# True Few-Shot Learning with Prompts – A Real-World Perspective

## Timo Schick and Hinrich Schütze

Center for Information and Language Processing (CIS), LMU Munich, Germany

schickt@cis.lmu.de,inquiries@cislmu.org

## **Abstract**

Prompt-based approaches are strong at fewshot learning. However, Perez et al. (2021) have recently cast doubt on their performance because they had difficulty getting good results in a "true" few-shot setting in which prompts and hyperparameters cannot be tuned on a dev set. In view of this, we conduct an extensive study of PET, a method that combines textual instructions with examplebased finetuning. We show that, if correctly configured, PET performs strongly in a true few-shot setting, i.e., without a dev set. Crucial for this strong performance is PET's ability to intelligently handle multiple prompts. We then put our findings to a real-world test by running PET on RAFT, a benchmark of tasks taken directly from realistic NLP applications for which no labeled dev or test sets are available. PET achieves a new state of the art on RAFT and performs close to nonexpert humans for 7 out of 11 tasks. These results demonstrate that prompt-based learners like PET excel at true few-shot learning and underpin our belief that learning from instructions will play an important role on the path towards human-like few-shot learning capabilities.

## 1 Introduction

With pretrained language models (LMs) getting ever larger (Radford et al., 2019; Raffel et al., 2020; Brown et al., 2020; Fedus et al., 2021), *instruction-based learning* has emerged as a powerful method for few-shot text classification (e.g., Jiang et al., 2020; Schick and Schütze, 2021a,c; Brown et al., 2020; Wei et al., 2021; Sanh et al., 2021). The key idea is to give an LM access to descriptive names for all possible outputs and to short prompts explaining the task to be solved. In settings where at most a few dozen examples are available, this simple idea leads to substantial improvements over

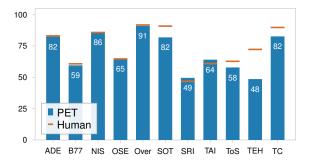


Figure 1: PET achieves near-human performance for 7 out of 11 tasks of the RAFT benchmark (Alex et al., 2021), for which labeled development and test sets are not available. This demonstrates that prompt-based learners like PET, if correctly configured, excel at true few-shot learning, i.e., without any tuning of instructions or hyperparameters on a development set.

various baselines (Schick and Schütze, 2021a,c; Gao et al., 2021; Tam et al., 2021).

However, recent work has questioned the strong few-shot performance of instruction-based approaches, arguing in particular that the considered settings are often not true few-shot settings (Perez et al., 2021; Logan IV et al., 2021) mainly for two reasons: For one, some approaches (e.g., Xie et al., 2019; Zhang et al., 2020; Chen et al., 2020; Tam et al., 2021) make use of large development sets to optimize hyperparameters. Beyond that, it is argued that manually designed instructions require manual tuning on development sets to achieve strong performance (Perez et al., 2021; Logan IV et al., 2021). Indeed, performance can vary largely - and in mostly unpredictable ways - across different instructions (Jiang et al., 2020; Schick and Schütze, 2021a); this issue even persists after finetuning a model on hundreds of instructions (Sanh et al., 2021). Even separate from this problem, the need for human involvement is generally seen as a huge drawback of manually designed instructions (Shin et al., 2020; Lester et al., 2021). Thus, several recent works abandon them in favor of automatically generated prompts (Shin et al., 2020; Gao et al., 2021; Hambardzumyan et al., 2021; Li and Liang, 2021; Lester et al., 2021).

Contrary to this trend, we argue that when correctly configured, prompt-based approaches achieve strong performance even in true few-shot settings and that there is no problem in using manually designed instructions per se. On the opposite, such instructions are often relatively easy to specify if one is familiar with the task to be solved, they provide an intuitive interface to convey task-specific knowledge, and if properly used, they consistently improve model performance in few-shot settings. To provide empirical support for these claims, we revisit PET (Schick and Schütze, 2021a) - a method for combining instructions with example-based finetuning whose key feature is that it allows users to specify multiple instructions for a single task - and thoroughly examine its performance with human-made instructions in true few-shot settings. In order to simulate a real-world scenario as best as possible, we proceed in two steps: First, we conduct an extensive study of PET using three English academic datasets to analyze its ability to perform true few-shot learning in a controlled environment and to derive best practices regarding the choice of instructions and other hyperparameters. We then put our findings to the test and evaluate PET on a large variety of different real-world tasks from the RAFT benchmark (Alex et al., 2021), for which no labeled development or test sets are available, enforcing a *true* few-shot setting (Perez et al., 2021). On average, PET clearly outperforms all baselines on this dataset and comes surprisingly close to the performance of non-expert humans (see Figure 1), demonstrating that instruction-based learning can successfully be applied to real-world tasks in true few-shot settings.

In summary, the main contributions of this work are as follows:

- We investigate the performance of PET for various models, tasks and training set sizes, its ability to cope with different instructions and its robustness to hyperparameter choices in true few-shot settings.
- We show how PET can be used when no unlabeled data is available and propose a variant for efficient classification in scenarios with many different classes.
- We apply PET to RAFT (Alex et al., 2021),

a benchmark of real-world tasks where it obtains a new state of the art and achieves near-human performance for 7 out of 11 tasks in true few-shot settings.

#### 2 Related Work

As a precursor to instruction-based learning, some works have investigated ways of informing classifiers about the meaning of different output classes both for text (Chang et al., 2008; Veeranna et al., 2016; Zhou et al., 2018) and image classification (Norouzi et al., 2014; Romera-Paredes and Torr, 2015); actually providing instructions in the form of short prompts was first proposed by Radford et al. (2019). This idea has since been applied to solve a wide range of different NLP tasks without any task-specific training data (Puri and Catanzaro, 2019; Opitz, 2019; Davison et al., 2019; Schick et al., 2021; Schick and Schütze, 2021; Wei et al., 2021; Sanh et al., 2021). While most approaches rephrase tasks as a language modeling problem, some use prompts to reformulate them as different tasks for which large amounts of training data are available (Levy et al., 2017; McCann et al., 2018; Yin et al., 2019; Sun et al., 2021; Sainz et al., 2021). Instruction-based learning has also been used in few-shot settings; popular variants include *in-context learning*, where the model's parameters are fixed and examples are provided as additional context (Brown et al., 2020; Lu et al., 2021; Kumar and Talukdar, 2021; Min et al., 2021), finetuning the entire model (Schick and Schütze, 2021a,c; Gao et al., 2021; Tam et al., 2021), and *prompt tuning*, where only the instruction itself is optimized (Shin et al., 2020; Hambardzumyan et al., 2021; Li and Liang, 2021; Lester et al., 2021).

Several works investigating the limitations and drawbacks of instruction-based few-shot approaches find that current LMs are mostly unable to understand complex instructions that go beyond short prompts or simple questions (Efrat and Levy, 2020; Weller et al., 2020; Webson and Pavlick, 2021) and that they are highly sensitive to the exact wording of the instructions provided (Jiang et al., 2020; Schick and Schütze, 2021a; Elazar et al., 2021). In a similar vein, Perez et al. (2021) and Logan IV et al. (2021) argue that prior work overestimates few-shot performance as manual prompt tuning is required to achieve good performance. Accordingly, some works attempt to obtain either prompts (Shin et al., 2020; Gao et al., 2021; Li and Liang, 2021; Lester et al., 2021) or meaningful



Figure 2: Different choices of patterns and corresponding verbalizers for classifying movie reviews as *positive* (+) or *negative* (-). The input is first converted into a cloze question using the pattern; classification is done by computing the output whose verbalization is the most likely substitute for the mask according to the MLM.

names for output classes (Schick et al., 2020; Gao et al., 2021) without any human involvement.

Finally, many benchmarks have been proposed for comparing few-shot approaches in a standardized way (e.g., Mishra et al., 2021; Bragg et al., 2021; Xu et al., 2021; Ye et al., 2021; Alex et al., 2021). As our focus is on the real-world applicability of few-shot methods, we evaluate PET on the RAFT benchmark (Alex et al., 2021), which measures performance in applied settings.

