It's Not Just Size That Matters: Small Language Models Are Also Few-Shot Learners

Timo Schick^{1,2} and Hinrich Schütze¹

¹ Center for Information and Language Processing, LMU Munich, Germany
² Sulzer GmbH, Munich, Germany

timo.schick@sulzer.de

Abstract

When scaled to hundreds of billions of parameters, pretrained language models such as GPT-3 (Brown et al., 2020) achieve remarkable few-shot performance. However, enormous amounts of compute are required for training and applying such big models, resulting in a large carbon footprint and making it difficult for researchers and practitioners to use them. We show that performance similar to GPT-3 can be obtained with language models that are much "greener" in that their parameter count is several orders of magnitude smaller. This is achieved by converting textual inputs into cloze questions that contain a task description, combined with gradient-based optimization; exploiting unlabeled data gives further improvements. We identify key factors required for successful natural language understanding with small language models.¹

1 Introduction

Pretraining ever-larger language models (LMs) on massive corpora has led to large improvements in NLP (Radford et al., 2018; Devlin et al., 2019; Liu et al., 2019; Raffel et al., 2020, *i.a.*). A standard approach is to replace the pretrained model's output layer with a task-specific head and finetune the entire model on a set of labeled training data. However, language modeling is not only a powerful pretraining objective, but many tasks can be reformulated as cloze questions (e.g., by appending phrases such as "the correct answer is __"), allowing pretrained LMs to solve them without any or with only very few labeled examples (Radford et al., 2019; Schick and Schütze, 2021).

Recently, Brown et al. (2020) introduced GPT-3, a pretrained LM with an enormous 175 billion parameters, and showed that it has amazing few-shot abilities: By reformulating tasks as LM problems,

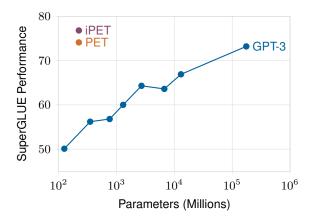


Figure 1: Performance on SuperGLUE with 32 training examples. **ALBERT with PET/iPET outperforms GPT-3 although it is much "greener" in that it has three orders of magnitude fewer parameters.**

GPT-3 achieves near state-of-the-art results for some SuperGLUE (Wang et al., 2019) tasks given just 32 labeled examples. This is achieved through *priming*: GPT-3 is given a few demonstrations of inputs and corresponding outputs as context for its predictions, but no gradient updates are performed. While being straightforward to use, this method has two major drawbacks:

- It requires a gigantic LM to work well, making it unusable in many real-world scenarios and resulting in a large carbon footprint (Strubell et al., 2019).
- It does not scale to more than a few examples as the context window of most LMs is limited to a few hundred tokens.²

An alternative to priming is *pattern-exploiting training* (PET) (Schick and Schütze, 2021), which combines the idea of reformulating tasks as cloze questions with regular gradient-based finetuning. While PET additionally requires unlabeled data, unlabeled data is much easier to obtain than labeled

¹Our implementation is publicly available at https://github.com/timoschick/pet.

²While GPT-3 can process up to 2,048 tokens, this is still not enough to fit \geq 32 examples for some SuperGLUE tasks.

examples for many real-world applications. Crucially, PET only works when the answers to be predicted by the LM correspond to a single token in its vocabulary; this is a severe limitation as many tasks cannot easily be worded that way.

In this work, we adapt PET for tasks that require predicting multiple tokens. We then show that in combination with ALBERT (Lan et al., 2020), PET and its iterative variant (iPET) both outperform GPT-3 on SuperGLUE with 32 training examples, while requiring only 0.1% of its parameters (Figure 1). Moreover, training with PET can be performed in several hours on a single GPU without requiring expensive hyperparameter optimization. Finally, we show that similar performance can also be achieved without unlabeled data and provide a detailed analysis of the factors contributing to PET's strong performance: its ability to combine multiple task formulations, its resilience to wordings that are hard to understand, its usage of labeled data, and characteristics of the underlying LM. Given PET's "green" properties, we see our work as an important contribution to an environmentally sound NLP.

2 Related Work

Enabling LMs to perform zero-shot learning by providing task descriptions was proposed by Radford et al. (2019) and has been applied to text classification (Puri and Catanzaro, 2019), commonsense knowledge mining (Davison et al., 2019) and argumentative relation classification (Opitz, 2019). It is also commonly used for probing the knowledge contained within LMs (Trinh and Le, 2018; Petroni et al., 2019; Talmor et al., 2020; Schick and Schütze, 2020; Ettinger, 2020, *i.a.*).

As finding ways to reformulate tasks as cloze questions that are understood well by LMs is difficult (Jiang et al., 2020), Schick and Schütze (2021) propose PET, a method that uses knowledge distillation (Hinton et al., 2015) and self-training (e.g., Scudder, 1965; Yarowsky, 1995; Brin, 1999; McClosky et al., 2006) to easily combine several reformulations. Our modified version of PET uses masked language models (Devlin et al., 2019) to assign probabilities to sequences of text; this is similar to using them in a generative fashion (Wang and Cho, 2019) and has previously been investigated by Salazar et al. (2020) and Ghazvininejad et al. (2019). In contrast to PET, which uses gradient-based optimization, Radford et al. (2019)

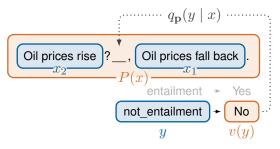


Figure 2: Application of a PVP $\mathbf{p} = (P, v)$ for recognizing textual entailment: An input $x = (x_1, x_2)$ is converted into a cloze question P(x); $q_{\mathbf{p}}(y \mid x)$ for each y is derived from the probability of v(y) being a plausible choice for the masked position.

and Brown et al. (2020) investigate priming, where examples are given as context but no parameter updates are performed.

Finally, our focus on reducing the amount of compute required for few-shot learning is closely related to other efforts in Green AI (Schwartz et al., 2020a) that aim to improve model efficiency, including techniques for knowledge distillation (e.g., Hinton et al., 2015; Sanh et al., 2019; Jiao et al., 2020; Mao et al., 2020; Anderson and Gómez-Rodríguez, 2020), pruning (Han et al., 2015, 2016; Sanh et al., 2020) and quantization (Gong et al., 2014; Zafrir et al., 2019; Stock et al., 2021) as well as early exit strategies for inference (Liu et al., 2020; Schwartz et al., 2020b; Xin et al., 2020).

