

Learning to Initialize: Can Meta Learning Improve Cross-task Generalization in Prompt Tuning?

Anonymous ACL submission

Abstract

Prompt tuning (PT) which only tunes the embeddings of an additional sequence of tokens per task, keeping the pre-trained language model (PLM) frozen, has shown remarkable performance in few-shot learning. Despite this, PT has been shown to rely heavily on good initialization of the prompt embeddings. In this work, we study *meta prompt tuning* (MPT) to systematically explore how meta-learning can help improve (if it can) cross-task generalization in PT through learning to initialize the prompt embeddings from other relevant tasks. We empirically analyze a representative set of meta learning algorithms in a wide range of adaptation settings with different source/target task configurations on a large set of few-shot tasks. With extensive experiments and analysis, we demonstrate the effectiveness of MPT. We find the improvement to be significant particularly on classification tasks. For other kinds of tasks such as question answering, we observe that while MPT can outperform PT in most cases, it does not always outperform multi-task learning. We further provide an in-depth analysis from the perspective of task similarity.

1 Introduction

Humans can easily learn to perform new tasks with only few data by leveraging previously acquired knowledge from other relevant tasks. Such capability is a hallmark of human intelligence (Carey and Bartlett, 1978). However, when it comes to the models, they often face over-fitting issues when they are tasked to learn from a few labeled examples (Lake et al., 2017; Linzen, 2020), a problem commonly termed as *few-shot learning* (FSL).

With the recent advancements in developing large-scale pre-trained language models (PLMs), prompt-based methods have shown promising results in FSL. Brown et al. (2020) show that by virtue of in-context (meta) learning, a frozen GPT-3 model can achieve good results on a variety of

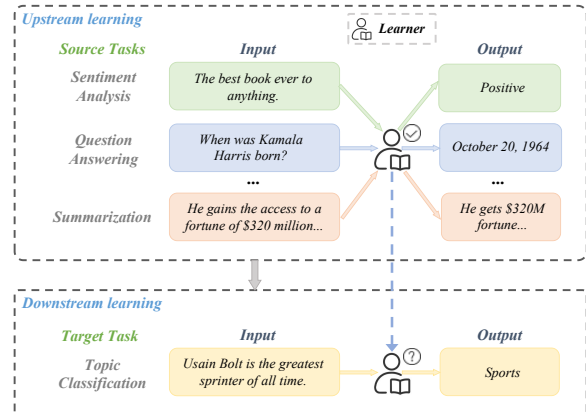


Figure 1: Illustration of cross-task generalization, where the model is expected to learn an unseen *target* task given the knowledge acquired from previously learned *source* tasks.

few-shot tasks through manually designed *prompts*, which are task instructions along with a few examples expressed in natural language. However, the performance of in-context learning has been shown to be highly sensitive to the design of such “discrete” prompts (Zhao et al., 2021). It is also limited by the maximum sequence length supported by the PLMs (Li and Liang, 2021). Down this line, efforts have been made on automatically searching and optimizing for discrete prompts (Shin et al., 2020; Schick and Schütze, 2021; Gao et al., 2021).

As an alternative to discrete prompts, recent efforts attempt to learn “soft” prompts that add additional trainable parameters (Liu et al., 2021b; Li and Liang, 2021; Lester et al., 2021), showing better results than discrete prompts (Liu et al., 2021a). Lester et al. (2021) introduce *prompt tuning* (PT) that prepends a sequence of *tunable* tokens to the input and optimize their embeddings keeping the PLM frozen. Despite its strong few-shot performance, PT has been shown to be sensitive to the initialization of the embeddings, which might limit its practical application (Qin and Joty, 2022b). To address this, Gu et al. (2022) propose *pre-trained*

066 *prompt tuning* (PPT) to pre-train soft prompts using
067 self-supervised tasks on unlabeled data. It relies
068 on carefully designed pre-training tasks tailored to
069 the downstream tasks, and the pre-training objec-
070 tives are only applicable to classification tasks. [Vu](#)
071 [et al. \(2022\)](#) introduce *soft prompt transfer* (SPoT),
072 which uses the soft prompts learned from a set
073 of source tasks through multi-task learning to ini-
074 tialize the prompt for a target task. Both PPT
075 and SPoT demonstrate *cross-task generalization*
076 (Fig. 1) – learning of a new task can benefit from
077 learning of other related tasks ([Ye et al., 2021](#)).