# 3 Pattern-Exploiting Training

We briefly review *pattern-exploiting training* (PET) (Schick and Schütze, 2021a,c), the method we use for instruction-based text classification. At its core, PET combines textual instructions with regular finetuning using labeled examples. To that end, users must specify one or more patterns which convert an input example x into a cloze question so that it can readily be processed by a masked language model (MLM) (Devlin et al., 2019). These patterns can take on very different forms; some examples are shown in Figure 2. In addition, users must inform the model about the meaning of all output classes; this is done with a verbalizer that assigns a natural language expression to each output y (see Figure 2, right). We refer to the combination of a pattern and verbalizer as a *pattern-verbalizer pair* (PVP).

Given a single PVP, let  $p(y \mid x)$  be the probability that an MLM assigns to y's verbalization in the cloze question obtained by applying the pattern to x, normalized over all y. The MLM is finetuned on labeled examples (x, y) by minimizing the crossentropy loss between  $p(y \mid x)$  and y.

If a user specifies multiple PVPs, individual mod-

els are trained for each pair. Similar to knowledge distillation (Hinton et al., 2015), they are then used to annotate unlabeled examples for training a final classifier with a regular sequence classification head (Devlin et al., 2019). We use the *weighted* variant of PET without auxiliary language modeling; see Schick and Schütze (2021a) for details.

Multi-Token Verbalizers When using encoder-only MLMs (Devlin et al., 2019; Liu et al., 2019; Lan et al., 2020), one limitation of PET is that the verbalization of each output class must correspond to a single token. Schick and Schütze (2021c) propose a fix for this that uses multiple mask tokens, but their approach is very inefficient. As discussed by Schick and Schütze (2021b), an alternative solution is to instead use an encoder-decoder LM (Lewis et al., 2020; Raffel et al., 2020). In Section 5, we propose yet another solution that allows us to stick with encoder-only MLMs while being much more efficient than the approach of Schick and Schütze (2021c).

Iterative PET We also experiment with iPET (Schick and Schütze, 2021a), an iterative variant of PET that employs self-training (e.g., Scudder, 1965; Yarowsky, 1995; Brin, 1999; McClosky et al., 2006) to train several generations of models on datasets of increasing size. To this end, an ensemble of MLMs is trained as in regular PET and then used to assign labels to unlabeled examples; a new ensemble is trained on the so-obtained data. This process is repeated for multiple iterations, where the number of annotated examples is increased in each iteration. We refer to Schick and Schütze (2021a) for more details.

## **4** True Few-Shot Learning with PET

We conduct a variety of experiments to answer six important research questions (Q1–Q6) regarding

<sup>&</sup>lt;sup>1</sup>We use the term *prompt* to refer to a short sequence of tokens that typically contains some form of instruction; the term *pattern* is used to denote the function that adds a prompt to an input.

the extent to which true few-shot learning is possible with PET. Beyond that, the purpose of our experiments is to establish a set of best practices for our real-world experiments on RAFT (Alex et al., 2021). Before discussing individual experiments, we describe the underlying setup.

Tasks and Datasets While they are heavily used in prior work (e.g., Brown et al., 2020; Schick and Schütze, 2021c; Logan IV et al., 2021; Webson and Pavlick, 2021), we decide against tasks and datasets from the GLUE (Wang et al., 2018) and SuperGLUE benchmarks (Wang et al., 2019) as they are very different from what we expect to see in real-world applications. Instead, we experiment with AG's News, Yelp Reviews Full Star and Yahoo Questions (Zhang et al., 2015) as these datasets represent classification tasks in three different domains that resemble real-world settings. Further, it is relatively easy to come up with a variety of instructions for each of these tasks, making it more straightforward to experiment with a large number of different patterns.

We consider settings with n=10 and n=100 training examples. For each n, we generate five different training sets per task by randomly sampling examples from the original training sets while ensuring that the number of examples is about the same for each possible output class. In addition, for both n=10 and n=100, we sample 1,000 unlabeled examples from the original training sets. We repeat all of our experiments for all five training sets and, by default, report average performance.

**PVPs** We manually write a total of 23 patterns per task, all of which can be categorized into one of the following groups:<sup>2</sup>

- NULL: Following Logan IV et al. (2021), these patterns simply insert a mask token.
- PUNC: Similar to some patterns of Schick and Schütze (2021a), these patterns only add punctuation characters and a mask token.
- PROMPTS: Patterns in this group add short prompts – typically consisting of no more than three words – to the input, similar to Radford et al. (2019) and Schick and Schütze (2021a).
- Q&A: These patterns rephrase the task as a question *q* and append

Question: q Answer: [MASK].

to the input, similar to Brown et al. (2020) and Schick et al. (2021).

For all patterns, we use only a single verbalizer which we adopt from Schick and Schütze (2021a). While varying the verbalizer may also lead to interesting insights, finding a large amount of reasonable verbalizers is challenging for some tasks as there is often a single, *natural* choice (e.g., the actual category names for AG's News and Yahoo Questions).

Hyperparameters We consider a setting similar to that of Schick and Schütze (2021a,c) and, unless otherwise specified, use the default settings of the PET library.<sup>3</sup> As our experiments require training hundreds of models, we make a few changes to reduce environmental impact (Strubell et al., 2019) and computational cost: We use the base variant of RoBERTa (Liu et al., 2019) as underlying LM, we train only one model per PVP, and we reduce the training steps for all individual models and the final classifier to 100 and 1,000, respectively.

**Monitoring** Finetuning pretrained LMs can be unstable on small datasets (Devlin et al., 2019; Dodge et al., 2020), sometimes leading to very poor performance. Luckily, we can detect such finetuning issues to some extent even without a labeled test set using the following two checks:

- TRAIN SET UNDERFITTING: We check for training runs that result in less than 50% accuracy on the training set. As finetuning on up to 100 examples typically leads to perfect predictions on the training set, this is a clear indicator of a failed finetuning run.
- CONSTANT PREDICTIONS: Another strong indicator of unsuccessful training is when the finetuned model predicts the same output class for all inputs; we check this both on the training set and on the unlabeled set.

Whenever one of these events occurs, we restart training using a different seed for initializing all random number generators.

# Q1: How can we find *the* best pattern – or do we even need to?

Even slightly different patterns can lead to very different performance (Jiang et al., 2020; Schick and

<sup>&</sup>lt;sup>2</sup>The full set of PVPs can be found in Appendix A.

<sup>&</sup>lt;sup>3</sup>See https://github.com/timoschick/pet

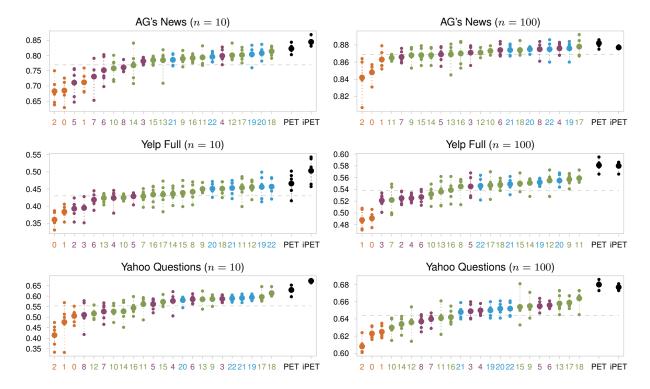


Figure 3: Performance of individual patterns, PET and iPET on all tasks considered. Accuracy is shown on the y-axis; the x-axis shows individual pattern ids where color is used to distinguish the different pattern categories (Null, Punc, Prompt, Q&A). Small bullets ( $\bullet$ ) correspond to individual training sets, large bullets ( $\bullet$ ) correspond to average performance. Average performance across all patterns is shown as a dashed gray line.

Schütze, 2021a; Schick et al., 2021; Webson and Pavlick, 2021; Sanh et al., 2021, i.a.) and popular model selection criteria can not reliably identify patterns that achieve similar performance to the best one (Perez et al., 2021). We thus investigate to what extent PET can eliminate the need to find the best instruction even in extreme settings where there are dozens of candidates to choose from.