3 Pattern-Exploiting Training

Let M be a masked language model (MLM), T its vocabulary and $_$ $\in T$ the mask token; we denote the set of all token sequences as T^* . For some $\mathbf{z} \in T^*$ containing at least k masks and $t \in T$, we denote with $q_M^k(t \mid \mathbf{z})$ the probability that M assigns to t at the kth masked position in \mathbf{z} ; the model's logits before applying softmax are denoted with $s_M^k(t \mid \mathbf{z})$. We consider the task of mapping inputs $x \in X$ to outputs $y \in Y$, for which PET requires a set of pattern-verbalizer pairs (PVPs). Each PVP $\mathbf{p} = (P, v)$ consists of

- a pattern $P: X \to T^*$ that maps inputs to cloze questions containing a single mask;
- a verbalizer v : Y → T that maps each output to a single token representing its task-specific meaning in the pattern.

As illustrated in Figure 2, the core idea of PET is to derive the probability of y being the correct output for x from the probability of v(y) being

the "correct" token at the masked position in P(x). Based on this intuition, a conditional probability distribution $q_{\mathbf{p}}$ of y given x is defined as

$$q_{\mathbf{p}}(y \mid x) = \frac{\exp s_{\mathbf{p}}(y \mid x)}{\sum_{y' \in Y} \exp s_{\mathbf{p}}(y' \mid x)} \qquad (1)$$

where $s_{\mathbf{p}}(y \mid x) = s_{M}^{1}(v(y) \mid P(x))$ is the raw score of v(y) at the masked position in P(x).

For a given task, identifying PVPs that perform well is challenging in the absence of a large development set. Therefore, PET enables a combination of multiple PVPs $\mathbf{P} = \{\mathbf{p}_1, \dots, \mathbf{p}_n\}$ as follows:

- 1. For each PVP \mathbf{p} , a MLM is finetuned on training examples (x,y) by minimizing the cross entropy between y and $q_{\mathbf{p}}(y \mid x)$. In practice, Schick and Schütze (2021) train three MLMs per pattern as performance can vary substantially between runs.
- 2. The ensemble of finetuned MLMs is used to annotate a set of unlabeled examples; each unlabeled example $x \in X$ is annotated with soft labels based on the probability distribution

$$q_{\mathbf{P}}(y \mid x) \propto \exp \sum_{\mathbf{p} \in \mathbf{P}} w_{\mathbf{p}} \cdot s_{\mathbf{p}}(y \mid x)$$
 (2)

similar to Eq. 1 where $w_{\mathbf{p}}$ is a weighting term that is proportional to the accuracy achieved with \mathbf{p} on the training set *before* training.

3. The resulting soft-labeled dataset is used to train a regular sequence classifier by minimizing cross entropy between its output and $q_{\mathbf{P}}$.

As steps (2) and (3) above closely resemble knowledge distillation (Hinton et al., 2015), we also refer to them simply as *distillation*. Importantly, this process does not require holding the entire ensemble of MLMs in memory at the same time as each model's predictions can be computed sequentially; therefore, it is not more memory expensive than using a single model.

To give MLMs trained on different patterns further opportunity to learn from one another, Schick and Schütze (2021) also propose iPET, an iterative variant of PET in which several generations of models are trained on datasets of increasing size that are labeled by previous generations. This is achieved as follows: First, an ensemble of MLMs is trained as in regular PET. For each model M_i , a random subset of other models is used to generate

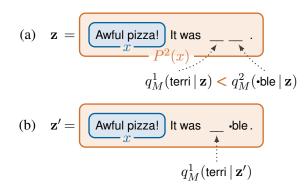


Figure 3: Inference for a verbalization consisting of the two tokens terri and •ble. (a) We first compute the probability of each token at its position in the cloze question $P^2(x)$ and identify the token with the highest probability. (b) We insert this token into the cloze question and compute the probability of the remaining token.

a new training set T_i by assigning labels to those unlabeled examples for which the selected subset of models is most confident in its prediction. Each M_i is then retrained on T_i ; this process is repeated several times, each time increasing the number of examples in T_i by a constant factor. For further details, we refer to Schick and Schütze (2021).

3.1 PET with Multiple Masks

An important limitation of PET is that the verbalizer v must map each output to a single token, which is impossible for many tasks. We thus generalize verbalizers to functions $v:Y\to T^*$; this requires some modifications to inference and training. We further generalize PET in that we do not assume the output space to be identical for each input: for each $x\in X$, we denote with $Y_x\subseteq Y$ the set of possible outputs given x as input. Given a PVP $\mathbf{p}=(P,v)$, we define $l(x)=\max_{y\in Y_x}|v(y)|$ to be the maximum number of tokens required to express any output in Y_x and $P^k(x)$ to be P(x) with the mask token replaced by k masks.

As a running example, we consider the task of binary sentiment classification for restaurant reviews with labels $Y = \{+1, -1\}$. We use the pattern P(x) = x. It was ___ and a verbalizer v that maps +1 to the single token great and -1 to the sequence terri ble, i.e., we assume that the MLM's tokenizer splits the word "terrible" into the two tokens terriand ble. For this example, l(x) = 2 for all x; $P^2(x)$ is illustrated in Figure 3 (a).

³While PET can easily be adapted to generative MLMs (e.g., Lewis et al., 2020; Raffel et al., 2020), we stick with regular MLMs as they are more lightweight and performed better on simple cloze tasks in preliminary experiments.

Inference For $x \in X$, $y \in Y_x$ and |v(y)| = k, we redefine $q_{\mathbf{p}}(y \mid x)$ in an autoregressive fashion: Starting from $P^k(x)$, we perform k consecutive predictions, where we always select the next token to predict based on the MLM's confidence. That is, we set $q_{\mathbf{p}}(y \mid x) = q(v(y) \mid P^k(x))$ where

$$q(t_1 \dots t_k | \mathbf{z}) = \begin{cases} 1 & \text{if } k = 0 \\ q_M^j(t_j | \mathbf{z}) \cdot q(t' | \mathbf{z}') & \text{if } k \ge 1 \end{cases}$$
(3)

with $j = \arg\max_{i=1}^k q_M^i(t_i \mid \mathbf{z})$, \mathbf{z}' is \mathbf{z} except $\mathbf{z}'_j = t_j$ and $t' = t_1 \dots t_{j-1} t_{j+1} \dots t_k$. Note that unlike in original PET (Eq. 1), $q_{\mathbf{p}}$ is not a probability distribution as its values do not sum to one.

