positions. Given  $x_{masked}$ , let  $h_{[MASK]_i}$  be the embedding of its *i*-th mask token from the last layer of the PLM encoder. Additionally, let  $f(.) : \mathbb{R}^H \to \mathbb{R}^K$  be a per-token classifier that computes the predictions by multiplying the mask token embedding with its corresponding label embedding. Formally defined as:

$$\boldsymbol{t_i} = f(\boldsymbol{h}_{[\text{MASK}]_i}) = \boldsymbol{L_i^T} \boldsymbol{h}_{[\text{MASK}]_i}, \qquad (3)$$

where  $L_i \in \mathbb{R}^{K \times H}$  shows the label embedding for the *i*-th mask position. Then, for each mask position, we optimize a multi-class hinge loss between their scores  $t_i$  and labels. Formally defined as:

$$\mathcal{L}(\boldsymbol{x}, y, i) = \frac{\sum_{k=1, k \neq y}^{K} \max(0, m - \boldsymbol{t}_{iy} + \boldsymbol{t}_{ik})}{K}, \quad (4)$$

where  $t_{ik}$  shows the k-th element of  $t_i$ , representing the score corresponding to class k, and m is the margin, which we fix to the default value of m = 1. Then, the final loss is computed by averaging the loss over all mask tokens and training samples:

$$\mathcal{L} = \frac{1}{M|\mathcal{D}|} \sum_{(\boldsymbol{x}, y) \in \mathcal{D}} \sum_{i=1}^{M} \mathcal{L}(\boldsymbol{x}, y, i)$$
(5)

#### **3.4 Inference with PERFECT**

During evaluation, instead of relying on the prior iterative autoregressive decoding schemes (Schick and Schütze, 2021b), we classify a query point by finding the nearest class prototype to the mask token embeddings:

$$y = \underset{y \in \mathcal{Y}}{\operatorname{argmax}} \max_{i \in \{1, \dots, M\}} \left( \exp^{-d(\boldsymbol{h}_{i}^{\boldsymbol{q}}, \boldsymbol{c}_{iy})} \right), \quad (6)$$

where d is squared euclidean distance,<sup>4</sup>  $h_i^q$  indicates the embedding of the *i*-th mask position for the

<sup>&</sup>lt;sup>2</sup>We discuss the general case with inserting multiple masks; for some datasets this improves performance (§4.3.1).

 $<sup>^{3}</sup>$ We insert mask tokens after the input string in singlesentence benchmarks, and after the first sentence in the case of sentence-pair datasets and encode both sentences as a single input, which we found to perform the best (Appendix C).

<sup>&</sup>lt;sup>4</sup>We also tried with cosine similarity but found a slight improvement with squared Euclidean distance (Snell et al., 2017).

409

query sample q, and  $c_{iy} \in \mathbb{R}^D$  is the prototype representation of the *i*-th mask token with class label y, i.e., the mean embedding of *i*-th mask position in all training samples with label y:

$$\boldsymbol{c_{iy}} = \frac{1}{|\mathcal{D}_y|} \sum_{\boldsymbol{b} \in \mathcal{D}_y} \boldsymbol{h_i^b}, \tag{7}$$

where  $h_i^b$  shows the embedding of *i*-th mask position for training sample *b*, and  $\mathcal{D}_y$  is the training instances with class *y*. This strategy closely follows prototypical networks (Snell et al., 2017), but applied across multiple tokens. We choose this form of inference because prototypical networks are known to be sample efficient and robust (Snell et al., 2017), and because it substantially speeds up evaluation compared to prior methods (§4.2).

### 4 Experiments

316

317

319

320

322

323

324

328

329

333

334

335

337

340

341

345

347

349

351

We conduct extensive experiments on a variety of NLP datasets to evaluate the performance of PERFECT and compare it with state-of-the-art few-shot learning.

**Datasets:** We consider 7 tasks and 12 datasets: 1) the sentiment analysis datasets SST-2 (Socher et al., 2013), SST-5 (Socher et al., 2013), MR (Pang and Lee, 2005), and CR (Hu and Liu, 2004), 2) the subjectivity classification dataset SUBJ (Pang and Lee, 2004), 3) the question classification dataset TREC (Voorhees and Tice, 2000), 4) the natural language inference datasets CB (De Marneffe et al., 2019) and RTE (Wang et al., 2019a), 5) the question answering dataset QNLI (Rajpurkar et al., 2016), 6) the word sense disambiguation dataset WiC (Pilehvar and Camacho-Collados, 2019), 7) the paraphrase detection datasets MRPC (Dolan and Brockett, 2005) and OOP.<sup>5</sup> See datasets statistics in Appendix A.

For MR, CR, SST-5, SUBJ, and TREC, we test on the original test sets, while for other datasets, since test sets are not publicly available, we test on the original validation set. We sample 16 instances per label from the training set to form training and validation sets.

**Baselines** We compare with the state-of-the-art few-shot learning of PET and fine-tuning:

**PET** (Schick and Schütze, 2021a,b) is the stateof-the-art few-shot learning method that employs carefully crafted verbalizers and patterns. We report the best (PET-best) and average (PET-average) results among all patterns and verbalizers.<sup>6</sup> **FINETUNE** The standard fine-tuning (Devlin et al., 2019), with adding a classifier on top of the [CLS] token and fine-tuning all parameters.

**Our method** We study the performance of PERFECT and perform an extensive ablation study to show the effectiveness of our design choices:

**PERFECT-rand** We randomly initialize the label embedding L from a normal distribution  $\mathcal{N}(0,\sigma)$  with  $\sigma = 10^{-4}$  (chosen based on validation performance, see Appendix D) without relying on any handcrafted patterns and verbalizers. As an ablation, we study the following two variants:

**PERFECT-init** We initialize the label embedding with the token embeddings of manually designed verbalizers in the PLM's vocabulary to study the impact of engineered verbalizers.