**Setup** Using our default setup, we train individual models for each PVP and a final PET model; we also train an iPET model for 3 iterations.

**Results** Performance of individual models for each pattern and of the distilled models obtained using PET and iPET is shown in Figure 3. Interestingly, sorting all pattern groups by their average performance gives the exact same order for each task and training set size: NULL patterns clearly perform worst, followed by PUNC and PROMPT; Q&A gives the best average results. Contrary to findings of Logan IV et al. (2021), this shows that LMs can benefit a lot from manually written instructions even if combined with finetuning.

Crucially, PET's performance is much higher than average performance of individual patterns; further, it consistently outperforms even the best pattern, verifying that PET indeed removes the need to find *the* best pattern. While iPET gives clear improvements for n=10, it performs worse than PET for n=100. The reason for this may be that we use a much smaller set of unlabeled examples than prior work (Schick and Schütze, 2021a,c).

# Q2: Does performance of different patterns transfer across models?

While our results for Q1 show a consistent order of pattern groups for different training set sizes and tasks, an important question for real-world applications is whether the same finding also holds for different model sizes, and even entirely different models.

**Setup** We consider BERT (Devlin et al., 2019) RoBERTa (Liu et al., 2019) and ALBERT (Lan et al., 2020) as underlying language models;<sup>4</sup> we experiment with the base and large variants. For each model and size, we repeat the exact same experiment as for Q1.

<sup>&</sup>lt;sup>4</sup>For Yahoo, we do not consider BERT as it uses a vocabulary that does not assign a single token to each verbalization.

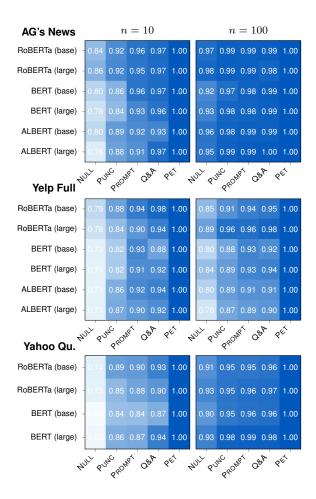


Figure 4: Relative performance of individual pattern groups and PET for different models and sizes. Scores are normalized so that the best performance for each task, number of examples, model and size is 1.0.

Results Figure 4 shows the performance of each pattern group (i.e., average performance of all individual patterns in this group) and PET; scores are normalized so that the best-performing approach for each task, training set size and model gets a score of 1.0. Except for very few exceptions, our findings from Q1 regarding the relative performance of pattern groups and PET (NULL < PUNC < PROMPT < Q&A < PET) also hold for different models and sizes. In general, the performance of individual patterns strongly correlates between different models and sizes (Spearman's  $\rho \geq 0.7$  except in one case).

# Q3: Does PET still work if some PVPs are not well understood?

While Q1 and Q2 have shown that PET performs even better than the best PVP if a large set of highquality PVPs is given, a potential concern is that

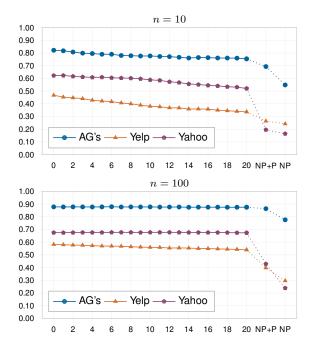


Figure 5: Performance of PET with three randomly selected patterns when adding noise PVPs; the *x*-axis shows the number of noise PVPs added. We also show performance of using *only* noise PVPs with PET (NP+P) and their average performance (NP).

performance may be much worse if the LM fails to understand a fair amount of patterns and verbalizers (e.g., because they are in a very different style from the model's pretraining data). For real-world scenarios, it is thus relevant to know how such "bad" PVPs affect the performance of PET.

**Setup** It is difficult to obtain large quantities of bad instructions as they might occur in real-world scenarios. As a proxy, we resort to *noise patterns* that add random tokens to the input, serving as a lower bound for truly bad patterns. In concrete terms, we add up to three randomly sampled tokens before and after the input. We also create *noise verbalizers* by assigning uniformly selected tokens to each output class. Using this process, we obtain 20 different intentionally bad PVPs per task. For each task, we start with 3 randomly selected, high-quality patterns from our original set of manually designed instructions, add noise PVPs one by one and investigate the effect on performance.

**Results** The effect of adding noise PVPs is shown in Figure 5. Interestingly, for both n=10 and n=100, performance remains almost con-

<sup>&</sup>lt;sup>5</sup>If there are multiple input texts, we shuffle their order and additionally add 0–3 tokens in between them.

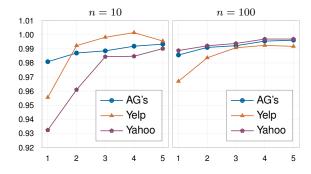


Figure 6: Relative performance of PET with only a subset of patterns compared to that achieved using all 23 manually designed patterns. The *x*-axis shows the number of patterns used.

stant even if more than half of the used PVPs are noise PVPs, demonstrating that PET is very robust to "bad" instructions. Figure 5 also shows performance when using *only* noise PVPs; except for AG's News with n=100, this leads to substantially worse performance than PET with a small amount of manually designed instructions.

# **Q4:** How many patterns are required for good performance?

Orthogonal to Q3, it is also important to know how many *high-quality* prompts are at least required to achieve satisfactory performance. This is of great practical importance because coming up with dozens of different PVPs for a task may take a significant amount of time, that could otherwise be spent annotating further examples.

**Setup** We generate 10 random permutations of all 23 patterns per task. For each permutation and training set, we use the same setup as in Q1 to compute the average performance obtained with PET when using only the first i patterns, where i ranges from 1 to 5.

**Results** Average performance of PET trained with the first i patterns is shown in Figure 6, relative to the performance of PET trained with all 23 patterns. For all tasks and training set sizes, as little as four different patterns are already sufficient to achieve performance very close to that of PET trained with all 23 patterns. Surprisingly, PET's performance is much higher than the average performance of a model trained on individual patterns even with i=1. This indicates that the process of knowledge distillation using unlabeled data is also beneficial when using only a single instruction.

# **Q5:** How important are values for other hyperparameters?

For true few-shot settings, the same set of hyperparameter values should ideally achieve good performance across different tasks; this enables us to adopt these values for new tasks without requiring manual tuning on task-specific validation sets. We therefore investigate how different choices for common hyperparameters affect the performance of PET; in concrete terms, we consider learning rate, training steps and batch size.

**Setup** Based on choices found in previous work, we consider different learning rates from the set  $\{10^{-4}, 5 \cdot 10^{-4}, 10^{-5}, 5 \cdot 10^{-5}, 10^{-6}\}$ , training steps from  $\{10, 25, 50, 100, 250, 500, 1000\}$  and batch sizes from  $\{1, 2, 4, 8, 16, 32\}$ . Learning rate and batch size are changed for the individual models and the final classifier simultaneously; the number of training steps is varied only for individual models. We modify each hyperparameter independently, keeping all other parameters at their default values (i.e., a learning rate of  $10^{-5}$ , 100 steps and a batch size of 4).

**Results** All results are shown in Figure 7. For training steps and batch size, performance is relatively stable across a wide range of different values, with more steps and larger batch sizes typically leading to slightly better performance (especially for n=100). Learning rate clearly has the strongest impact on performance, but values of  $10^{-5}$  and  $5 \cdot 10^{-5}$  consistently give the best results across all tasks considered; those are also the values typically used for finetuning in prior work (Devlin et al., 2019; Liu et al., 2019).

# Q6: Do we really need large amounts of unlabeled data?

One drawback of PET compared to using individual PVPs is the additional need for unlabeled data; while such data is often available in real-world settings, there are cases in which even unlabeled examples are hard to obtain. As recent work shows that pretrained LMs are capable of generating data of reasonable quality themselves (Anaby-Tavor et al., 2020; Papanikolaou and Pierleoni, 2020; Yang et al., 2020; Mohapatra et al., 2020; Kumar et al., 2021; Schick and Schütze, 2021), we investigate whether unlabeled data can simply be replaced by synthetic examples.