For our sentiment classification example, Figure 3 illustrates how $q_{\mathbf{p}}(-1 \mid x)$ is computed: As $|v(y)| = |\{\text{terri}, \cdot \text{ble}\}| = 2$, we first use $\mathbf{z} = P^2(x)$ to compute the probability of each token in v(y) (Figure 3a). We then choose the token with the highest probability, put it in place of the corresponding mask token, and use the resulting cloze question \mathbf{z}' to compute the probability of the remaining token (Figure 3b). The overall score for y = -1 is then computed as

$$q_{\mathbf{p}}(-1 \mid x) = q_M^2(\bullet \mathsf{ble} \mid \mathbf{z}) \cdot q_M^1(\mathsf{terri} \mid \mathbf{z}')$$

Training Computing $q_{\mathbf{p}}(y \mid x)$ as in Eq. 3 for each training example (x,y) would be prohibitively expensive. To enable computation of all required probabilities in a single forward pass, we approximate $q_{\mathbf{p}}(y \mid x)$ by (i) always inserting the maximum number of mask tokens required to express any output and (ii) for each $y' \in Y_x$, predicting all tokens in $v(y') = t_1 \dots t_k$ in parallel, where we simply ignore the model's predictions for all l(x) - k superfluous mask tokens:

$$\tilde{q}_{\mathbf{p}}(y' \mid x) = \prod_{i=1}^{k} q_{M}^{i}(t_{i} \mid P^{l(x)}(x))$$
 (4)

For our running example, this means we approximate the scores $q_{\mathbf{p}}(y \mid x)$ by computing

$$\begin{split} \tilde{q}_{\mathbf{p}}(+1\mid x) &= q_M^1(\text{great}\mid \mathbf{z}) \\ \tilde{q}_{\mathbf{p}}(-1\mid x) &= q_M^1(\text{terri}\mid \mathbf{z}) \cdot q_M^2(\text{•ble}\mid \mathbf{z}) \end{split}$$

which can be done in a single forward pass as it only requires processing the cloze question $\mathbf{z} = P^2(x)$ shown in Figure 3 (a) once.

As $\tilde{q}_{\mathbf{p}}$ is not a probability distribution over Y_x , cross entropy is not an ideal training objective as it

can also be minimized by reducing the probability assigned to sequences $\mathbf{z} \notin v(Y_x)$ that are not part of the output space, despite this having no effect on the model's prediction. We instead opt for multiclass hinge loss (Weston and Watkins, 1999; Dogan et al., 2016) and minimize:

$$\sum_{y' \in Y_x} \max \left(0; 1 - \log \tilde{q}_{\mathbf{p}}(y|x) + \log \tilde{q}_{\mathbf{p}}(y'|x) \right) \tag{5}$$

That is, we require the difference between the log probability of y and the log probability of any output $y' \in Y_x \setminus \{y\}$ to be at least 1.

4 Experiments

We compare PET and GPT-3 on SuperGLUE (Wang et al., 2019), a natural language understanding benchmark consisting of eight challenging tasks. We cannot evaluate PET using the exact same training data as GPT-3 because for most tasks, GPT-3 uses a different set of training examples for each test example and for the other tasks, training sets were not available upon request; however, the exact choice of examples has little impact on GPT-3's performance.⁴ We thus create new training sets by randomly selecting 32 examples for each task using a fixed random seed.

We additionally create sets of up to 20,000 unlabeled examples for each task; this is done by removing all labels from the original training sets. We refer to the resulting sets of training examples and unlabeled examples as *FewGLUE*.⁵

4.1 Tasks

Below, we describe each of the SuperGLUE tasks and our corresponding PVPs. We use a vertical bar (|) to mark boundaries between text segments. Of the eight tasks considered, only COPA, WSC and ReCoRD require the use of PET with multiple masks as introduced in Section 3.1.

BoolQ (Clark et al., 2019) is a QA task where each example consists of a passage p and a yes/no question q. We use the following patterns:

- p. Question: q? Answer: ___.
- p. Based on the previous passage, q? ___.
- Based on the following passage, q? ___. p

⁴Based on personal correspondence with the authors.

⁵FewGLUE is publicly available at https://github.com/timoschick/fewglue.

We define two verbalizers mapping questions containing a true statement to yes/true and others to no/false, respectively, for a total of 6 PVPs.

CB (De Marneffe et al., 2019) and **RTE** (Dagan et al., 2006) are textual entailment tasks like MNLI, so we use PVPs similar to Schick and Schütze (2021). For a premise p and hypothesis h, we use

and a verbalizer that maps entailment to yes, disagreement to no and neutral to maybe.

Given a premise p, the task in **COPA** (Gordon et al., 2012) is to determine the *cause* or *effect* of the premise given two options c_1 and c_2 . For determining the *effect*, we use the following patterns:

"
$$c_1$$
" or " c_2 "? p , so $_$. , c_1 or c_2 ? p , so $_$.

For determining the *cause*, we use the same patterns but replace so with because. The verbalizer for c_1 and c_2 is the identity function.

For **WiC** (Pilehvar and Camacho-Collados, 2019), given a word w and two sentences s_1 and s_2 in which it occurs, the task is to decide if w is used with the same sense in both sentences. We use:

- " s_1 " / " s_2 ". Similar sense of "w"? ___.
- $s_1 \ s_2$ Does w have the same meaning in both sentences?
- w. Sense (1) (a) " s_1 " (__) " s_2 "

For the first two patterns, we use yes as verbalization for words used in the same sense and no for other words; for the third pattern, we use b and 2.

For WSC (Levesque et al., 2011), each example consists of a sentence s with a marked pronoun p and noun n, and the task is to determine whether p refers to n. We follow (Raffel et al., 2020; Brown et al., 2020) and treat WSC as a generative task. We highlight p in s by putting it in asterisks and use the following patterns:

- s The pronoun '*p*' refers to ___.
- s In the previous sentence, the pronoun '*p*' refers to $\ .$
- s In the passage above, what does the pronoun '*p*' refer to? Answer: .

We use the identity function as verbalizer for n. Note that WSC is different from other tasks in that it requires free-form completion. This in

turn requires some modifications during training and inference that are discussed in Appendix A.

MultiRC (Khashabi et al., 2018) is a QA task. Given a passage p, a question q and an answer candidate a, the task is to decide whether a is a correct answer for q. We use the same verbalizer as for BoolQ and similar patterns:

- p. Question: q? Is it a? ___.
- p. Question: q? Is the correct answer "a"? ___.
- p. Based on the previous passage, q? Is "a" a correct answer? ___.

For **ReCoRD** (Zhang et al., 2018), given a passage p and a cloze question q, the task is to decide which of a given set of answer candidates is the correct replacement for the placeholder in the cloze question. As this task is already presented in the form of a cloze question, there is little room for designing PVPs, so we only use a trivial one: the concatenation of p and q as pattern and the identity function as verbalizer. With only one PVP, there is no need to perform knowledge distillation so we directly use the resulting model as our final classifier.

4.2 Setup

As underlying LM for PET we choose ALBERTxxlarge-v2 (Lan et al., 2020), the best-performing MLM on SuperGLUE when training is performed on the regular, full size training sets. We use the same model, supplemented by a sequence classification head, as our final classifier. We run PET on the FewGLUE training sets for all SuperGLUE tasks. We do not use any development set to optimize hyperparameters; instead we use the exact same setup and hyperparameters as Schick and Schütze (2021). For COPA, WSC and ReCoRD, we use our proposed modification of PET to support verbalizers mapping labels to multiple tokens; for all other tasks, we use regular PET. We train iPET on all tasks except COPA and WSC, as their unlabeled sets contain well below 1,000 examples, as well as ReCoRD, for which iPET makes no sense as we only use a single PVP. For these three tasks, we simply reuse the results of regular PET.