**PERFECT-prompt** We compare using adapters versus soft prompt-tuning with removing adapters and appending trainable continuous prompt embeddings to the input (Lester et al., 2021). Then, we only tune the soft prompt and the label embedding.

**Experimental details:** We use the RoBERTa large model (Liu et al., 2019) (355M parameters) as the underlying PLM for all methods. We use the Hugging-Face PyTorch implementation (Wolf et al., 2020). For the baselines, we used the carefully manually designed patterns and verbalizers in Gao et al. (2021), Min et al. (2021), and Schick and Schütze (2021b) (usually 5 different options per datasets; see Appendix B).

We evaluate all methods using 5 different random samples to create the training/validation sets and 4 different random seeds for training. Therefore, for PET-average, we report the results on 20 x 5 (number of patterns and verbalizers) = 100 runs, while for PET-best and our method, we report the results over 20 runs. The variance in few-shot learning methods is usually high (Perez et al., 2021; Zhao et al., 2021; Lu et al., 2021). Therefore, we report average, worst-case performance, and standard deviation across all runs, where the last two values can be important for risk-sensitive applications (Asri et al., 2016).

#### 4.1 Experimental Results

Table 1 shows the performance of all methods. PERFECT obtains state-of-the-art results, improving the performance compared to PET-average by +1.1 and +4.6 points for single-sentence and sentence-pair datasets respectively. It even outperforms PET-best, where we report the best performance of PET across multiple manually engineered patterns and verbalizers.

<sup>&</sup>lt;sup>5</sup>https://quoradata.quora.com/

<sup>&</sup>lt;sup>6</sup>For a controlled study, we use the MLM variant shown in (2), which has been shown to perform the best (Tam et al., 2021).

Method	SST-2	CR	MR	SST-5	Subj	TREC	Avg
Single-Sentence Benchmarks							
Finetune	81.4/70.0/4.0	80.1/72.9/4.1	77.7/66.8/4.6	39.2/34.3/2.5	90.2/84.1/1.8	87.6/75.8/3.7	76.0/67.3/3.4
PET-Average	89.7/81.0/2.4	88.4/68.8/3.0	85.9/79.0/2.1	45.9/40.3/2.4	88.1/79.6/2.4	85.0/70.6/4.5	80.5/69.9/2.8
PET-Best	89.1/81.0/2.6	88.8/85.8/1.9	86.4/82.0/1.6	46.0/41.2/2.4	88.7/84.6/1.8	85.8/70.6/4.4	80.8/74.2/2.4
PERFECT-rand	90.7/88.2/1.2	<b>90.0</b> /85.5/1.4	86.3/81.4/1.6	42.7/35.1/2.9	89.1/82.8/2.1	<b>90.6</b> /81.6/3.2	81.6/75.8/2.1
			Ablation	ı			
PERFECT-init	<b>90.9</b> /87.6/1.5	89.7/87.4/1.2	85.4/75.8/3.3	42.8/35.9/3.5	87.6/81.6/2.8	90.4/86.6/1.8	81.1/75.8/2.4
PERFECT-prompt	70.6/56.0/8.3	71.0/55.8/8.2	66.6/49.6/7.3	32.2/26.5/3.2	82.7/69.6/3.9	79.6/66.8/6.5	67.1/54.0/6.2
Method	СВ	RTE	QNLI	MRPC	QQP	WiC	Avg
Sentence-Pair Benchmarks							
		Senter	ice-Pair De	ncnmarks			
FINETUNE	72.9/67.9/2.5	56.8/50.2/3.5	62.7/51.4/7.0	<b>70.1</b> /62.7/4.7	65.0/59.8/3.6	52.4/46.1/3.7	63.3/56.4/4.2
FINETUNE PET-Average	72.9/67.9/2.5 86.9/73.2/5.1	56.8/50.2/3.5 60.1/49.5/4.7	62.7/51.4/7.0 66.5/55.7/6.2	<b>70.1</b> /62.7/4.7 62.1/38.2/6.8	65.0/59.8/3.6 63.4/44.7/7.9	52.4/46.1/3.7 51.0/46.1/2.6	63.3/56.4/ <b>4.2</b> 65.0/51.2/5.6
FINETUNE PET-Average PET-Best	72.9/67.9/ <b>2.5</b> 86.9/73.2/5.1 90.0/78.6/3.9	<b>Senter</b> 56.8/50.2/3.5 60.1/49.5/4.7 <b>62.3</b> /51.3/4.5	62.7/51.4/7.0 66.5/55.7/6.2 70.5/57.9/6.4	<b>70.1/62.7/4.7</b> 62.1/38.2/6.8 63.4/49.3/6.5	65.0/59.8/3.6 63.4/44.7/7.9 70.7/55.2/5.8	52.4/46.1/3.7 51.0/46.1/2.6 51.6/ <b>47.2/2.3</b>	63.3/56.4/ <b>4.2</b> 65.0/51.2/5.6 68.1/56.6/4.9
FINETUNE PET-Average PET-Best PERFECT-rand	72.9/67.9/2.5 86.9/73.2/5.1 90.0/78.6/3.9 <b>90.3/83.9/</b> 3.5	<b>Senter</b> 56.8/50.2/3.5 60.1/49.5/4.7 <b>62.3</b> /51.3/4.5 60.4/53.1/4.7	62.7/51.4/7.0 66.5/55.7/6.2 70.5/57.9/6.4 74.1/60.3/4.6	<b>70.1</b> /62.7/4.7 62.1/38.2/6.8 63.4/49.3/6.5 67.8/54.7/5.7	65.0/59.8/3.6 63.4/44.7/7.9 70.7/55.2/5.8 <b>71.2</b> /64.2/ <b>3.5</b>	52.4/46.1/3.7 51.0/46.1/2.6 51.6/47.2/2.3 53.8/47.0/3.0	63.3/56.4/4.2 65.0/51.2/5.6 68.1/56.6/4.9 <b>69.6/60.5/4.2</b>
FINETUNE PET-Average PET-Best PERFECT-rand	72.9/67.9/2.5 86.9/73.2/5.1 90.0/78.6/3.9 <b>90.3/83.9/</b> 3.5	<b>Senter</b> 56.8/50.2/3.5 60.1/49.5/4.7 <b>62.3</b> /51.3/4.5 60.4/53.1/4.7	62.7/51.4/7.0 66.5/55.7/6.2 70.5/57.9/6.4 74.1/60.3/4.6 Ablation	<b>70.1</b> /62.7/4.7 62.1/38.2/6.8 63.4/49.3/6.5 67.8/54.7/5.7	65.0/59.8/3.6 63.4/44.7/7.9 70.7/55.2/5.8 <b>71.2</b> /64.2/ <b>3.5</b>	52.4/46.1/3.7 51.0/46.1/2.6 51.6/47.2/2.3 53.8/47.0/3.0	63.3/56.4/ <b>4.2</b> 65.0/51.2/5.6 68.1/56.6/4.9 <b>69.6/60.5/4.2</b>
FINETUNE PET-Average PET-Best PERFECT-rand PERFECT-init	72.9/67.9/2.5 86.9/73.2/5.1 90.0/78.6/3.9 <b>90.3/83.9</b> /3.5 87.9/75.0/4.9	<b>Senter</b> 56.8/50.2/3.5 60.1/49.5/4.7 <b>62.3</b> /51.3/4.5 60.4/53.1/4.7 60.7/52.7/4.5	62.7/51.4/7.0 66.5/55.7/6.2 70.5/57.9/6.4 74.1/60.3/4.6 Ablation 72.8/56.7/6.8	<b>70.1</b> /62.7/4.7 62.1/38.2/6.8 63.4/49.3/6.5 67.8/54.7/5.7 <b>6</b> 65.9/56.6/6.0	65.0/59.8/3.6 63.4/44.7/7.9 70.7/55.2/5.8 71.2/64.2/3.5 71.1/65.6/3.5	52.4/46.1/3.7 51.0/46.1/2.6 51.6/47.2/2.3 53.8/47.0/3.0 51.7/46.6/2.8	63.3/56.4/4.2 65.0/51.2/5.6 68.1/56.6/4.9 69.6/60.5/4.2 68.4/58.9/4.8