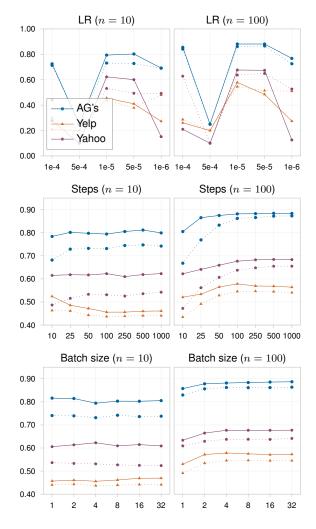


Figure 7: Performance of PET (solid lines) and average performance of individual models (dotted lines) for different learning rates (LR), training steps (Steps), and batch sizes. For reasons of readability, the legend is shown in the top left plot only.

**Setup** We use GPT2-XL (Radford et al., 2019) to generate synthetic unlabeled data. To this end, we provide one or two randomly chosen training examples without labels as *in-context examples* (Brown et al., 2020) and let the model generate one additional example. For two inputs  $x_1$  and  $x_2$ , the input given to the model is

Example 1:  $x_1 \leftarrow$  Example 2:  $x_2 \leftarrow$  Example 3:

where ← denotes two consecutive line breaks. If an input consists of two texts, we simply concatenate them using the sequence +++ as a separator.

We generate 10,000 examples for n=10 and 30,000 examples for n=100 using top-p sampling (Holtzman et al., 2020) with p=0.9. For each input, we stop the generation process as soon as the model generates two consecutive line breaks.

We discard all examples for which the model does not generate two consecutive line breaks within 128 tokens; for datasets with text pairs, we also discard examples where the model fails to generate the sequence separator (+++).

As the datasets obtained with this method may be highly imbalanced regarding the distribution of (unknown) labels, we also experiment with a *balanced* variant: We use the ensemble of models trained on individual PVPs to assign labels to each example and only keep so many examples per label that the resulting dataset – which is used for training the final classifier – is balanced.

**Results** Figure 8 shows the performance of individual patterns as well as PET and iPET both with real and synthetic unlabeled data. Except for iPET on Yahoo Questions with n=10, using synthetic data consistently achieves accuracies within one point of using real data, with our balanced version performing slightly better. Moreover, for n=10 using synthetic data even improves accuracy in some cases. This shows that in the absence of unlabeled examples, synthetic data obtained from generative language models can serve as a drop-in replacement without substantially degrading performance.

## 5 PET for Real-World Tasks

We use our insights from Section 4 to apply PET to the RAFT benchmark, a collection of 11 diverse real-world tasks whose automated solution has inherent value to someone (Alex et al., 2021). These tasks pose various challenges to few-shot approaches as they require some amount of domain expertise and the ability to understand detailed instructions, to process long inputs and to handle a large number of output classes.

Tasks and Datasets The RAFT benchmark includes 11 tasks from different domains: ADE Corpus V2 (ADE), Banking77 (B77), NeurIPS impact statement risks (NIS), OneStopEnglish (OSE), Overruling (Over), Semiconductor org types (SOT), Systematic review inclusion (SRI), TAI safety research (TAI), Terms of Service (ToS), TweetEval Hate (TEH) and Twitter complaints (TC). For a detailed overview of all tasks, we refer to Alex et al. (2021). Each task comes with 50 labeled training examples; in accordance with the RAFT rules, we additionally make use of the unlabeled data (ranging from 150 to 5,000 examples) for PET's distil-

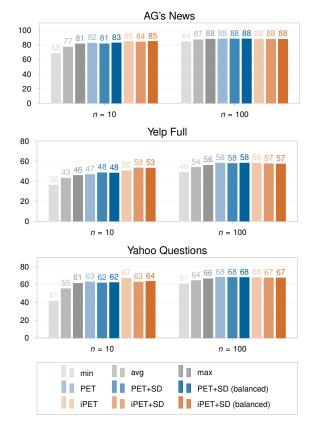


Figure 8: Minimum, average and maximum performance of individual patterns compared to regular PET and iPET as well as PET and iPET with synthetic data (+SD). Accuracies are scaled by a factor of 100 for better readability.

lation step. Note that the unlabeled set for RAFT is the same as the test set, meaning that unlike in Section 4, our final classifier is directly trained on (unlabeled) test examples.

**PVPs** Based on our findings from Q1 and Q2, we exclusively specify Q&A prompts for each task. To obtain the question q, we make minimal changes to the original instructions of Alex et al. (2021); we rephrase all binary classification tasks as yes/no questions to facilitate finding a suitable verbalizer. For example, we rephrase the instruction "Label the sentence based on whether it is related to an adverse drug effect (ADE)." as the question "Is this sentence related to an adverse drug effect (ADE)?" Following our results from Q4, we specify 4 PVPs per task. In case of binary classification, we use two different patterns that either include or omit the full task specification of Alex et al. (2021)<sup>6</sup> and combine them with both a yes/no verbalizer and a

true/false verbalizer. The full set of PVPs for all tasks can be found in Appendix B.

**Hyperparameters** We mostly keep the default values for hyperparameters used throughout Section 4 as our experiments for Q5 show that these perform well for all tasks considered. However, we make the following changes:

- We replace RoBERTa (base) with ALBERT (xxlarge, v2). While being much slower to train, ALBERT was shown to outperform RoBERTa both in regular and few-shot settings (Lan et al., 2020; Schick and Schütze, 2021c; Logan IV et al., 2021).
- For tasks with more than 4,000 unlabeled examples, we finetune the distilled model for 2,000 steps instead of 1,000 steps. Otherwise, some examples would not be seen at all during training.<sup>7</sup>
- Following Schick and Schütze (2021a,c) we train three individual models per PVP. This improves robustness as performance can vary largely between individual finetuning runs.

**Handling Many Labels** The B77 dataset contains online banking customer service queries annotated with one of 77 possible intents. This large amount of different outputs leads to several issues for PET: First, it is impossible to specify a meaningful verbalizer that maps each intent to a single token. We initially experimented with the multimask version of PET (Schick and Schütze, 2021c), but found it to be too inefficient to get results in a reasonable amount of time. Therefore, we tried the following solution: We rephrase the task as binary classification, where for each pair of query x and intent y, the task is to decide whether y is the correct intent for x. For each original training example (x, y), we generate one example (x, y, True) and four examples (x, y', False) with randomly sampled, wrong intents y'. As this increases the amount of data fivefold, we finetune each individual model for 500 steps instead of 100 steps.

While this approach solves our problem, it is still not particularly efficient: Reframing the task as binary classification means that for each input, 77 forward passes are required to find the correct

<sup>&</sup>lt;sup>6</sup>For a full list of all task specifications, see https://github.com/oughtinc/raft-baselines.

<sup>&</sup>lt;sup>7</sup>Training for 1,000 steps with a batch size of 4 means that at most 4,000 examples can be processed.

| Method   | ADE  | B77         | NIS         | OSE         | Over        | SOT  | SRI         | TAI         | ToS         | TEH         | TC   | Avg         |
|----------|------|-------------|-------------|-------------|-------------|------|-------------|-------------|-------------|-------------|------|-------------|
| GPT-2    | 60.0 | 12.1        | 56.1        | 24.5        | 49.8        | 38.0 | 49.2        | 61.2        | 49.8        | 31.1        | 72.3 | 45.8        |
| GPT-Neo  | 45.2 | 14.9        | 40.8        | 34.3        | 68.1        | 40.6 | 49.3        | 60.5        | 56.5        | 55.4        | 63.6 | 48.1        |
| AdaBoost | 54.3 | 02.3        | 62.6        | 47.5        | 83.8        | 45.5 | 50.6        | 55.6        | 56.0        | 44.3        | 62.5 | 51.4        |
| snlt     | 60.3 | 24.8        | 58.5        | 30.2        | 83.1        | 33.6 | 49.2        | 62.6        | 54.0        | 44.9        | 79.1 | 52.8        |
| GPT-3    | 68.6 | 29.9        | 67.9        | 43.1        | <u>93.7</u> | 76.9 | <u>51.6</u> | <u>65.6</u> | 57.4        | 52.6        | 82.1 | 62.7        |
| SetFit   | 72.6 | 53.8        | <u>87.2</u> | 52.1        | 90.7        | 68.2 | 49.3        | 62.8        | 62.0        | 53.2        | 83.7 | 66.9        |
| PET      | 82.2 | 59.3        | 85.7        | <u>64.6</u> | 90.8        | 81.6 | 49.3        | 63.8        | 57.6        | 48.3        | 82.4 | 69.6        |
| Human    | 83.0 | <u>60.7</u> | 85.7        | <u>64.6</u> | 91.7        | 90.8 | 46.8        | 60.9        | <u>62.7</u> | <u>72.2</u> | 89.7 | <u>73.5</u> |

Table 1: Performance of various baselines and PET on the RAFT benchmark (Alex et al., 2021); shown numbers are macro F1 scores multiplied by 100. Best model performance is shown in bold, best overall performance (including human annotators) is underlined. The final column shows average performance across all 11 tasks.

intent. We thus train the final model as a regular classifier with 77 different output classes; for training this classifier on an input x, we set the target probability of each output y proportional to the probability of True being the correct output for (x,y) according to our ensemble of binary classifiers.