078 In a recent survey, [Lee et al. \(2022\)](#) claim that
079 *meta learning* ([Schmidhuber, 1987](#)) can play an im-
080 portant role for cross-task generalization in NLP.¹
081 Different from multi-task learning which consid-
082 ers the performance on the source tasks to learn
083 the initial parameters, meta learning aims to find
084 initial parameters suitable for adapting to a target
085 few-shot task. Hence, it could outperform multi-
086 task learning in several scenarios with *full-model*
087 finetuning ([Dou et al., 2019](#); [Chen et al., 2020b](#)).
088 However, to our knowledge, there is no systematic
089 study on the role of meta learning on PT. In a recent
090 work, [Huang et al. \(2022\)](#) adopt MAML ([Finn et al.,](#)
091 [2017](#)) for pre-training soft prompts. One major lim-
092 itation of their study is that it is limited to only one
093 type of meta learning algorithm and only sentiment
094 classification tasks, lacking comprehensive under-
095 standing of cross-task generalization. [Min et al.](#)
096 [\(2022\)](#) and [Chen et al. \(2022\)](#) show the effective-
097 ness of in-context learning for PLMs, whereas we
098 mainly focus on optimization-based meta learning.

099 To systematically study meta prompt tuning
100 (MPT) for cross-task generalization, we conduct
101 experiments on a large collection of few-shot tasks
102 involving different types of datasets with a unified
103 text-to-text format ([Ye et al., 2021](#)). We investigate
104 a wide range of adaptation settings with different
105 source/target task types, which helps better under-
106 stand the capability and limitation of meta learning
107 in PT. With extensive experiments, we aim to ad-
108 dress the following research questions:

- 109 • **Q1.** Can MPT improve cross-task generalization
110 in PT? Is it better than multi-task learning?
- 111 • **Q2.** What happens with more labelled data for
112 source/target tasks (beyond few-shot settings)?

¹Unless otherwise specified, by meta learning in this paper we generally refer to the optimization-based meta learning algorithms, and use more specific names for the other kinds such as *in-context learning* for black-box meta learning and *metric learning* for non-parametric meta learning.

- **Q3.** Does it help with more diverse source tasks? 113
- **Q4.** Is the performance gain of MPT consistent 114
across different backbone models? 115

116 To answer these questions, we empirically an-
117alyze MAML ([Finn et al., 2017](#)), FoMAML and
118 Reptile ([Nichol et al., 2018](#)), which constitute a
119 representative set of meta learning methods. Ex-
120 perimental results show that MPT can indeed help
121 cross-task generalization, *e.g.*, MAML improves
122 the performance of PT by more than 20% on clas-
123 sification tasks. However, we also notice that MPT
124 does not always outperform multi-task learning, es-
125 pecially on non-classification tasks. We provide an
126 in-depth analysis from the perspective of task sim-
127 ilarity. As for Q2, we find that MPT does benefit
128 cross-task generalization beyond few-shot settings.
129 For Q3, we observe that increasing the diversity
130 of source tasks does not necessarily improve cross-
131 task generalization. Finally, the consistent gain of
132 MPT across different models shows its robustness
133 to model type and size. In summary, the two main
134 contributions of this work are:

- To the best of our knowledge, we are the first 135
to extensively explore how meta learning helps 136
cross-task generalization in prompt tuning. 137
- With extensive experiments and analysis, we 138
show the effectiveness and limitation of meta 139
prompt tuning in various source/target settings. 140
Our code base is available at <redacted>. 141

142 2 Related Work

143 **Few-shot Learning (FSL)** FSL aims to learn a
144 task with only a few labeled examples, which often
145 leads to the over-fitting problem. Existing methods
146 to address this problem mainly focus on optimizing
147 the hypothesis space of the few-shot tasks ([Tri-](#)
148 [antafillou et al., 2017](#); [Finn et al., 2017](#); [Hu et al.,](#)
149 [2018](#)) or augmenting the few-shot data ([Gao et al.,](#)
150 [2020](#); [Qin and Joty, 2022a](#)). Recently, large-scale
151 pre-trained language models (PLMs) have demon-
152 strated strong FSL ability through prompt-based
153 methods, including both discrete ([Brown et al.,](#)
154 [2020](#)) and soft prompts ([Lester et al., 2021](#)).