Table 1: Performance of all methods on single-sentence and sentence-pair benchmarks. We report average/worst-case accuracy/standard deviation. PERFECT obtains the state-of-the-art results. Bold fonts indicate the best results.

Moreover, PERFECT improves the minimum performance and reduces standard deviation substantially. Finally, PERFECT is also significantly more efficient: reducing the training and inference time, memory usage, and storage costs (see §4.2).

410

411

412 413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

PET-best improves the results over PET-average showing that PET is unstable to the choice of patterns and verbalizers; this difference is more highlighted for sentence-pair benchmarks. This can be because the position of the mask highly impacts the results, and patterns used in sentence-pair datasets in Schick and Schütze (2021b) well leverages this by putting the mask in multiple locations (see Appendix B).

As an ablation, even if we initialize the label embedding with handcrafted verbalizers, PER-FECT-init consistently obtains lower performance, demonstrating that PERFECT is able to obtain state-ofthe-art performance with learning from *pure random initialization*. We argue that initializing randomly close to zero (with low variance  $\sigma = 10^{-4}$ ), as done in our case, slightly improves performance, which perhaps is not satisfied when initializing from the manually engineered verbalizers (see Appendix D).

As a second ablation, when learning patterns with optimizing soft prompts in PERFECT-prompt, we observe high sensitivity to learning rate, as also confirmed in Li and Liang (2021) and Mahabadi et al. (2021a). We experimented with multiple learning rates but performance consistently lags behind PERFECT-rand. This can be explained by the low flexibility of such methods as all the information regarding specifying patterns needs to be contained in the prefixes. As a result, the method only allows limited interaction with the rest of the model parameters, and obtaining good performance requires very large models (Lester et al., 2021). In addition, increasing the sequence length leads to memory overhead (Mahabadi et al., 2021a), and the number of prompt tokens is capped by the number of tokens that can fit in the maximum input length, which can be a limitation for tasks requiring large contexts.

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

#### 4.2 Efficiency Evaluation

In this section, we compare the efficiency of PERFECT with the state-of-the-art few-shot learning method, PET. To this end, we train all methods for ten epochs on the 500-sampled QNLI dataset. We select the largest batch size for each method that fits a fixed budget of the GPU memory (40 GB).

Due to the auto-regressive inference strategy of PET (Schick and Schütze, 2021b), every prior work implemented it with a batch size of 1 (Perez et al.,

Metric	РЕТ	PERFECT	$\Delta\%$
Trained params (M)	355.41	3.28	-99.08%
Peak memory (GB)	20.93	16.34	-21.93%
Training time (min)	23.42	0.65	-97.22%
+ PET in batch	0.94	0.65	-30.85%
Inference time (min)	9.57	0.31	-96.76%

Table 2: Percentage of trained parameters, average peak memory, training, and inference time.  $\Delta\%$  is the relative difference with respect to PET. Lower is better.

2021; Schick and Schütze, 2021b; Tam et al., 2021). Additionally, since PET deals with verbalizers of variable lengths, it is hard to implement their training phase in a batch mode. We specifically choose QNLI to have verbalizers of the same length and enable batching for comparison purposes (referred to as *PET in batch*). However, verbalizers are still not of the fixed-length for most other tasks, and this speed-up does not apply generally to PET.

In Table 2, we report the percentage of trained parameters, memory usage of each method, training, and inference time. PERFECT reduces the number of trained parameters, and therefore the storage requirement, by 99.08%. It additionally reduces the memory requirement by 21.93% compared to PET. PERFECT speeds up training substantially, by 97.22% relative to the original PET's implementation, and 30.85% to our implementation of PET. This is because adapter-based tuning saves on memory and allows training with larger batch sizes. On the other hand, PERFECT is significantly faster during inference time (96.76% less inference time relative to PET).