Finally, another issue is that with 50 labeled examples, at least 27 labels are not covered in the training set at all; this may bias a model to never predict any of these labels. To alleviate this issue, we train two generations of models using iPET. For training the second generation, we obtain training data covering all possible labels as follows: For each label, we pick the two examples from our set of unlabeled data for which this label is most likely according to the first generation; the same approach was used by Schick and Schütze (2021a) for applying iPET in zero-shot settings.

Of course, the nature of RAFT makes it impossible to measure the impact of any of these choices. While we could conduct experiments similar to those in Section 4, none of the datasets considered therein has a similar structure to B77; as our modifications affect only one out of 11 tasks, we thus decided to not perform any further analysis.

**Monitoring** We checked for TRAIN SET UNDERFITTING and CONSTANT PREDICTIONS as in Section 4 to detect finetuning issues. Unlike for our experiments in Section 4, on RAFT we encountered some issues that could *not* be resolved simply by retraining with a different seed:

 We observed TRAIN SET UNDERFITTING for the final classifier on B77. This may be due to the classification head for 77 classes introducing many new parameters; we tried training the final model for 5,000 steps instead of 2,000 steps, which fixed this issue.

- We observed CONSTANT PREDICTIONS for the ToS training set. Doubling the number of training steps resolved this problem.
- Finally, we also observed CONSTANT PRE-DICTIONS on the unlabeled data of SRI. Upon manually inspecting the training set, we observed that all but one out of 50 examples have the same label. As all models already classified the training set perfectly, we left the setup for our SRI submission unchanged.

**Results** For all 11 tasks, results of PET and various baselines are shown in Table 1.8 As can be seen, PET performs better than all other approaches on average, achieving near-human performance for 7 out of 11 tasks. Note however that non-expert humans perform worse than a majority baseline SRI, so results on this task should be taken with a grain of salt. PET also clearly outperforms a GPT-3 model (Brown et al., 2020) by almost 7 points, despite the latter being larger by several orders of magnitude. 9 While PET is particularly successful on ADE, B77 and OSE (where it outperforms GPT-3 by 13.6, 21.5 and 29.4 points, respectively), it performs comparably bad on datasets in the law (Over, ToS) and social media (TEH, TC) domain. Our approach for handling many labels performs surprisingly well on B77 without any tuning of

<sup>&</sup>lt;sup>8</sup>All results are taken directly from the leaderboard at https://huggingface.co/spaces/ought/raft-leaderboard.

<sup>&</sup>lt;sup>9</sup>We were unable to find any information on SetFit, the second-best performing method (e.g., which underlying LM is used or whether it makes use of any additional data), so we cannot make a fair comparison.

its parameters. Due to the nature of the RAFT benchmark, we cannot perform further analysis or ablation studies.

### 6 Discussion

Our experimental results in Section 4 and 5 show that strong performance in few-shot settings is clearly possible without manual prompt tuning or hyperparameter optimization on large development sets; in other words, PET can successfully be applied in true few-shot settings. While we believe that it should be an important goal of future work to make LMs more robust to different instructions, even with current models it is relatively easy to successfully apply PET when following a few simple principles - such as rephrasing the task in a Q&A format, using simple vocabulary and singletoken verbalizers where possible, and specifying at least a handful of different patterns. In light of these findings, we also hope that future work will not view human involvement in prompt design as a drawback of instruction-based approaches, but rather as an exciting possibility to communicate with models in ways other than exclusively through examples.

There are various limitations to our study. First, a major obstacle to actually applying PET in realworld applications is that we do not know a priori how well it performs for a given task; we therefore believe an important next step is to investigate methods for estimating performance without access to large test sets – for example, through model calibration (Desai and Durrett, 2020; Jiang et al., 2021) - in real-world settings. In addition, we did not fully explore the capabilities of PET; for example, we did not investigate domain-adaptive pretraining (Gururangan et al., 2020) and auxiliary language modeling (Chronopoulou et al., 2019), both of which were shown to be helpful by Schick and Schütze (2021a). We also did not quantify the impact of our decisions regarding B77 and the effectiveness of our monitoring and only considered English models and datasets. Finally, we did not examine PET's performance beyond aggregate scores. While this is not feasible on RAFT due to the nature of this dataset, performing such analysis either with other datasets or with methods such as the ones proposed by Ribeiro et al. (2020) would be relevant future work to understand real-world capabilities of instruction-based approaches more comprehensively.

## 7 Conclusion

In light of recent work casting doubt on the performance of prompt-based approaches in true few-shot settings (Perez et al., 2021), we have conducted an extensive study of PET. In a controlled environment, we found that manually designed instructions consistently outperform null prompts, with Q&Astyle prompts performing best (Q1, Q2). Across different tasks, models and training set sizes, PET consistently outperforms even the best individual prompt (Q1, Q2). We have also shown that PET is robust to uninformative prompts and to different choices of hyperparameters (Q3, Q5), that as little as four prompts are sufficient to reach good performance (O4), and that synthetic examples can be used to replace large amounts of unlabeled data (Q6). On the basis of these insights, we applied PET to a benchmark of real-world tasks, where it achieves near-human performance for 7 out of 11 tasks without any tuning on a development set, demonstrating the power of instruction-based approaches in true few-shot settings.

#### References

Neel Alex, Eli Lifland, Lewis Tunstall, Abhishek Thakur, Pegah Maham, C. Jess Riedel, Emmie Hine, Carolyn Ashurst, Paul Sedille, Alexis Carlier, Michael Noetel, and Andreas Stuhlmüller. 2021. RAFT: A real-world few-shot text classification benchmark. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.

Ateret Anaby-Tavor, Boaz Carmeli, Esther Goldbraich, Amir Kantor, George Kour, Segev Shlomov, Naama Tepper, and Naama Zwerdling. 2020. Do not have enough data? Deep learning to the rescue! *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):7383–7390.

Jonathan Bragg, Arman Cohan, Kyle Lo, and Iz Beltagy. 2021. FLEX: Unifying evaluation for few-shot NLP. In *Thirty-Fifth Conference on Neural Information Processing Systems*.

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.

Ming-Wei Chang, Lev Ratinov, Dan Roth, and Vivek Srikumar. 2008. Importance of semantic representation: Dataless classification. In *Proceedings of the 23rd National Conference on Artificial Intelligence - Volume 2*, AAAI'08, page 830–835. AAAI Press.

Jiaao Chen, Zichao Yang, and Diyi Yang. 2020. MixText: Linguistically-informed interpolation of hidden space for semi-supervised text classification. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2147–2157, Online. Association for Computational Linguistics.

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.

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.

Shrey Desai and Greg Durrett. 2020. Calibration of pre-trained transformers. In *Proceedings of* 

the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 295–302, Online. Association for Computational Linguistics.

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.

Avia Efrat and Omer Levy. 2020. The turking test: Can language models understand instructions? *Computing Research Repository*, arXiv:2010.11982.

Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhilasha Ravichander, Eduard Hovy, Hinrich Schütze, and Yoav Goldberg. 2021. Measuring and improving consistency in pretrained language models. *Computing Research Repository*, arXiv:2102.01017.

William Fedus, Barret Zoph, and Noam Shazeer. 2021. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Computing Research Repository*, arXiv:2101.03961.

Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. Making pre-trained language models better fewshot learners. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3816–3830, Online. Association for Computational Linguistics.

Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. 2020. Don't stop pretraining: Adapt language models to domains

- and tasks. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8342–8360, Online. Association for Computational Linguistics.
- Karen Hambardzumyan, Hrant Khachatrian, and Jonathan May. 2021. WARP: Word-level Adversarial ReProgramming. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4921–4933, Online. Association for Computational Linguistics.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *Computing Research Repository*, arXiv:1503.02531.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration. In *International Conference on Learning Representations*.
- Zhengbao Jiang, Jun Araki, Haibo Ding, and Graham Neubig. 2021. How Can We Know When Language Models Know? On the Calibration of Language Models for Question Answering. *Transactions of the Association for Computational Linguistics*, 9:962–977.
- 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.
- Sawan Kumar and Partha Talukdar. 2021. Reordering examples helps during priming-based few-shot learning. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4507–4518, Online. Association for Computational Linguistics.
- Varun Kumar, Ashutosh Choudhary, and Eunah Cho. 2021. Data augmentation using pretrained transformer models. *Computing Research Repository*, arXiv:2003.02245.
- 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.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. *Computing Research Repository*, arXiv:2104.08691.
- Omer Levy, Minjoon Seo, Eunsol Choi, and Luke Zettlemoyer. 2017. Zero-shot relation extraction via reading comprehension. In *Proceedings* of the 21st Conference on Computational Natural Language Learning (CoNLL 2017), pages 333–342, Vancouver, Canada. Association for Computational Linguistics.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training 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.
- Xiang Lisa Li and Percy Liang. 2021. Prefixtuning: Optimizing continuous prompts for generation.
- 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.
- 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. *Computing Research Repository*, 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. *Computing Research Repository*, arXiv:2104.08786.
- Bryan McCann, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. 2018. The natural

- language decathlon: Multitask learning as question answering. *Computing Research Repository*, arXiv:1806.08730.
- 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.
- Sewon Min, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2021. Noisy channel language model prompting for few-shot text classification. *Computing Research Repository*, arXiv:2108.04106.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2021. Cross-task generalization via natural language crowdsourcing instructions. *Computing Research Repository*, arXiv:2104.08773.
- Biswesh Mohapatra, Gaurav Pandey, Danish Contractor, and Sachindra Joshi. 2020. Simulated chats for task-oriented dialog: Learning to generate conversations from instructions. *Computing Research Repository*, arXiv:2010.10216.
- Mohammad Norouzi, Tomas Mikolov, Samy Bengio, Yoram Singer, Jonathon Shlens, Andrea Frome, Greg S. Corrado, and Jeffrey Dean. 2014. Zero-shot learning by convex combination of semantic embeddings.
- 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.
- Yannis Papanikolaou and Andrea Pierleoni. 2020. DARE: Data augmented relation extraction with GPT-2. *Computing Research Repository*, arXiv:2004.13845.
- Ethan Perez, Douwe Kiela, and Kyunghyun Cho. 2021. True few-shot learning with language models. In *Thirty-Fifth Conference on Neural Information Processing Systems*.

- Raul Puri and Bryan Catanzaro. 2019. Zeroshot text classification with generative language models. *Computing Research Repository*, arXiv:1912.10165.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. Technical report, Open AI.
- 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.
- Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. Beyond accuracy: Behavioral testing of NLP models with CheckList. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4902–4912, Online. Association for Computational Linguistics.
- Bernardino Romera-Paredes and Philip Torr. 2015. An embarrassingly simple approach to zero-shot learning. In *International Conference on Machine Learning*, pages 2152–2161.
- Oscar Sainz, Oier Lopez de Lacalle, Gorka Labaka, Ander Barrena, and Eneko Agirre. 2021. Label verbalization and entailment for effective zero and few-shot relation extraction. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1199–1212, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Fevry, Jason Alan Fries, Ryan Teehan, Stella Biderman, Leo Gao, Tali Bers, Thomas Wolf, and Alexander M. Rush. 2021. Multitask

- prompted training enables zero-shot task generalization. *Computing Research Repository*, arXiv:2110.08207.
- Timo Schick, Helmut Schmid, and Hinrich Schütze. 2020. Automatically identifying words that can serve as labels for few-shot text classification. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5569–5578, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Timo Schick and Hinrich Schütze. 2021. Generating datasets with pretrained language models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics.
- Timo Schick and Hinrich Schütze. 2021a. 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.
- Timo Schick and Hinrich Schütze. 2021b. Fewshot text generation with pattern-exploiting training. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics.
- Timo Schick and Hinrich Schütze. 2021c. It's not just size that matters: Small language models are also few-shot learners. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2339–2352, Online. Association for Computational Linguistics.
- Timo Schick, Sahana Udupa, and Hinrich Schütze. 2021. Self-diagnosis and self-debiasing: A proposal for reducing corpus-based bias in NLP. *Transactions of the 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.

- Taylor Shin, Yasaman Razeghi, Robert L. Logan IV,
  Eric Wallace, and Sameer Singh. 2020. AutoPrompt: Eliciting Knowledge from Language
  Models with Automatically Generated Prompts.
  In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 4222–4235, Online. Association for Computational Linguistics.
- Emma Strubell, Ananya Ganesh, and Andrew Mc-Callum. 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.
- Yi Sun, Yu Zheng, Chao Hao, and Hangping Qiu. 2021. NSP-BERT: A prompt-based zero-shot learner through an original pre-training tasknext sentence prediction. *Computing Research Repository*, arXiv:2109.03564.
- Derek Tam, Rakesh R Menon, Mohit Bansal, Shashank Srivastava, and Colin Raffel. 2021. Improving and simplifying pattern exploiting training. *Computing Research Repository*, arXiv:2103.11955.
- Sappadla Prateek Veeranna, Jinseok Nam, Eneldo Loza Mencia, and Johannes Fürnkranz. 2016. Using semantic similarity for multi-label zero-shot classification of text documents. In *Proceeding of European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning. Bruges, Belgium: Elsevier*, pages 423–428.
- 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.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, Brussels,

- Belgium. Association for Computational Linguistics.
- Albert Webson and Ellie Pavlick. 2021. Do prompt-based models really understand the meaning of their prompts? *Computing Research Repository*, arXiv:2109.01247.
- Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2021. Finetuned language models are zero-shot learners.
- Orion Weller, Nicholas Lourie, Matt Gardner, and Matthew E. Peters. 2020. Learning from task descriptions. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1361–1375, Online. Association for Computational Linguistics.
- Qizhe Xie, Zihang Dai, Eduard Hovy, Minh-Thang Luong, and Quoc V. Le. 2019. Unsupervised data augmentation for consistency training. *Computing Research Repository*, arXiv:1904.12848.
- Liang Xu, Xiaojing Lu, Chenyang Yuan, Xuanwei Zhang, Huilin Xu, Hu Yuan, Guoao Wei, Xiang Pan, Xin Tian, Libo Qin, and Hu Hai. 2021. Few-CLUE: A chinese few-shot learning evaluation benchmark. *Computing Research Repository*, arXiv:2107.07498.
- Yiben Yang, Chaitanya Malaviya, Jared Fernandez, Swabha Swayamdipta, Ronan Le Bras, Ji-Ping Wang, Chandra Bhagavatula, Yejin Choi, and Doug Downey. 2020. Generative data augmentation for commonsense reasoning. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1008–1025, Online. Association for Computational Linguistics.
- 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.
- Qinyuan Ye, Bill Yuchen Lin, and Xiang Ren. 2021. CrossFit: A few-shot learning challenge for cross-task generalization in NLP. *Computing Research Repository*, arXiv:2104.08835.