155 **Prompt-based Learning (PL)** PL is a new
156 paradigm which prepends a task-specific template
157 or prompt to the input for learning new tasks ([Liu](#)
158 [et al., 2021a](#)). Initial PL methods mainly focus
159 on designing, searching or optimizing discrete
160 prompts ([Brown et al., 2020](#); [Shin et al., 2020](#);
161 [Gao et al., 2021](#)). However, discrete prompts are
162 hard to optimize. To solve this, recent PL methods

attempt to optimize prompts in a continuous space, *i.e.*, learn soft prompts (Li and Liang, 2021; Liu et al., 2021b; Lester et al., 2021), showing impressive FSL performance (Qin and Joty, 2022b). In addition to prompt design, several recent studies have explored the applications (Zhu et al., 2022; Li et al., 2022) and analysis (Zhong et al., 2021; Le Scao and Rush, 2021) of PL.

Meta Learning Meta Learning or learning to learn, has been applied to boost few-shot performance on various NLP tasks, *e.g.*, relation extraction (Han et al., 2018) and machine translation (Gu et al., 2018). Meta learning algorithms can be divided into three main categories. First, *black-box* methods adopt additional meta learners to help adaptation (Santoro et al., 2016; Garnelo et al., 2018; Mishra et al., 2018; Brown et al., 2020). Second, *non-parametric* methods explore how to learn metrics that can compare the distances between different samples, *i.e.*, learning to compare (Koch et al., 2015; Vinyals et al., 2016; Snell et al., 2017). Finally, *optimization-based* methods aim to learn better parameter initialization to effectively and efficiently adapt to unseen tasks, *i.e.*, learning to initialize (Finn et al., 2017; Nichol et al., 2018; Kedia et al., 2021). Lee et al. (2022) claim that meta learning can be effective for cross-task generalization, especially the optimization-based methods. They can be applied to various problems in a model-agnostic way to improve FSL on target tasks with model fine-tuning (Ye et al., 2021).

Summary. Existing work shows that meta learning can improve cross-task few-shot generalization with full model fine-tuning. However, there is no systematic study on whether (and how) meta learning can do so with prompt tuning of PLMs. To fill this research gap, our work provides a comprehensive understanding of the effectiveness and limitation of meta learning in prompt tuning.

3 Preliminaries

In this section, we revisit the basics about prompt tuning and optimization-based meta learning.

3.1 Prompt Tuning

Following Lester et al. (2021), we reframe all tasks into a text-to-text format. Given a training dataset $\mathcal{D}^{tr} = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ for a task \mathcal{T} , different from traditional model fine-tuning, prompt tuning (PT) is a parameter-efficient learning method which freezes the PLM θ and prepends

the input text X_i with a sequence of *tunable* soft tokens P , parameterized by prompt embeddings ϕ . The prompt embeddings ϕ are initialized from the vocabulary of the PLM and optimized through gradient descent with the following objective:

$$\mathcal{L}_\phi^T = \mathcal{L}(\phi, \mathcal{D}^{tr}) = - \sum_{i=1}^n \log p(Y_i | [P, X_i], \phi, \theta) \quad (1)$$

3.2 Optimization-based Meta Learning

The main goal of optimization-based meta learning (or learning to initialize), is to learn better initial parameters that can effectively and efficiently adapt to a new task \mathcal{T}^{new} with limited data. We denote the initial parameters (meta-parameters) as ϕ^* .