4.3 Results

Our main results are shown in Table 1. As can be seen, ALBERT with PET performs similar to the largest GPT-3 model, which is larger by a factor

	Model	Params (M)	BoolQ Acc.	CB Acc. / F1	COPA Acc.	RTE Acc.	WiC Acc.	WSC Acc.	MultiRC EM / F1a	ReCoRD Acc. / F1	Avg
	GPT-3 Small	125	43.1	42.9 / 26.1	67.0	52.3	49.8	58.7	6.1 / 45.0	69.8 / 70.7	50.1
	GPT-3 Med	350	60.6	58.9 / 40.4	64.0	48.4	55.0	60.6	11.8 / 55.9	77.2 / 77.9	56.2
	GPT-3 Large	760	62.0	53.6 / 32.6	72.0	46.9	53.0	54.8	16.8 / 64.2	81.3 / 82.1	56.8
	GPT-3 XL	1,300	64.1	69.6 / 48.3	77.0	50.9	53.0	49.0	20.8 / 65.4	83.1 / 84.0	60.0
>	GPT-3 2.7B	2,700	70.3	67.9 / 45.7	83.0	56.3	51.6	62.5	24.7 / 69.5	86.6 / 87.5	64.3
dev	GPT-3 6.7B	6,700	70.0	60.7 / 44.6	83.0	49.5	53.1	67.3	23.8 / 66.4	87.9 / 88.8	63.6
	GPT-3 13B	13,000	70.2	66.1 / 46.0	86.0	60.6	51.1	75.0	25.0 / 69.3	88.9 / 89.8	66.9
	GPT-3	175,000	77.5	82.1 / 57.2	92.0	72.9	55.3	75.0	32.5 / 74.8	89.0 / 90.1	73.2
	PET	223	79.4	85.1 / 59.4	95.0	69.8	52.4	80.1	37.9 / 77.3	86.0 / 86.5	74.1
	iРЕТ	223	80.6	92.9 / 92.4	95.0	74.0	52.2	80.1	33.0 / 74.0	86.0 / 86.5	76.8
	GPT-3	175,000	76.4	75.6 / 52.0	92.0	69.0	49.4	80.1	30.5 / 75.4	90.2 / 91.1	71.8
;t	PET	223	79.1	87.2 / 60.2	90.8	67.2	50.7	88.4	36.4 / 76.6	85.4 / 85.9	74.0
test	iPET	223	81.2	88.8 / 79.9	90.8	70.8	49.3	88.4	31.7 / 74.1	85.4 / 85.9	75.4
	SotA	11,000	91.2	93.9 / 96.8	94.8	92.5	76.9	93.8	88.1 / 63.3	94.1 / 93.4	89.3

Table 1: Results on SuperGLUE for GPT-3 primed with 32 randomly selected examples and for PET / iPET with ALBERT-xxlarge-v2 after training on FewGLUE. State-of-the-art results when using the regular, full size training sets for all tasks (Raffel et al., 2020) are shown in italics.

of 785. On average, PET performs 18 points better compared to GPT-3 Med, a model of similar size. iPET brings further improvements for 3 out of the 5 tasks that we use iPET for, most notably for CB, but results in a slight performance drop for MultiRC. Despite PET's strong performance, it still clearly performs worse than a state-of-the-art model trained on the regular, full size SuperGLUE training set.

5 Analysis

We investigate the importance of several factors for few-shot performance: the choice of patterns and verbalizers, the usage of both unlabeled and labeled data, and properties of the underlying language model. We also look into our proposed modification for PET to work with multiple masks and compare it to various baselines. Finally, we measure how choosing different sets of training examples affects performance. Our analysis focuses on PET as GPT-3 is not publicly available.⁶

5.1 Patterns

The way in which tasks are reformulated as cloze questions can have a huge impact on performance (Jiang et al., 2020; Schick and Schütze, 2021). These reformulations can be arbitrarily complex; for example, the pattern used by GPT-3 for WSC contains an introductory section of almost 30 words; it is unclear if and how this formulation has been optimized.⁷ To investigate the importance

of patterns and verbalizers, we compare three sets of PVPs: our initial set as defined in Section 4.1 (denoted \mathbf{p}_{ours}), the single PVP used by GPT-3 (\mathbf{p}_{GPT-3}), and the combination of both (\mathbf{p}_{comb}).

We train ALBERT using PET with all three sets of patterns; results for selected SuperGLUE tasks are shown in Table 2 (top). As can be seen, the PVP used by GPT-3 outperforms our PVPs on RTE whereas our initial set of patterns performs much better on MultiRC. These large differences in performance highlight the importance of finding good ways to express tasks as cloze questions. As it is difficult to ascertain which patterns perform well without trying them on a large set of examples, a key challenge for few-shot approaches is to compensate for PVPs that the LM fails to understand well. As seen in the performance of the model trained with p_{comb}, PET is able to do so: not only does combining all PVPs compensate for the worse performance of \mathbf{p}_{ours} on RTE and of p_{GPT-3} on MultiRC, it even further improves average performance across the three tasks compared to the best-performing set of patterns. This clearly demonstrates the potential of carefully engineering a set of suitable patterns as opposed to just choosing a single formulation without means of evaluating its effectiveness.

5.2 Unlabeled Data Usage

Unlike GPT-3, PET requires unlabeled data to distill the knowledge of all models based on individual PVPs into a single classifier; for iPET, unlabeled data is additionally used to generate training sets for future generations. The underlying assumption

⁶We could not obtain access to OpenAI's GPT-3 API.

⁷While the authors use a different terminology, GPT-3 also makes use of PVPs (Brown et al., 2020, pp. 50–61).