- Wenpeng Yin, Jamaal Hay, and Dan Roth. 2019. Benchmarking zero-shot text classification: Datasets, evaluation and entailment approach. 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 3914–3923, Hong Kong, China. Association for Computational Linguistics.
- Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter Liu. 2020. PEGASUS: Pre-training with extracted gap-sentences for abstractive summarization. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 11328–11339, Virtual. PMLR.
- Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015.
  Character-level convolutional networks for text classification. In C. Cortes, N. D. Lawrence,
  D. D. Lee, M. Sugiyama, and R. Garnett, editors, Advances in Neural Information Processing Systems 28, pages 649–657. Curran Associates, Inc.
- Ben Zhou, Daniel Khashabi, Chen-Tse Tsai, and Dan Roth. 2018. Zero-shot open entity typing as type-compatible grounding. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2065–2076, Brussels, Belgium. Association for Computational Linguistics.

## A PVPs for AG's, Yelp, Yahoo

Below, we list the 23 patterns and the single verbalizer used for each task throughout Section 4. We group all patterns by their category.

#### A.1 AG's News

## **Inputs**

- $x_1$ : The headline of the news article to be classified.
- $x_2$ : The text body of the news article to be classified.

## **NULL Patterns**

- [MASK]  $x_1 x_2$
- $x_1 x_2$  [MASK]
- $x_1$  [MASK]  $x_2$

## **PUNC Patterns**

- [MASK] :  $x_1 x_2$
- [MASK]  $x_1 x_2$
- [MASK] .  $x_1 x_2$
- ( [MASK] )  $x_1 x_2$
- $x_1$  ( [MASK] )  $x_2$
- $x_1 x_2$  ( [MASK] )

#### **PROMPT Patterns**

- $x_1 x_2$  Category: [MASK].
- *x*<sub>1</sub> *x*<sub>2</sub> Class: [MASK].
- $x_1 x_2$  Topic: [MASK].
- $x_1 x_2$  Theme: [MASK].
- x<sub>1</sub> x<sub>2</sub> Category: [MASK]
- *x*<sub>1</sub> *x*<sub>2</sub> Class: [MASK]
- *x*<sub>1</sub> *x*<sub>2</sub> Topic: [MASK]
- $x_1 x_2$  Theme: [MASK]
- [MASK] News:  $x_1 x_2$
- [MASK] NEWS:  $x_1 x_2$

#### **Q&A Patterns**

- x<sub>1</sub> x<sub>2</sub> Question: What is the topic of this article? Answer: [MASK].
- x<sub>1</sub> x<sub>2</sub> Question: What is the category of this article? Answer: [MASK].
- x<sub>1</sub> x<sub>2</sub> Question: What is the topic of this article? Answer: [MASK]
- x<sub>1</sub> x<sub>2</sub> Question: What is the category of this article? Answer: [MASK]

**Verbalizer** World  $\mapsto$  World, Sports  $\mapsto$  Sports, Business  $\mapsto$  Business, Science/Tech  $\mapsto$  Tech

## A.2 Yelp Reviews Full Star

# **Inputs**

• x: The review to be classified.

## **NULL Patterns**

- [MASK] x
- *x* [MASK]

## **PUNC Patterns**

- [MASK] : x
- [MASK] x
- [MASK] . x
- ( [MASK] ) x
- x ( [MASK] )

#### **PROMPT Patterns**

- x It was [MASK]
- x All in all, it was [MASK]
- x In summary, it was [MASK]
- x It was [MASK].
- x All in all, it was [MASK].
- x In summary, it was [MASK].
- x The restaurant was [MASK]
- x All in all, the restaurant was [MASK]
- x In summary, the restaurant was [MASK]
- x The restaurant was [MASK].
- x All in all, the restaurant was [MASK].
- x In summary, the restaurant was [MASK].

## **Q&A Patterns**

- x Question: What does the customer think of this restaurant? Answer: It is [MASK].
- x Question: What does the customer think of this place? Answer: It is [MASK].
- x Question: What does the customer think of this restaurant? Answer: It is [MASK]
- x Question: What does the customer think of this place? Answer: It is [MASK]

**Verbalizer**  $1 \mapsto \text{terrible}, 2 \mapsto \text{bad}, 3 \mapsto \text{okay}, 4 \mapsto \text{good}, 5 \mapsto \text{great}$ 

#### A.3 Yahoo Questions

## **Inputs**

- $x_1$ : The question to be classified.
- $x_2$ : The top answer to the question.

#### **NULL Patterns**

- [MASK]  $x_1 x_2$
- $x_1 x_2$  [MASK]
- $x_1$  [MASK]  $x_2$

#### **PUNC Patterns**

- [MASK] :  $x_1 x_2$
- [MASK] x<sub>1</sub> x<sub>2</sub>
- [MASK] .  $x_1 x_2$
- ( [MASK] )  $x_1 x_2$
- $x_1$  ( [MASK] )  $x_2$
- $x_1 x_2$  ( [MASK])

#### **PROMPT Patterns**

- $x_1 x_2$  Category: [MASK].
- *x*<sub>1</sub> *x*<sub>2</sub> Class: [MASK].
- *x*<sub>1</sub> *x*<sub>2</sub> Topic: [MASK].
- $x_1 x_2$  Theme: [MASK].
- x<sub>1</sub> x<sub>2</sub> Category: [MASK]
- *x*<sub>1</sub> *x*<sub>2</sub> Class: [MASK]
- x<sub>1</sub> x<sub>2</sub> Topic: [MASK]
- $x_1 x_2$  Theme: [MASK]
- [MASK] Question:  $x_1 x_2$
- [MASK] QUESTION:  $x_1 x_2$

## **Q&A Patterns**

- x<sub>1</sub> x<sub>2</sub> Question: What is the topic of this question? Answer: [MASK].
- x<sub>1</sub> x<sub>2</sub> Question: What is the category of this question? Answer: [MASK].
- $x_1 x_2$  Question: What is the topic of this question? Answer: [MASK]
- x<sub>1</sub> x<sub>2</sub> Question: What is the category of this question? Answer: [MASK]

**Verbalizer** Society & Culture  $\mapsto$  Society, Science & Mathematics  $\mapsto$  Science, Health  $\mapsto$  Health, Education & Reference  $\mapsto$  Education, Computers & Internet  $\mapsto$  Computer, Sports  $\mapsto$  Sports, Business & Finance  $\mapsto$  Business, Entertainment & Music  $\mapsto$  Entertainment, Family & Relationships  $\mapsto$  Relationship, Politics & Government  $\mapsto$  Politics

## B PVPs for RAFT

Below, we list the PVPs used for all tasks in the RAFT benchmark. For each task t, we make use of the original task description  $D_t$  that we copy verbatim from Alex et al. (2021). We use two vertical bars (II) to mark boundaries between text segments. If multiple patterns and verbalizers are specified, each verbalizer is used in combination with each pattern.

#### B.1 ADE

**Task Description** Label the sentence based on whether it is related to an adverse drug effect (ADE). Details are described below:

Drugs: Names of drugs and chemicals that include brand names, trivial names, abbreviations and systematic names were annotated. Mentions of drugs or chemicals should strictly be in a therapeutic context. This category does not include the names of metabolites, reaction byproducts, or hospital chemicals (e.g. surgical equipment disinfectants).

Adverse effect: Mentions of adverse effects include signs, symptoms, diseases, disorders, acquired abnormalities, deficiencies, organ damage or death that strictly occur as a consequence of drug intake.

#### **Inputs**

• x: The text to be classified.

#### **Patterns**

- D<sub>t</sub> || x Question: Is this sentence related to an adverse drug effect (ADE)? Answer: [MASK].
- x Question: Is this sentence related to an adverse drug effect (ADE)? Answer: [MASK].

<sup>10</sup>See https://github.com/oughtinc/raftbaselines/tree/master/example\_prompts.

#### **Verbalizers**

- not ADE-related  $\mapsto$  No, ADE-related  $\mapsto$  Yes
- not ADE-related  $\mapsto$  False, ADE-related  $\mapsto$  True

#### B.2 B77

**Task Description** The following is a banking customer service query. Classify the query into one of the 77 categories available.