To obtain ϕ^* , the model needs to learn from a series of *meta-training* tasks $\mathcal{T}^{\text{meta}} = \{\mathcal{T}_1, \dots, \mathcal{T}_n\}$. The dataset \mathcal{D}_i of each task \mathcal{T}_i is divided into two disjoint sets: a *support set* \mathcal{S}_i and a *query set* \mathcal{Q}_i . The objective for learning ϕ^* is

$$\phi^* = \arg \min_{\phi} \sum_{\mathcal{T}_i \in \mathcal{T}^{\text{meta}}} \underbrace{\mathcal{L}(\phi - \alpha \nabla_{\phi} \mathcal{L}(\phi, \mathcal{S}_i), \mathcal{Q}_i)}_{\text{inner update}} \quad (2)$$

where \mathcal{L} is the objective function defined in Eq. (1), ϕ is the set of parameters to meta-learn and α is the inner learning rate. Denoting the overall loss as $\mathcal{L}_\phi^{\mathcal{T}^{\text{meta}}} = \sum_{\mathcal{T}_i \in \mathcal{T}^{\text{meta}}} \mathcal{L}(\phi', \mathcal{Q}_i)$ with ϕ' being the inner-updated value of ϕ , we use gradient descent to update ϕ further in the meta-training stage:

$$\phi = \phi - \beta \nabla_{\phi} \mathcal{L}_\phi^{\mathcal{T}^{\text{meta}}} \quad (3)$$

where β is the outer learning rate. This is actually the Model-Agnostic Meta-Learning or MAML (Finn et al., 2017). Notice that optimizing Eq. (3) requires calculating second-order gradients, which can be quite memory-consuming. To alleviate this, First-order MAML (FoMAML) and Reptile (Nichol et al., 2018) are proposed to use first-order approximations, allowing lower memory costs.

After the meta-training stage, ϕ^* serves as the initial parameters for learning an unseen *meta-testing* task \mathcal{T}^{new} which is usually few-shot.

4 Approach

In this section, we first introduce the problem setting and evaluation metric. Then, we illustrate the key methods for meta prompt tuning (MPT).

4.1 Problem Setting

To evaluate cross-task generalization in prompt tuning, we select a large and diverse collection of few-shot tasks from Ye et al. (2021), covering various

types including classification, question answering and generation. We partition the set of all tasks \mathcal{T}^{all} into two disjoint parts: source tasks \mathcal{T}^{src} and target tasks \mathcal{T}^{tgt} . Details of the tasks and partitions are provided later in our experiment setup (§5).

Following Min et al. (2022), we can divide the whole learning process into two stages (Fig. 1):

- Upstream learning on source tasks** In this stage, the model has access to \mathcal{T}^{src} , which is regarded as *meta-training* tasks $\mathcal{T}^{\text{meta}}$ in Eq. (2). We divide the dataset \mathcal{D}_i of every source task \mathcal{T}_i into training (or support) and validation (or query) sets, and conduct optimization-based meta learning or multi-task learning on these sets to obtain meta-parameters ϕ^* . Note that we use both support and query sets for model training in multi-task learning to ensure fair data access for both methods.

- Downstream learning on target tasks** After the upstream learning stage, we use the learned meta-parameters ϕ^* as the initial point for learning target tasks \mathcal{T}^{tgt} . Every target task \mathcal{T}_k has its own training set $\mathcal{D}_k^{\text{tr}}$, validation set $\mathcal{D}_k^{\text{val}}$, and test set $\mathcal{D}_k^{\text{test}}$. The model is required to learn from $\mathcal{D}_k^{\text{tr}}$ via prompt tuning and will be evaluated on $\mathcal{D}_k^{\text{test}}$. The performance on $\mathcal{D}_k^{\text{val}}$ is used for hyper-parameters tuning and model selection.

This two-stage learning paradigm can naturally reflect cross-task generalization where the model needs to learn an unseen task given previously acquired knowledge from other tasks.

4.2 Evaluation Metric

We evaluate the model performance on a set of target tasks \mathcal{T}^{tgt} . As \mathcal{T}^{tgt} may cover various task types, simply averaging the performance of different target tasks is unreasonable. Following Ye et al. (2021), we use *average relative gain* (ARG) as the main evaluation metric. We first calculate *relative gain* (RG) for each target task, *i.e.*, relative performance improvement before and after applying the upstream (meta or multi-task) learning on the source tasks. Then we average the relative gains of all target tasks to obtain the final result which indicates the overall performance improvement.