Model	CB	RTE	MultiRC	Avg
	Acc. / F1	Acc.	EM / F1a	-
PET (p _{ours}) PET (p _{GPT-3}) PET (p _{comb})	85.1 / 59.4	69.8	37.9 / 77.3	66.6
	83.3 / 58.1	71.8	25.4 / 68.3	63.1
	84.5 / 59.0	74.7	39.1 / 77.7	68.3
PET (p _{ours}) ¬dist	83.9 / 76.2	66.4	38.9 / 76.2	68.0
PET (p _{comb}) ¬dist	83.9 / 76.2	72.9	39.6 / 76.6	70.4

Table 2: Results on selected tasks for various sets of PVPs for regular PET and for an ensemble of PET models with no knowledge distillation ("¬dist")

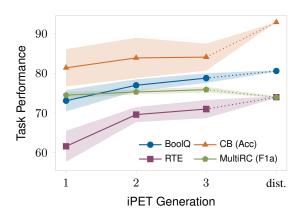


Figure 4: Average performance (± standard deviation) of all MLMs trained on individual patterns for three generations and of the distilled classifier ("dist.") across three individual training runs

is that unlabeled data can easily be obtained, which may not always be the case in real-world settings. We thus investigate the importance of unlabeled data for regular PET. To this end, we compare the performance of the final classifier in PET to that of directly using the ensemble of models corresponding to individual PVPs. While using this ensemble entirely removes the need for unlabeled data, the ensemble for k PVPs is larger than the distilled model by a factor of $3 \cdot k$ as we follow the default setting of PET and train three models per PVP. However, even for a large number of PVPs the ensemble is smaller than GPT-3 by two orders of magnitude.

Results without distillation can be seen in Table 2 (bottom). Averaged across the three tasks, the ensemble performs even better than the distilled classifier. This shows that if the goal is only to achieve good performance, then unlabeled data is not necessary; however, it is required to obtain a single, lightweight model as final classifier.

Figure 4 illustrates the benefit of training multiple generations with iPET. For all tasks except MultiRC, there are substantial improvements from

Model	CB Acc. / F1	RTE Acc.	MultiRC EM / F1a	Avg -
PET unsupervised supervised	85.1 / 59.4 33.5 / 23.1 60.7 / 42.5	69.8 55.0 50.2	37.9 / 77.3 3.9 / 60.3 4.3 / 49.8	66.6 38.5 43.0
PET (XLNet) Priming (XLNet)	88.7 / 83.0 56.3 / 37.7	60.4 49.5	21.4 / 66.6 - / -	63.4

Table 3: Results on selected tasks for various ways of using the labeled examples available in FewGLUE

the first to the second generation, whereas the third generation achieves only slight additional improvements. On average, standard deviation is reduced in later generations, illustrating that the models learn from each other and their predictions converge. The final distillation step brings further improvements for all tasks except MultiRC and reduces standard deviation across three training runs to almost zero, illustrating that PET and iPET are effective means of reducing finetuning instability (Dodge et al., 2020).

Of course, there are further ways to leverage unlabeled data such as keeping an auxiliary language modeling objective during finetuning (Chronopoulou et al., 2019). While we leave investigating the impact of additionally using such methods to future work, we note that they can easily be applied to PET while there is no straightforward way to combine them with priming.

5.3 Labeled Data Usage

We next investigate the effect of how labeled data is used, which is one of the key differences between priming and PET. We first compare PET with regular supervised training (i.e., without using any patterns), and with a fully unsupervised model (i.e., an ensemble using all PVPs but no labeled training examples). Given 32 examples, PET clearly outperforms both baselines (Table 3).

We next compare PET directly to priming. However, we cannot do so using ALBERT as it is only able to process sequences of up to 512 tokens, which is not enough for a set of 32 examples; we instead use XLNet (Yang et al., 2019) for this comparison. As shown in Table 3, XLNet in general performs worse than ALBERT. More importantly, XLNet with PET performs much better than priming. We were not able to obtain results with priming on MultiRC because the 32 examples in FewGLUE would require more than 10,000 tokens, so processing them with a standard Transformer (Vaswani

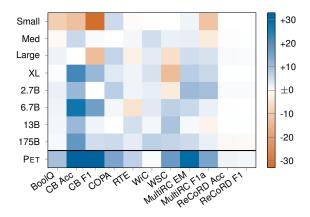


Figure 5: Accuracy differences between priming with 32 examples and one-shot priming for all GPT-3 models as well as between ALBERT with PET (without distillation) and unsupervised ALBERT (bottom row)

et al., 2017) is infeasible due to the quadratic complexity of self-attention. This highlights another important issue with priming: It does not scale well to more than a few examples; even GPT-3 is only able to process sequences of up to 2,048 tokens. While there are some Transformer variants that can deal with much longer contexts (e.g., Kitaev et al., 2020; Beltagy et al., 2020), it has yet to be investigated to what extent such models make good use of priming examples over long context spans.

We further investigate the effectiveness of priming by looking at results obtained with GPT-3 more closely. To this end, Figure 5 shows the performance difference between priming GPT-3 with 32 examples and priming it with just a single example for each task and model size. As can be seen, priming with 32 examples only slightly improves performance for most tasks and model sizes. For some tasks, adding more examples even leads to worse performance, especially for smaller models. For ReCoRD, even the largest model's performance slightly drops when adding more examples.

The bottom row of Figure 5 shows the performance difference between ALBERT trained with PET (without distillation) and a fully unsupervised ALBERT model on all tasks. While results are not directly comparable due to different underlying models and PVPs, PET results in much stronger performance improvements compared to priming and does not worsen results for any task.

Model	Params			MultiRC EM / F1a	0
ALBERT	223M	87.5 / 78.7	74.7	38.9 / 76.2	71.8
RoBERTa	355M	85.7 / 77.5	62.8	23.3 / 70.0	63.7
GPT-2	345M	73.2 / 73.7	47.7	12.4 / 57.4	52.0

Table 4: Results on selected tasks for PET without knowledge distillation combined with various LMs using p_{GPT-3} for CB/RTE and p_{ours} for MultiRC

5.4 Model Type

We next look into the impact of the underlying LM on PET by comparing ALBERT with RoBERTa large (Liu et al., 2019) and GPT-2 medium (Radford et al., 2019). As GPT-2 is a unidirectional model similar to GPT-3, it can only process patterns where the mask token is the very last token. We therefore use p_{GPT-3} for CB and RTE; for MultiRC, we stick with our original set of patterns as they already fulfill this requirement. We also do not perform distillation and instead report the ensemble's performance as there is no established way of equipping GPT-2 with a sequence classification head.

Results for training all three LMs with PET in Table 4 show that using ALBERT as underlying LM is crucial for PET's strong performance; exchanging ALBERT with RoBERTa results in an average performance drop of 8 points. However, RoBERTa still clearly outperforms GPT-3 13B, which is larger by two orders of magnitude. Importantly, PET with GPT-2 performs much worse than with the two other models. As anticipated by Brown et al. (2020), a reason for this drop in performance may be that like GPT-3, GPT-2 is unidirectional, making tasks that require comparing two sequences a challenge. However, it is important to note that there are also other substantial differences between GPT-2 and the other two models, most notably the pretraining dataset. Regardless of whether unidirectionality is the reason for GPT-2's bad performance, bidirectionality of the underlying LM is important for PET as it removes the need for the mask token to be at the very end and thus allows for more flexibility in the creation of patterns.