Overall, given the size of PLMs with millions and billions of parameters (Liu et al., 2019; Raffel et al., 2020), efficient few-shot learning methods are of paramount importance for practical applications. PER-FECT not only outperforms the state-of-the-art in terms of accuracy and stability (Table 1), but also is significantly more efficient in runtime, storage, and memory.

#### 4.3 Analysis

**Can task-specific adapters replace manually engineered patterns?** PERFECT is a pattern-free approach and employs adapters to provide the PLMs with task descriptions implicitly. In this section, we study the contribution of replacing manual patterns with adapters in isolation without considering our other contributions in representing labels, training, and inference. In PET (Schick and Schütze, 2021a,b), we replace the handcrafted patterns with task-specific adapters (*Pattern-Free*) while keeping the verbalizers

Dataset	PET-Average	Pattern-Free
SST-2	89.7/81.0/2.4	90.5/87.8/1.2
CR	88.4/68.8/3.0	89.8/87.0/1.4
MR	85.9/79.0/2.1	86.4/83.0/1.8
SST-5	45.9/40.3/2.4	44.8/40.0/2.4
SUBJ	88.1/79.6/2.4	85.3/74.7/3.8
TREC	85.0/70.6/4.5	87.9/84.6/1.8
CB	86.9/73.2/5.1	93.0/89.3/1.9
RTE	60.1/49.5/4.7	63.7/56.3/4.1
QNLI	66.5/55.7/6.2	71.3/65.8/2.5
MRPC	62.1/38.2/6.8	66.0/54.4/5.6
QQP	63.4/44.7/7.9	71.8/64.3/3.7
WiC	51.0/46.1/2.6	53.7/50.3/2.0
Avg	72.8/60.6/4.2	75.4/69.8/2.7

Table 3: Average performance of *PET* with five different patterns vs. *Pattern-Free* that replaces handcrafted patterns with task-specific adapters. We report the average/worst-case performance/and the standard deviation.

and the training and inference intact<sup>7</sup> and train it with a similar setup as in §4. Table 3 shows the results. While PET is very sensitive to the choice of prompts, adapters provide an efficient alternative to learn patterns robustly by improving the performance (average and worst-case) and reducing the standard deviation. This finding demonstrates that task-specific adapters can effectively replace manually engineered prompts. Additionally, they also save on the training budget by at least 1/number of patterns (normally 1/5) by not requiring running the method for different choices of patterns, and by freezing most parameters, this saves on memory and offers additional speed-up.

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

### 4.3.1 Ablation Study

**Impact of Removing Adapters** To study the impact of adapters in learning patterns, we remove adapters, while keeping the label embedding. Handcrafted patterns are not included and we tune all parameters of the model. Table 4 shows the results. Adding adapters for learning patterns contributes to the performance by improving the average performance, and making the model robust by improving the minimum performance and reducing the standard deviation. This is because training PLMs with millions of parameters is sample-inefficient and unstable on resource-limited datasets (Dodge et al., 2020; Zhang et al., 2020). However, by using adapters, we substantially reduce the number of trainable parameters, allowing the model to be better

488

489

490

491

492

493

494

495

496

497

498

499

461

462

<sup>&</sup>lt;sup>7</sup>Since we don't have patterns, in the case of multiple sets of verbalizers, we use the first set of verbalizers as a random choice.

Dataset	PERFECT	-Adapters
SST-2	90.7/88.2/1.2	88.2/81.9/2.3
CR	90.0/85.5/1.4	89.2/83.1/1.7
MR	86.3/81.4/1.6	82.5/78.2/2.5
SST-5	42.7/35.1/2.9	40.6/33.6/3.3
SUBJ	89.1/82.8/2.1	89.7/85.0/1.9
TREC	90.6/81.6/3.2	89.8/74.2/4.3
CB	<b>90.3/83.9</b> /3.5	89.6 <b>/83.9/2.8</b>
RTE	60.4/53.1/ <b>4.7</b>	<b>61.7/53.8</b> /5.1
QNLI	74.1/60.3/4.6	73.2/56.3/5.8
MRPC	67.8 <b>/54.7/5.7</b>	<b>68.0</b> /54.2/6.1
QQP	71.2/64.2/3.5	71.0/62.0/3.7
WiC	53.8/47.0/3.0	52.5/46.9/3.0
Avg	75.6/68.1/3.1	74.7/66.1/3.5

Table 4: Performance of PERFECT w/o adapters, *-Adapters*. We report the average performance/worst-case performance/and the standard deviation.

tuned in a few-shot setting.

530

531

532

533

534

535

536

537

538

539

540

541

542

545

546

547

548

550

551

552

554

555 556 **Impact of the number of masks** In Table 1, to compare our design with PET in isolation, we fixed the number of mask tokens as the maximum number inserted by PET. In table 5, we study the impact of varying the number of inserted mask tokens for a random selection of six tasks. For most tasks, having two mask tokens performs the best, while for MR and RTE, having one, and for MRPC, inserting ten masks improves the results substantially. The number of required masks might be correlated with the difficulty of tasks. PERFECT is designed to be general, enabling having multiple mask tokens.

#### 5 Related Work

Adapter Layers: Mahabadi et al. (2021b) and Üstün et al. (2020) proposed to generate adapters' weights using hypernetworks (Ha et al., 2017), where Mahabadi et al. (2021b) proposed to share a small hypernetwork to generate conditional adapter weights efficiently for each transformer layer and task. Mahabadi et al. (2021a) proposed compacter layers by building on top of ideas of parameterized hypercomplex layers (Zhang et al., 2021) and low-rank methods (Li et al., 2018; Aghajanyan et al., 2021), as an efficient fine-tuning method for PLMs. We are the first to employ adapters to replace handcrafted patterns for few-shot learning.