## **Inputs**

- x: The text to be classified.
- y: The correct intent for the given text.

#### **Patterns**

- D<sub>t</sub> || x Question: Is y the correct category for this query? Answer: [MASK].
- x Question: Is y the correct category for this query? Answer: [MASK].

#### Verbalizers

- False  $\mapsto$  No, True  $\mapsto$  Yes
- False  $\mapsto$  False, True  $\mapsto$  True

#### B.3 NIS

Task Description Label the impact statement based on whether it mentions a harmful application of the research done in the paper. Make sure the statement is sufficient to conclude there are harmful applications of the research being done, not a past risk that this research is solving.

## **Inputs**

• x: The text to be classified.

## **Patterns**

- D<sub>t</sub> || x Question: Does this impact statement mention a harmful application? Answer: [MASK].
- x Question: Does this impact statement mention a harmful application? Answer: [MASK].

#### **Verbalizers**

- doesn't mention a harmful application 

  No, mentions a harmful application 

  Yes
- doesn't mention a harmful application  $\mapsto$  False, mentions a harmful application  $\mapsto$  True

#### B.4 OSE

Task Description The following is an article sourced from The Guardian newspaper, and rewritten by teachers to suit three levels of adult English as Second Language (ESL) learners: elementary, intermediate, and advanced. Predict the level of the article.

## **Inputs**

• x: The text to be classified.

#### **Patterns**

- D<sub>t</sub> || x Question: What is the level of this article? Answer: [MASK].
- x Question: What is the level of this article? Answer: [MASK].
- D<sub>t</sub> || x Question: Is the level of this article "elementary", "intermediate" or "advanced"? Answer: [MASK].
- x Question: Is the level of this article "elementary", "intermediate" or "advanced"?
   Answer: [MASK].

## Verbalizers

elementary → elementary, intermediate
 → intermediate, advanced → advanced

#### B.5 Over

Task Description In law, an overruling sentence is a statement that nullifies a previous case decision as a precedent, by a constitutionally valid statute or a decision by the same or higher ranking court which establishes a different rule on the point of law involved. Label the sentence based on whether it is overruling or not.

## **Inputs**

• x: The text to be classified.

#### **Patterns**

- D<sub>t</sub> || x Question: Is this sentence overruling? Answer: [MASK].
- x Question: Is this sentence overruling? Answer: [MASK].

#### Verbalizers

- not overruling  $\mapsto$  No, overruling  $\mapsto$  Yes
- not overruling  $\mapsto$  False, overruling  $\mapsto$  True

## B.6 SOT

Task Description The dataset is a list of institutions that have contributed papers to semiconductor conferences in the last 25 years, as catalogued by IEEE and sampled randomly. The goal is to classify the institutions into one of three categories: "university", "company" or "research institute".

### **Inputs**

- $x_1$ : The title of the paper to be classified.
- $x_2$ : The name of the organization to be classified.

### **Patterns**

- D<sub>t</sub> ∥ Organization name: x<sub>1</sub> Paper title: x<sub>2</sub> Question: What is the category of this institution? Answer: [MASK].
- Organization name:  $x_1$  Paper title:  $x_2$  Question: What is the category of this institution? Answer: [MASK].
- D<sub>t</sub> ∥ Paper title: x<sub>2</sub> Organization name: x<sub>1</sub> Question: What is the category of this institution? Answer: [MASK].
- Paper title: x<sub>2</sub> Organization name: x<sub>1</sub>
   Question: What is the category of this institution? Answer: [MASK].

#### Verbalizers

• company  $\mapsto$  company, research institute  $\mapsto$  institute, university  $\mapsto$  university

#### B.7 SRI

Task Description Identify whether this paper should be included in a meta-review which includes the findings of systematic reviews on interventions designed to promote charitable donations.

Included reviews should describe monetary charitable donations, assess any population of participants in any context, and be peer reviewed and written in English.

They should not report new data, be nonsystematic reviews, consider cause-related marketing or other kinds of prosocial behaviour.

## **Inputs**

- $x_1$ : The title of the paper to be classified.
- $x_2$ : The abstract of the paper to be classified.
- $x_3$ : The journal of the paper to be classified.

#### **Patterns**

- $D_t$  || Title:  $x_1$  Abstract:  $x_2$  Journal:  $x_3$  Question: Should this paper be included in a meta-review which includes the findings of systematic reviews on interventions designed to promote charitable donations? Answer: [MASK].
- Title: x<sub>1</sub> Abstract: x<sub>2</sub> Journal: x<sub>3</sub> Question: Should this paper be included in a meta-review which includes the findings of systematic reviews on interventions designed to promote charitable donations? Answer: [MASK].

### Verbalizers

- not included  $\mapsto$  No, included  $\mapsto$  Yes
- not included  $\mapsto$  False, included  $\mapsto$  True

## B.8 TAI

Task Description Transformative AI (TAI) is defined as AI that precipitates a transition comparable to (or more significant than) the agricultural or industrial revolution. Label a paper as "TAI safety research" if:

- 1. The contents of the paper are directly motivated by, and substantively inform, the challenge of ensuring good outcomes for TAI,
- 2. There is substantive content on AI safety,

not just AI capabilities,

- 3. The intended audience is the community of researchers.
- 4. It meets a subjective threshold of seriousness/quality,
- 5. Peer review is not required.

## **Inputs**

- $x_1$ : The title of the paper to be classified.
- $x_2$ : The abstract of the paper to be classified.

#### **Patterns**

- D<sub>t</sub> || Title: x<sub>1</sub> Abstract: x<sub>2</sub> Question: Is this paper a TAI safety research paper? Answer: [MASK].
- Title: x<sub>1</sub> Abstract: x<sub>2</sub> Question: Is this paper a TAI safety research paper? Answer: [MASK].

#### **Verbalizers**

- not TAI safety research → No, TAI safety research → Yes
- not TAI safety research  $\mapsto$  False, TAI safety research  $\mapsto$  True

## B.9 ToS

**Task Description** Label the sentence from a Terms of Service based on whether it is potentially unfair. If it seems clearly unfair, mark it as potentially unfair.

According to art. 3 of the Directive 93/13 on Unfair Terms in Consumer Contracts, a contractual term is unfair if: 1) it has not been individually negotiated; and 2) contrary to the requirement of good faith, it causes a significant imbalance in the parties rights and obligations, to the detriment of the consumer.

## **Inputs**

• x: The text to be classified.

## Patterns

- D<sub>t</sub> || x Question: Is this sentence potentially unfair? Answer: [MASK].
- x Question: Is this sentence potentially unfair? Answer: [MASK].

#### Verbalizers

- not potentially unfair  $\mapsto$  No, potentially unfair  $\mapsto$  Yes
- not potentially unfair  $\mapsto$  False, potentially unfair  $\mapsto$  True

## **B.10 TEH**

Task Description Label whether the following tweet contains hate speech against either immigrants or women. Hate Speech (HS) is commonly defined as any communication that disparages a person or a group on the basis of some characteristic such as race, color, ethnicity, gender, sexual orientation, nationality, religion, or other characteristics.

## **Inputs**

• x: The text to be classified.

#### **Patterns**

- D<sub>t</sub> || x Question: Does this tweet contain hate speech against either immigrants or women? Answer: [MASK].
- x Question: Does this tweet contain hate speech against either immigrants or women? Answer: [MASK].

## Verbalizers

- not hate speech  $\mapsto$  No, hate speech  $\mapsto$  Yes
- not hate speech  $\mapsto$  False, hate speech  $\mapsto$  True

#### **B.11** TC

**Task Description** A complaint presents a state of affairs which breaches the writer's favorable expectation. Label the tweet text based on whether it contains a complaint.

# **Inputs**

• x: The text to be classified.

## Patterns

- D<sub>t</sub> || x Question: Does this tweet text contain a complaint? Answer: [MASK].
- x Question: Does this tweet text contain a complaint? Answer: [MASK].

## **Verbalizers**

- no complaint  $\mapsto$  No, complaint  $\mapsto$  Yes
- no complaint  $\mapsto$  False, complaint  $\mapsto$  True