5.5 PET with Multiple Masks

We modified PET to work for outputs that require more than a single token. To investigate the impact of this modification, we look at the three tasks for which this is required: COPA, WSC and ReCoRD. We compare our decoding strategy of predicting to-

⁸We do not compare priming to zero-shot performance as for unknown reasons, zero-shot GPT-3 performs well below random guessing for some tasks (e.g., 0.0% accuracy for WiC). To not overestimate the benefit of priming, we therefore show gains from providing 32 examples compared to just one.

Model	COPA Acc.	WSC Acc.	ReCoRD Acc. / F1	Avg
PET	95.0	80.1	86.0 / 86.5	87.1
PET ¬dist (max-first)	90.0	80.8	86.0 / 86.5	85.7
PET ¬dist (ltr)	89.0	79.8	84.7 / 85.3	84.6
PET ¬dist (parallel)	77.0	80.8	82.5 / 83.1	80.2
untrained	72.5	59.9	84.7 / 85.4	72.5

Table 5: Results on selected tasks for our proposed variant of PET as well as other decoding strategies and for untrained ALBERT

kens in order of the probability assigned to them, to which we refer as *max-first*, with two alternatives: decoding left-to-right (ltr) as is common for many autoregressive language models, and decoding all tokens simultaneously (parallel) as is done during training. Additionally, we compare PET with untrained ALBERT to measure the effectiveness of our proposed training loss.

Results are shown in Table 5. PET clearly outperforms untrained ALBERT for the three tasks. Not performing distillation hurts performance for COPA, but leads to slight improvements on WSC; for ReCoRD, we did not perform distillation in the first place as we only use a single PVP. Our decoding strategy is clearly superior to parallel decoding except for WSC, for which most predictions consist only of one or two tokens, and performs slightly better than left-to-right decoding.

5.6 Training Examples

Recall that we conduct our experiments with training examples from FewGLUE, a randomly selected subset of the original SuperGLUE training examples. We used a fixed random seed s_0 to generate FewGLUE. Let Σ_i be the randomly selected subset of SuperGLUE for random seed s_i , so $\Sigma_0 =$ FewGLUE. In this subsection, we create two additional subsets of SuperGLUE, Σ_1 and Σ_2 , based on different seeds. This allows us to investigate how different sets of training examples affect performance. To this end, we run PET for CB, RTE and MultiRC using the three Σ_i . To measure only the effect of varying the training set while ignoring unlabeled examples, we do not use distillation.

Table 6 shows that for all tasks, changing the set of training examples can result in large performance differences for PET. This highlights the importance of using the same set of examples when comparing different few-shot approaches, which is why we make the particular set of examples in FewGLUE publicly available. However, we note

Model	CB Acc. / F1	RTE Acc.	MultiRC EM / F1a	Avg
GPT-3	82.1 / 57.2	72.9	32.5 / 74.8	65.4
PET \neg dist (Σ_0)	83.9 / 76.2	66.4	38.9 / 76.2	68.0
PET \neg dist (Σ_1)	82.1 / 57.4	61.4	39.2 / 77.9	63.2
PET \neg dist (Σ_2)	87.5 / 84.0	61.4	34.7 / 76.3	67.6

Table 6: Results on selected tasks for GPT-3 and for PET using training sets Σ_0 , Σ_1 , Σ_2

that the average performance of PET is similar to that of GPT-3 for all seeds.

While our results may seem contrary to the insight that for GPT-3, the exact choice of examples does not play a major role, we suspect this to be due to the fact that priming benefits much less from training examples than PET (cf. Section 5.3); accordingly, the influence of the exact set of training examples on the model's performance is smaller.

6 Conclusion

We have proposed a simple yet effective modification of PET, enabling us to use it for tasks that require predicting multiple tokens. In extensive experiments, we have identified several factors responsible for the strong performance of PET combined with ALBERT: the possibility to concurrently use multiple patterns for transforming examples into cloze questions, the ability to compensate for patterns that are difficult to understand, the usage of labeled data to perform parameter updates, and the underlying LM itself.

We have shown that using PET, it is possible to achieve few-shot text classification performance similar to GPT-3 on SuperGLUE with LMs that have three orders of magnitude fewer parameters. This not only lowers financial cost, but above all reduces environmental impact immensely and leads to a much smaller carbon footprint. We see this as an important contribution to achieving the goal of an environmentally more friendly NLP. To enable comparisons with our work, we make our code, models and datasets publicly available.

For future work, it would be interesting to see whether PET also works for generative tasks when combined with generative LMs and whether further improvements are possible in multi-task settings.

Acknowledgments This work was funded by the European Research Council (ERC #740516). We thank the anonymous reviewers for their helpful comments.

References

- Mark Anderson and Carlos Gómez-Rodríguez. 2020. Distilling neural networks for greener and faster dependency parsing. In *Proceedings of the 16th International Conference on Parsing Technologies and the IWPT 2020 Shared Task on Parsing into Enhanced Universal Dependencies*, pages 2–13, Online. Association for Computational Linguistics.
- Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *Computing Research Repository*, arXiv:2004.05150.
- Sergey Brin. 1999. Extracting patterns and relations from the world wide web. In *The World Wide Web and Databases*, pages 172–183, Berlin, Heidelberg. Springer Berlin Heidelberg.
- 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 Advances in Neural Information Processing Systems, volume 33, pages 1877–1901. Curran Associates, Inc.
- Alexandra Chronopoulou, Christos Baziotis, and Alexandros Potamianos. 2019. An embarrassingly simple approach for transfer learning from pretrained language models. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2089–2095, Minneapolis, Minnesota. Association for Computational Linguistics.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. BoolQ: Exploring the surprising difficulty of natural yes/no questions. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2924–2936, Minneapolis, Minnesota. Association for Computational Linguistics.
- Ido Dagan, Oren Glickman, and Bernardo Magnini. 2006. The PASCAL recognising textual entailment challenge. In Machine learning challenges. evaluating predictive uncertainty, visual object classification, and recognising tectual entailment, pages 177–190. Springer.
- Joe Davison, Joshua Feldman, and Alexander Rush. 2019. Commonsense knowledge mining from pretrained models. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language*

- Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 1173–1178, Hong Kong, China. Association for Computational Linguistics.
- Marie-Catherine De Marneffe, Mandy Simons, and Judith Tonhauser. 2019. The CommitmentBank: Investigating projection in naturally occurring discourse. In *Proceedings of Sinn und Bedeutung 23*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- 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. *Computing Research Repository*, arXiv:2002.06305.
- Ürün Dogan, Tobias Glasmachers, and Christian Igel. 2016. A unified view on multi-class support vector classification. *J. Mach. Learn. Res.*, 17(45):1–32.
- Allyson Ettinger. 2020. What BERT is not: Lessons from a new suite of psycholinguistic diagnostics for language models. *Transactions of the Association for Computational Linguistics*, 8:34–48.
- Marjan Ghazvininejad, Omer Levy, Yinhan Liu, and Luke Zettlemoyer. 2019. Mask-predict: Parallel decoding of conditional masked language models. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 6112–6121, Hong Kong, China. Association for Computational Linguistics.
- Yunchao Gong, Liu Liu, Ming Yang, and Lubomir Bourdev. 2014. Compressing deep convolutional networks using vector quantization. *Computing Research Repository*, arXiv:1412.6115.
- Andrew Gordon, Zornitsa Kozareva, and Melissa Roemmele. 2012. SemEval-2012 task 7: Choice of plausible alternatives: An evaluation of commonsense causal reasoning. In *SEM 2012: The First Joint Conference on Lexical and Computational Semantics Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012), pages 394–398, Montréal, Canada. Association for Computational Linguistics.
- Song Han, Huizi Mao, and William J Dally. 2016. Deep compression: Compressing deep neural net-

- works with pruning, trained quantization and huffman coding. *International Conference on Learning Representations (ICLR)*.
- Song Han, Jeff Pool, John Tran, and William Dally. 2015. Learning both weights and connections for efficient neural network. In *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *Computing Research Repository*, arXiv:1503.02531.
- Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham Neubig. 2020. How can we know what language models know? *Transactions of the Association for Computational Linguistics*, 8:423–438.
- Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. 2020. TinyBERT: Distilling BERT for natural language understanding. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4163–4174, Online. Association for Computational Linguistics.
- Daniel Khashabi, Snigdha Chaturvedi, Michael Roth, Shyam Upadhyay, and Dan Roth. 2018. Looking beyond the surface: A challenge set for reading comprehension over multiple sentences. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 252–262, New Orleans, Louisiana. Association for Computational Linguistics.
- Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. 2020. Reformer: The efficient transformer. In *International Conference on Learning Representations*.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. ALBERT: A lite BERT for self-supervised learning of language representations. In *International Conference on Learning Representations*.
- Hector J Levesque, Ernest Davis, and Leora Morgenstern. 2011. The Winograd schema challenge. In *AAAI Spring Symposium: Logical Formalizations of Commonsense Reasoning*, volume 46, page 47.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pretraining for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Weijie Liu, Peng Zhou, Zhiruo Wang, Zhe Zhao, Haotang Deng, and Qi Ju. 2020. FastBERT: a self-distilling BERT with adaptive inference time. In

- Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 6035–6044, Online. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A robustly optimized BERT pretraining approach. *Computing Research Repository*, arXiv:1907.11692.
- Yihuan Mao, Yujing Wang, Chufan Wu, Chen Zhang, Yang Wang, Quanlu Zhang, Yaming Yang, Yunhai Tong, and Jing Bai. 2020. LadaBERT: Lightweight adaptation of BERT through hybrid model compression. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 3225–3234, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- David McClosky, Eugene Charniak, and Mark Johnson. 2006. Effective self-training for parsing. In *Proceedings of the Human Language Technology Conference of the NAACL, Main Conference*, pages 152–159, New York City, USA. Association for Computational Linguistics.
- Juri Opitz. 2019. Argumentative relation classification as plausibility ranking. In Preliminary proceedings of the 15th Conference on Natural Language Processing (KONVENS 2019): Long Papers, pages 193– 202, Erlangen, Germany. German Society for Computational Linguistics & Language Technology.
- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic differentiation in PyTorch. In NIPS Autodiff Workshop.
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP).
- Mohammad Taher Pilehvar and Jose Camacho-Collados. 2019. WiC: the word-in-context dataset for evaluating context-sensitive meaning representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1267–1273, Minneapolis, Minnesota. Association for Computational Linguistics.
- Raul Puri and Bryan Catanzaro. 2019. Zero-shot text classification with generative language models. *Computing Research Repository*, arXiv:1912.10165.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training.

- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. Technical report.
- 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. *Journal of Machine Learning Research*, 21(140):1–67.
- Julian Salazar, Davis Liang, Toan Q. Nguyen, and Katrin Kirchhoff. 2020. Masked language model scoring. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2699–2712, Online. Association for Computational Linguistics.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. In Proceedings of the 5th Workshop on Energy Efficient Machine Learning and Cognitive Computing, NeurIPS 2019.
- Victor Sanh, Thomas Wolf, and Alexander Rush. 2020. Movement pruning: Adaptive sparsity by finetuning. In Advances in Neural Information Processing Systems, volume 33, pages 20378–20389. Curran Associates, Inc.
- Timo Schick and Hinrich Schütze. 2020. Rare words: A major problem for contextualized embeddings and how to fix it by attentive mimicking. In *Proceedings of the Thirty-Fourth AAAI Conference on Artificial Intelligence*.
- Timo Schick and Hinrich Schütze. 2021. Exploiting cloze questions for few shot text classification and natural language inference. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics*, Kyiv, Ukraine (Online). International Committee on Computational Linguistics.
- Roy Schwartz, Jesse Dodge, Noah A. Smith, and Oren Etzioni. 2020a. Green AI. *Commun. ACM*, 63(12):54–63.
- Roy Schwartz, Gabriel Stanovsky, Swabha Swayamdipta, Jesse Dodge, and Noah A. Smith. 2020b. The right tool for the job: Matching model and instance complexities. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6640–6651, Online. Association for Computational Linguistics.
- H Scudder. 1965. Probability of error of some adaptive pattern-recognition machines. *IEEE Transactions on Information Theory*, 11(3):363–371.
- Pierre Stock, Angela Fan, Benjamin Graham, Edouard Grave, Rémi Gribonval, Herve Jegou, and Armand Joulin. 2021. Training with quantization noise for extreme model compression. In *International Conference on Learning Representations*.

- Emma Strubell, Ananya Ganesh, and Andrew McCallum. 2019. Energy and policy considerations for deep learning in NLP. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3645–3650, Florence, Italy. Association for Computational Linguistics.
- Alon Talmor, Yanai Elazar, Yoav Goldberg, and Jonathan Berant. 2020. oLMpics on what language model pre-training captures. *Transactions of the Association for Computational Linguistics*, 8:743–758.
- Trieu H. Trinh and Quoc V. Le. 2018. A simple method for commonsense reasoning. *Computing Research Repository*, arXiv:1806.02847.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems 30*, pages 5998–6008. Curran Associates, Inc.
- Alex Wang and Kyunghyun Cho. 2019. BERT has a mouth, and it must speak: BERT as a Markov random field language model. In *Proceedings of the Workshop on Methods for Optimizing and Evaluating Neural Language Generation*, pages 30–36, Minneapolis, Minnesota. Association for Computational Linguistics.
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. SuperGLUE: A stickier benchmark for general-purpose language understanding systems. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Jason Weston and Chris Watkins. 1999. Support vector machines for multi-class pattern recognition. In *ESANN*, volume 99, pages 219–224.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 38–45, Online. Association for Computational Linguistics.
- Ji Xin, Rodrigo Nogueira, Yaoliang Yu, and Jimmy Lin. 2020. Early exiting BERT for efficient document ranking. In *Proceedings of SustaiNLP: Workshop on Simple and Efficient Natural Language Processing*, pages 83–88, Online. Association for Computational Linguistics.

Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 5753–5763. Curran Associates, Inc.

David Yarowsky. 1995. Unsupervised word sense disambiguation rivaling supervised methods. In 33rd Annual Meeting of the Association for Computational Linguistics, pages 189–196, Cambridge, Massachusetts, USA. Association for Computational Linguistics.

Ofir Zafrir, Guy Boudoukh, Peter Izsak, and Moshe Wasserblat. 2019. Q8BERT: Quantized 8bit BERT. In *NeurIPS EMC2 Workshop*.

Sheng Zhang, Xiaodong Liu, Jingjing Liu, Jianfeng Gao, Kevin Duh, and Benjamin Van Durme. 2018. ReCoRD: Bridging the gap between human and machine commonsense reading comprehension. *Computing Research Repository*, arXiv:1810.12885.

A Training Details

Our implementation can be found in the supplementary material. It extends the original implementation of PET by Schick and Schütze (2021) which, in turn, is based on the Transformers library (Wolf et al., 2020) and PyTorch (Paszke et al., 2017). All dependencies are listed in requirements.txt. Detailed instructions on how our results can be reproduced using this implementation can be found in README.md.

Unless explicitly stated differently, we use the exact same set of hyperparameters as Schick and Schütze (2021) (Table 7) with the only difference that for iPET, we only train 3 generations of models to speed up training. All of our experiments were conducted using a single GPU with 11GB RAM (NVIDIA GeForce GTX 1080 Ti). With this GPU, training a single PET model for 250 steps took approximately 45 minutes. Depending on the task, labeling unlabeled examples took 0.2–1.5 hours per model. Training the final classifier for 5,000 steps on the soft-labeled dataset took 2.5 hours on average. Below, we list task-specific implementation details for all tasks in SuperGLUE.

COPA For COPA, we randomly switch the two options c_1 and c_2 during training with a probability of 50% to make the input more diverse; for inference, we always keep the original order. For distilling the final PET model, we obtain logits for unlabeled examples x from individual PVPs \mathbf{p} as

 $s_{\mathbf{p}}(y \mid x) = \log q_{\mathbf{p}}(y \mid x)$; we use the input format proposed by Liu et al. (2019).

WiC Similar to COPA, we randomly switch the input sentences s_1 and s_2 during training. Given a word w and two sentences s_1 and s_2 , we use the sequence w: $s_1 \mid s_2$ as input for the final sequence classification model, where \mid marks the boundary between two text segments.

WSC Unlike other SuperGLUE tasks, the WSC formulation of Raffel et al. (2020) and Brown et al. (2020) requires free-form completion, meaning that for each sentence s and pronoun p, we only have a single correct choice n that the model needs to predict, but we do not provide any alternatives. During training, we thus use regular cross entropy loss between n and $\tilde{q}_{\mathbf{p}}(n \mid s, p)$ as defined in Eq. 4. However, in many cases this would allow the LM to easily identify the correct target based on the number of masks provided, so we modify each target by randomly adding up to three additional mask tokens, for which we require the model to predict a special <pad> token. For inference, we always just add a single mask token to ensure consistent results across multiple evaluations and perform greedy decoding as described in Section 3. We then follow Raffel et al. (2020) to map the output produced by the LM to a label $y \in \{\text{true}, \text{false}\}.$ For distillation, given an unlabeled example x we set $s_{\mathbf{p}}(y \mid x) = 1$ if the model's output for x was mapped to y and $s_{\mathbf{p}}(y \mid x) = 0$ otherwise. We provide inputs to the final PET model in the format $s \mid n$ where \mid is the boundary between two text segments and mark p in s with asterisks.

MultiRC Deviating from the hyperparameters used by Schick and Schütze (2021), we use a maximum sequence length of 512 tokens for MultiRC both during training and inference because we found many passages to be much longer than 256 tokens. Input for the final sequence classification model is of the form $p \mid q \mid a$ where p is the passage, q is the question, a is the answer candidate and we use \mid to mark boundaries between text segments.

ReCoRD For ReCoRD, we again use a maximum sequence length of 512 because many passages require more than 256 tokens. For some questions q, the ReCoRD training set contains a huge number of answer candidates. To facilitate training, we split each example into multiple examples as follows: let C be the set of answer candidates

Parameter	Value
adam_epsilon	1e-8
gradient_accumulation_steps	8
learning_rate	1e-5
max_grad_norm	1.0
max_seq_length	256
pet_max_steps	250
sc_max_steps	5,000
per_gpu_train_batch_size	2
distillation_temperature	2
weight_decay	0.01

Table 7: Hyperparameters for PET from Schick and Schütze (2021)

Dataset	Metrics	Unlabeled	Dev	Test
BoolQ	Acc.	9,427	3,270	3,245
СВ	Acc./F1	20,000	57	250
COPA	Acc.	400	100	500
MultiRC	F1a/EM	5,100	953	1,800
ReCoRD	F1/EM	20,000	10,000	10,000
RTE	Acc.	20,000	278	300
WiC	Acc.	6,000	638	1,400
WSC	Acc.	554	104	146

Table 8: Important statistics for all datasets used

with $C^+ \subset C$ being the set of correct answers. We create a training example for each $c \in C^+$ by randomly selecting up to 9 negative examples from $C \setminus C^+$ for a total of 10 answer candidates.

B Dataset Details

For each task and number of examples t, we create the FewGLUE training set \mathcal{T} by shuffling the entire original training set with a fixed random seed and collecting the first 32 examples of the shuffled dataset. Following (Raffel et al., 2020; Brown et al., 2020), we select only positive examples for WSC; for both MultiRC and ReCoRD, we follow Brown et al. (2020) and select a total of 32 questions – which corresponds to more than 32 training examples – to enable a fair comparison with GPT-3.

The unlabeled datasets for all tasks are obtained by collecting up to 20,000 examples from their training sets and removing the labels. As the training sets for RTE and CB are very small, for both tasks we additionally select random unlabeled examples from the MNLI training set for a total of 20,000 examples. For evaluation, we use the official validation and test sets for all tasks that are available at https://super.gluebenchmark.com/tasks.All datasets included in SuperGLUE are in English. Additional details for each dataset are given in Table 8.

Preprocessing We do not perform any preprocessing, except shortening all examples to the maximum sequence length. This is done using the *longest first* strategy implemented in the Transformers library. All input sequences are truncated *before* applying patterns.