557 Few-shot Learning with PLMs: Researchers con558 tinuously tried to address the challenges of manually
559 engineered patterns and verbalizers: a) Learning the

Datasets	1	2	5	10
CR	90.1	90.2	89.0	87.8
MR	86.9	86.1	85.4	85.6
MRPC	67.4	68.2	70.1	72.3
QNLI	73.7	73.9	73.0	65.1
RTE	60.0	57.3	56.2	56.0
TREC	90.0	90.9	88.9	88.8
Avg	78.0	77.8	77.1	75.9

Table 5: Test performance for the varying number of mask tokens. Bold fonts indicate the best results in each row.

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

590

591

592

593

patterns in a continuous space (Li and Liang, 2021; Qin and Eisner, 2021; Lester et al., 2021), while freezing PLM for efficiency, has the problem that, in most cases, such an approach only works with very large scale PLMs (Lester et al., 2021), and lags behind full fine-tuning in a general setting, while being inefficient and not as effective compared to adapters (Mahabadi et al., 2021a). b) Optimizing patterns in a discrete space (Shin et al., 2020; Jiang et al., 2020; Gao et al., 2021) has the problem that such methods are computationally costly. c) Automatically finding verbalizers in a discrete way (Schick et al., 2020; Schick and Schütze, 2021a) is computationally expensive and does not perform as well as manually designed ones. d) Removing manually designed patterns (Logan IV et al., 2021) substantially lags behind the expertdesigned ones. Our proposed method, PERFECT, does not rely on any handcrafted patterns and verbalizers.

## 6 Conclusion

We proposed PERFECT, a simple and efficient method for few-shot learning with pre-trained language models without relying on handcrafted patterns and verbalizers. PERFECT employs task-specific adapters to learn task descriptions implicitly, replacing previous handcrafted patterns, and a continuous multi-token label embedding to represent the output classes. Through extensive experiments over 12 NLP benchmarks, we demonstrate that PERFECT, despite being far simpler and more efficient than recent few-shot learning methods, produces state-of-the-art results. Overall, the simplicity and effectiveness of PERFECT make it a promising approach for few-shot learning with PLMs.

## References

Armen Aghajanyan, Luke Zettlemoyer, and Sonal Gupta.5942021. Intrinsic dimensionality explains the effectiveness595

- 596 601 608 610 611 612 614 615 618 619 622 623 627 628

633 634

637

641

646

of language model fine-tuning. ACL.

- Hiba Asri, Hajar Mousannif, Hassan Al Moatassime, and Thomas Noel. 2016. Using machine learning algorithms for breast cancer risk prediction and diagnosis. Procedia Computer Science.
- Roy Bar-Haim, Ido Dagan, Bill Dolan, Lisa Ferro, and Danilo Giampiccolo. 2006. The second pascal recognising textual entailment challenge. Second PASCAL Challenges Workshop on Recognising Textual Entailment.
- Luisa Bentivogli, Ido Dagan, Hoa Trang Dang, Danilo Giampiccolo, and Bernardo Magnini. 2009. The fifth pascal recognizing textual entailment challenge. In TAC.
  - Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In NeurIPS.
- Ido Dagan, Oren Glickman, and Bernardo Magnini. 2005. The pascal recognising textual entailment challenge. In Machine Learning Challenges Workshop.
  - Marie-Catherine De Marneffe, Mandy Simons, and Judith Tonhauser. 2019. The commitmentbank: Investigating projection in naturally occurring discourse. In proceedings of Sinn und Bedeutung.
  - Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In NAACL.
- Jesse Dodge, Gabriel Ilharco, Roy Schwartz, Ali Farhadi, Hannaneh Hajishirzi, and Noah Smith. 2020. Fine-tuning pretrained language models: Weight initializations, data orders, and early stopping. arXiv preprint arXiv:2002.06305.
- William B Dolan and Chris Brockett. 2005. Automatically constructing a corpus of sentential paraphrases. In IWP.
- Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. Making pre-trained language models better few-shot learners. ACL.
- Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and Bill Dolan. 2007. The third PASCAL recognizing textual entailment challenge. In ACL-PASCAL Workshop on Textual Entailment and Paraphrasing.
- David Ha, Andrew Dai, and Quoc V. Le. 2017. Hypernetworks. In ICLR.
  - Dan Hendrycks and Kevin Gimpel. 2016. Gaussian error linear units (gelus). arXiv preprint arXiv:1606.08415.

Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for nlp. In ICML.

648

649

650

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

702

- Minging Hu and Bing Liu. 2004. Mining and summarizing customer reviews. In SIGKDD.
- Zhengbao Jiang, Frank F Xu, Jun Araki, and Graham Neubig. 2020. How can we know what language models know? TACL.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. EMNLP.
- Quentin Lhoest, Albert Villanova del Moral, Patrick von Platen, Thomas Wolf, Mario Šaško, Yacine Jernite, Abhishek Thakur, Lewis Tunstall, Suraj Patil, Mariama Drame, Julien Chaumond, Julien Plu, Joe Davison, Simon Brandeis, Victor Sanh, Teven Le Scao, Kevin Canwen Xu, Nicolas Patry, Steven Liu, Angelina McMillan-Major, Philipp Schmid, Sylvain Gugger, Nathan Raw, Sylvain Lesage, Anton Lozhkov, Matthew Carrigan, Théo Matussière, Leandro von Werra, Lysandre Debut, Stas Bekman, and Clément Delangue. 2021a. huggingface/datasets: 1.15.1.
- Quentin Lhoest, Albert Villanova del Moral, Yacine Jernite, Abhishek Thakur, Patrick von Platen, Suraj Patil, Julien Chaumond, Mariama Drame, Julien Plu, Lewis Tunstall, Joe Davison, Mario Šaško, Gunjan Chhablani, Bhavitvya Malik, Simon Brandeis, Teven Le Scao, Victor Sanh, Canwen Xu, Nicolas Patry, Angelina McMillan-Major, Philipp Schmid, Sylvain Gugger, Clément Delangue, Théo Matussière, Lysandre Debut, Stas Bekman, Pierric Cistac, Thibault Goehringer, Victor Mustar, François Lagunas, Alexander Rush, and Thomas Wolf. 2021b. Datasets: A community library for natural language processing. In EMNLP.
- Chunyuan Li, Heerad Farkhoor, Rosanne Liu, and Jason Yosinski. 2018. Measuring the intrinsic dimension of objective landscapes. In ICLR.
- Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. arXiv preprint arXiv:2101.00190.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.
- Robert L Logan IV, Ivana Balažević, Eric Wallace, Fabio Petroni, Sameer Singh, and Sebastian Riedel. 2021. Cutting down on prompts and parameters: Simple few-shot learning with language models. arXiv preprint arXiv:2106.13353.
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2021. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. arXiv preprint arXiv:2104.08786.

- Timo Schick, Helmut Schmid, and Hinrich Schütze. 2020. 756 Automatically identifying words that can serve as labels for few-shot text classification. In COLING. 758 Timo Schick and Hinrich Schütze. 2021a. Exploiting 759 cloze-questions for few-shot text classification and 760 natural language inference. In EACL. 761 Timo Schick and Hinrich Schütze. 2021b. It's not just size 762 that matters: Small language models are also few-shot 763 learners. In NAACL. Karin Kipper Schuler. 2005. Verbnet: A broad-coverage, 765 comprehensive verb lexicon. PhD Thesis. 766 arXiv Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric 767 Wallace, and Sameer Singh. 2020. Eliciting knowledge 768 from language models using automatically generated 769 prompts. In EMNLP. 770 Jake Snell, Kevin Swersky, and Richard Zemel. 2017. Pro-771 totypical networks for few-shot learning. In NeurIPS. Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In EMNLP. Derek Tam, Rakesh R Menon, Mohit Bansal, Shashank 777 Srivastava, and Colin Raffel. 2021. Improving and 778 simplifying pattern exploiting training. arXiv preprint 779 arXiv:2103.11955. 780 Wilson L Taylor. 1953. "cloze procedure": A new tool for 781 measuring readability. Journalism quarterly. 782 Ahmet Üstün, Arianna Bisazza, Gosse Bouma, and Gertjan 783 van Noord. 2020. Udapter: Language adaptation for 784 truly universal dependency parsing. In EMNLP. 785 Ellen M Voorhees and Dawn M Tice. 2000. Building a 786 question answering test collection. In SIGIR. Alex Wang, Yada Pruksachatkun, Nikita Nangia, Aman-788 preet Singh, Julian Michael, Felix Hill, Omer Levy, 789 and Samuel R Bowman. 2019a. Superglue: a stickier 790 benchmark for general-purpose language understanding 791 systems. In NeurIPS. 792 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, 793 Omer Levy, and Samuel R. Bowman. 2019b. GLUE: A 794 multi-task benchmark and analysis platform for natural 795 language understanding. In ICLR. 796 Albert Webson and Ellie Pavlick. 2021. Do prompt-based 797 models really understand the meaning of their prompts? 798 arXiv preprint arXiv:2109.01247. 799 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien 800 Chaumond, Clement Delangue, Anthony Moi, Pierric 801 Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, 802 Joe Davison, Sam Shleifer, Patrick von Platen, Clara 803 Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le 804 Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-art natural language processing. 807 In EMNLP: System Demonstrations. 808
- Rabeeh Karimi Mahabadi, James Henderson, and Sebastian Ruder. 2021a. Compacter: Efficient low-rank hypercomplex adapter layers. In *NeurIPS*.

705

707

710

711

712

713

714

715

716

717

718

719

720

721

723

726

727

731

735

736

737

738

740

741

742

743

744

745

746

747

750

751

752

753

- Rabeeh Karimi Mahabadi, Sebastian Ruder, Mostafa Dehghani, and James Henderson. 2021b. Parameterefficient multi-task fine-tuning for transformers via shared hypernetworks. In *ACL*.
- George A Miller. 1995. Wordnet: a lexical database for english. *Communications of the ACM*.
- Sewon Min, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2021. Noisy channel language model prompting for few-shot text classification. *arXiv preprint arXiv:2108.04106*.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2021. Cross-task generalization via natural language crowdsourcing instructions.
- Bo Pang and Lillian Lee. 2004. A sentimental education: sentiment analysis using subjectivity summarization based on minimum cuts. In *ACL*.
- Bo Pang and Lillian Lee. 2005. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *ACL*.
- Ethan Perez, Douwe Kiela, and Kyunghyun Cho. 2021. True few-shot learning with language models. *NeurIPS*.
- Matthew E Peters, Sebastian Ruder, and Noah A Smith. 2019. To tune or not to tune? adapting pretrained representations to diverse tasks. In *RepL4NLP*.
- Jonas Pfeiffer, Aishwarya Kamath, Andreas Rückle, Cho Kyunghyun, and Iryna Gurevych. 2021. AdapterFusion: Non-destructive task composition for transfer learning. In *EACL*.
- Mohammad Taher Pilehvar and Jose Camacho-Collados. 2019. Wic: the word-in-context dataset for evaluating context-sensitive meaning representations. In *NAACL*.
- Guanghui Qin and Jason Eisner. 2021. Learning how to ask: Querying lms with mixtures of soft prompts. In *NAACL*.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *JMLR*.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. In *EMNLP*.
- Sylvestre-Alvise Rebuffi, Hakan Bilen, and Andrea Vedaldi. 2018. Efficient parametrization of multidomain deep neural networks. In *CVPR*.

Aston Zhang, Yi Tay, SHUAI Zhang, Alvin Chan,
Anh Tuan Luu, Siu Hui, and Jie Fu. 2021. Beyond fully-
connected layers with quaternions: Parameterization
of hypercomplex multiplications with 1/n parameters.
In ICLR.

- Tianyi Zhang, Felix Wu, Arzoo Katiyar, Kilian Q Weinberger, and Yoav Artzi. 2020. Revisiting few-sample
  bert fine-tuning. In *ICLR*.
- 817 Tony Z. Zhao, Eric Wallace, Shi Feng, Dan Klein, and
  818 Sameer Singh. 2021. Calibrate before use: Improving
  819 few-shot performance of language models. *ICML*.

Dataset	Task	#Train	#Test	K	
	Single-Sentence Benchmarks				
MR	Sentiment analysis	8662	2000	2	
CR	Sentiment analysis	1774	2000	2	
SST-2	Sentiment analysis	6920	872	2	
SST-5	Sentiment analysis	8544	2210	5	
SUBJ	Subjectivity classification	8000	2000	2	
TREC	Question classification	5452	500	6	
	Sentence-Pair Benchm	arks			
СВ	Natural language inference	250	56	3	
RTE	Natural language inference	2490	277	2	
WiC	Word sense disambiguation	5428	638	2	
MRPC	Paraphrase detection	3668	408	2	
QNLI	Question answering	104743	5463	2	
QQP	Paraphrase detection	363846	40430	2	

Table 6: Statistics of datasets used in this work. We sample  $N \times |\mathcal{Y}|$  instances (with multiple seeds) from the original training set to form the few-shot training and validation sets. The test column shows the size of the test set.

#### **A** Experimental Details

820

821

822

825

827

832

833

834

837

838

839

841

844

845

847

**Datasets** Table 6 shows the stastistics of the datasets used. We download SST-2, MR, CR, SST-5, and SUBJ from Gao et al. (2021), while the rest of the datasets are downloaded from the HuggingFace Datasets library (Lhoest et al., 2021b,a). RTE, CB, WiC datasets are from SuperGLUE benchmark (Wang et al., 2019a), while QQP, MRPC and QNLI are from GLUE benchmark (Wang et al., 2019a) with Creative Commons license (CC BY 4.0). RTE (Wang et al., 2019a) is a combination of data from RTE1 (Dagan et al., 2005), RTE2 (Bar-Haim et al., 2006), RTE3 (Giampiccolo et al., 2007), and RTE5 (Bentivogli et al., 2009). For WiC (Pilehvar and Camacho-Collados, 2019) sentences are selected from VerbNet (Schuler, 2005), WordNet (Miller, 1995), and Wiktionary.

**Computing infrastructure** We run all the experiments on one NVIDIA A100 with 40G of memory.

**Training hyper-parameters** We set the maximum sequence length based on the recommended values in the HuggingFace repository (Wolf et al., 2020) and prior work (Min et al., 2021; Schick and Schütze, 2021b), i.e., we set it to 256 for SUBJ, CR, CB, RTE, and WiC, and 128 for other datasets. For all methods, we use a batch size of 32. For FINETUNE and PET, we use the default learning rate of  $10^{-5}$ , while for our method, as required by adapter-based methods (Mahabadi et al., 2021a), we set the learning rate to

a higher value of  $10^{-4}$ .<sup>8</sup> Through all experiments, 848 we fix the adapter bottleneck size to 64. Following 849 Pfeiffer et al. (2021), we experimented with keeping 850 one of the adapters in each layer for better training 851 efficiency and found keeping the adapter after the 852 feed-forward module in each layer to perform the best. 853 For tuning label embedding, we use the learning rate 854 of  $\{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$  and choose the 855 one obtaining the highest validation performance. For 856 PERFECT-prompt, we tune the continuous prompt 857 for learning rate of  $\{10^{-1}, 10^{-2}, 10^{-3}\}$ . Following 858 Lester et al. (2021), for PERFECT-prompt, we set 859 the number of prompt tokens to 20, and initialize 860 them with a random subset of the top 5000 token's 861 embedding of the PLM. We train all methods for 862 6000 steps. Based on our results, this is sufficient to 863 allow the models to converge. We save a checkpoint 864 every 100 steps for all methods and report the results 865 for the hyper-parameters performing the best on the 866 validation set for each task.

## **B** Choice of Patterns and Verbalizers

For SST-2, MR, CR, SST-5, and TREC, we used 4 different patterns and verbalizers from Gao et al. (2021). For CB, WiC, RTE datasets, we used the designed patterns and verbalizers in Schick and Schütze (2021b). For QQP, MRPC, and QNLI, we wrote the patterns and verbalizers inspired by the ones in Schick and Schütze (2021b). The used patterns and verbalizers are as follows:

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

• For sentiment analysis tasks (MR, CR, SST-2, SST-5), given a sentence *s*:

s A <MASK> one. s It was <MASK>. s All in all <MASK>. s A <MASK> piece.

with "great" as a verbalizer for positive, "terrible" for negative. In case of SST-5 with five labels, we expand it to "great", "good", "okay", "bad", and "terrible".

 $<sup>^8 \</sup>rm We$  have also tried to tune the baselines with the learning rate of  $10^{-4}$  but it performed worst.

<sup>&</sup>lt;sup>9</sup>We also tried tuning prompts with learning rates of  $\{10^{-4}, 10^{-5}\}$  but it performed worst, as also observed in prior work (Mahabadi et al., 2021a; Min et al., 2021).

887	• For <b>SUBJ</b> , given a sentence <i>s</i> :	" <i>h</i> " ?   <mask>, "<i>p</i>"</mask>
888	<i>s</i> This is <mask>.</mask>	h?   < MASK >, p
889	<i>s</i> It's all <mask>.</mask>	
890	s It's <mask>.</mask>	with "Yes" as a verbalizer for entailment. "No"
891	s Is it <mask>?</mask>	for contradiction, "Maybe" for neutral.
892	with "subjective" and "objective" as verbalizers.	p question: h true, false or neither? answer: <mask></mask>
893 894	• For <b>TREC</b> , given a question q, the task is to classify the type of it:	with "true" as a verbalizer for entailment, "false" for contradiction, "neither" for neutral.
395	<i>q</i> <mask>:</mask>	• For <b>QNLI</b> , given a sentence <i>s</i> and question <i>q</i> :
396	<i>q</i> Q: <mask>:</mask>	s. Question: q? Answer: <mask>.</mask>
897	q why <mask>?</mask>	with "Yes" or "true" as verbalizers for entailment and "No" or "false" for not entailment.
98	q Answer: <mask>.</mask>	s. Based on the previous sentence, q? <mask>.</mask>
99 00 01	with "Description", "Entity", "Expression", "Human", "Location", "Number" as verbalizers for question types of "Description", "Entity",	with "Yes" or "true" as verbalizers for entailment and "No" or "false" for not entailment.
02 03	"Abbreviation", "Human", "Location", and "Numeric".	Based on the following sentence, q? <mask>.s</mask>
04 05	• For entailment task ( <b>RTE</b> ) given a premise <i>p</i> and hypothesis <i>h</i> :	with "Yes" and "No" as verbalizers for entailment and not entailment respectively.
06	" <i>h</i> " ?   <mask>, "<i>p</i>"</mask>	• For <b>QQP</b> , given two questions $q_1$ and $q_2$ :
	<i>h</i> ?   <mask>, <i>p</i></mask>	Do $q_1$ and $q_2$ have the same meaning? <mask>.</mask>
07	"h" ?   <mask>. p</mask>	with "Yes" or "true" as verbalizers for duplicate and "No" or "false" for not duplicate.
08 09 10	with "Yes" as a verbalizer for entailment, "No" for contradiction.	$q_1$ . Based on the previous question, $q_2$ ? <mask>.</mask>
11	<i>p</i> question: <i>h</i> True or False? answer: <mask></mask>	with "Yes" or "true" as verbalizers for duplicate and "No" or "false" for not duplicate.
12 13	with "true" as a verbalizer for entailment, "false" for contradiction.	Based on the following question, $q_1$ ? <mask>.<math>q_2</math></mask>
)14 )15	• For entailment task ( <b>CB</b> ) given a premise <i>p</i> and a hypothesis <i>h</i> :	with "Yes" and "No" as verbalizers for duplicate and not duplicate respectively.

	Do $s_1$ and $s_2$ have the same meaning? <mask>.</mask>
945	
946	with "Yes" or "true" as verbalizers for equivalent
947	and "No" or "false" for not equivalent.
948	s <sub>1</sub> . Based on the previous sentence, $s_2$ ? MASK>.
0.10	
949 950	and "No" or "false" for not equivalent.
951	Based on the following sentence, $s_1?.s_2$
050	with "Vac" and "No" as verbalizers for acuivalant
952 953	and not equivalent respectively.
954	• For <b>WiC</b> , given two sentences $s_1$ and $s_2$ and a
955	word $w$ , the task is to classify whether $w$ is used
956	in the same sense.
957	" $s_1$ " / " $s_2$ ". Similar sense of " $w$ "? <mask>.</mask>
958	$s_1 \ s_2$ Does w have the same meaning in both sentences? <mask></mask>
000	With UNI-U and UX-U as such them. for Esta-
959 960	and True.
961	w . Sense (1) (a) "s <sub>1</sub> " ( <mask>) "s<sub>2</sub>"</mask>
000	
962 963	True.
964	C Impact of the Position
965	of Masks in Sentence-pair Datasets
966	We evaluate the impact of the position of mask tokens
967	in sentence-pair benchmarks. Given two sentences $s_1$
968	and $s_2$ , we consider the following four locations for
969	inserting mask tokens, where in the case of encoding
970	as two sentences, input parts to the encoder are
971	separated with  :
972	1. $s_1 s_2 < MASK >$
973	2. $s_1 \leq MASK > s_2$

 $s_1 \mid \langle MASK \rangle s_2$ 

 $s_1 \mid s_2 < MASK >$ 

3.

4.

974

975

• For **MRPC**, given two sentences  $s_1$  and  $s_2$ :

944

Datasets	1	2	3	4
CB	89.8	91.6	88.9	86.5
RTE	69.1	69.1	64.5	65.3
QNLI	72.0	83.3	77.7	73.1
MRPC	71.6	69.5	66.4	72.0
QQP	79.2	82.8	72.5	70.2
WiC	60.3	59.5	60.2	59.5
Avg	73.7	76.0	71.7	71.1

Table 7: Validation performance for sentence-pair benchmarks for different locations of mask tokens. Bold fonts indicate the best results in each row.

Datasets	$10^{-2}$	$10^{-3}$	$10^{-4}$	$10^{-5}$
СВ	90.0/82.5	92.2/85.0	91.6/87.5	91.6/87.5
MRPC	69.8/56.2	70.8/56.2	69.5/56.2	70.8/56.2
QNLI	83.3/71.9	82.7/71.9	83.3/71.9	83.1/68.8
QQP	82.8/78.1	82.7/75.0	82.8/75.0	83.0/75.0
RTE	69.8/62.5	69.2/59.4	69.1/62.5	68.3/62.5
WiC	62.2/50.0	59.7/46.9	<b>59.5</b> /53.1	58.9/50.0
Avg	76.3/66.9	76.2/65.7	76.0/67.7	76.0/66.7
Total Avg	71.6	71.0	71.8	71.3

Table 8: Validation performance for different values of  $\sigma$ . We show mean performance/worst-case performance across 20 runs. The last row shows the average of mean performance/worst-case performance.

Table 7 shows how the position of masks impact the results. As demonstrated, pattern 2, inserting mask tokens between the two sentences and encoding both as a single sentence obtains the highest validation performance. We use this choice in all the experiments when removing handcrafted patterns.

#### D **Impact of Initialization**

We initialize the label embedding matrix with random initialization from a normal distribution  $\mathcal{N}(0,\sigma)$ . In table 8, we show the development results for different values of  $\sigma$ . We choose the  $\sigma$  obtaining the highest performance on average over average and worst case performance, i.e.,  $\sigma = 10^{-4}$ .

977 978 979 980 981

982

983

984

985

986

987

988