4.3 Meta Prompt Tuning (MPT)

As shown in Fig. 2, the key idea of MPT is to apply optimization-based meta-training as upstream learning to a set of source tasks in order to learn meta parameters, which in this case are prompt embeddings. The learned prompt embeddings serve

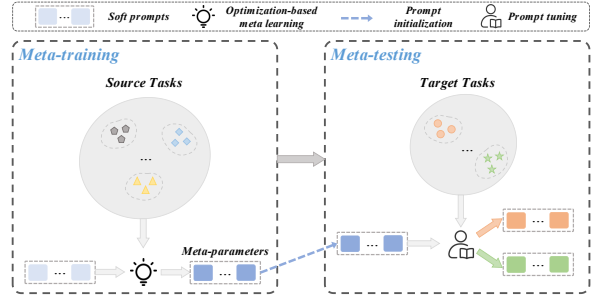


Figure 2: Overview of Meta Prompt Tuning (MPT). In the meta-training stage, we conduct optimization-based meta learning on source tasks to obtain meta-parameters (*i.e.*, soft prompts). The meta-parameters will then be used to initialize prompt embeddings for learning unseen target tasks in the meta-testing stage.

as the initialization for learning unseen target tasks, referred to as meta-testing or downstream learning.

4.3.1 Meta-training

We meta-train the prompt embeddings on source tasks \mathcal{T}^{src} . Without loss of generality, we take MAML (Finn et al., 2017) as an example. For every iteration, we first sample one source task \mathcal{T}_i which has a support set \mathcal{S}_i and a query set \mathcal{Q}_i . Then we sample a support batch \mathcal{B}_s from \mathcal{S}_i and a query batch \mathcal{B}_q from \mathcal{Q}_i . Denoting the trainable prompt embeddings as ϕ , \mathcal{B}_s and \mathcal{B}_q are used for one gradient update with the following objective:

$$\begin{aligned} \mathcal{L}_\phi^i &= \mathcal{L}(\phi - \alpha \nabla_\phi \mathcal{L}(\phi, \mathcal{B}_s), \mathcal{B}_q) \\ \phi &= \phi - \beta \nabla_\phi \mathcal{L}_\phi^i \end{aligned} \quad (4)$$

where \mathcal{L} is the task loss defined in Eq. (1), and α and β are inner and outer learning rates, respectively. During the meta-training stage, we iterate over tasks in \mathcal{T}^{src} to update prompt embeddings ϕ for a fixed number of steps. The learned meta-parameters ϕ^* is used in the meta-testing stage.

4.3.2 Meta-testing

In meta-testing, the model is expected to learn unseen target tasks \mathcal{T}^{tgt} . For each target task \mathcal{T}_k , we use the learned meta-parameters ϕ^* to initialize the prompt embeddings for the task. Denoting the training set of \mathcal{T}_k as $\mathcal{D}_k^{\text{tr}}$, the learning objective during meta testing is defined as:

$$\mathcal{L}_{\phi^*}(\mathcal{D}_k^{\text{tr}}) = - \sum_{i=1}^n \log p(Y_i | [P^*, X_i], \phi^*, \theta) \quad (5)$$

where θ is the frozen PLM, $(X_i, Y_i) \sim \mathcal{D}_k^{\text{tr}}$ is a training sample and P^* are the prompt tokens.

Source		Target	
Setting	#tasks	Setting	#tasks
Random	114	Random	20
Classification (Cls)	45	Classification	10
Both (Cls + Non-Cls)	23 + 22		
Non-Classification	45		
Classification	45	Non-Classification	12
Both (Cls + Non-Cls)	23 + 22		
Non-Classification	45		
QA	22	QA	15
Non-QA	33		
Non-Paraphrase Cls	60	Paraphrase	4

Table 1: Statistics of ten distinct source/target task partitions. Appendix A.1 for details about each partition.

We evaluate the model with the best validation performance on the test set and calculate average relative gain on the test sets of \mathcal{T}^{tgt} .

5 Experimental Setup

We first describe the source/target task partitions, and then introduce methods compared in our work. Finally, we present the implementation details.

5.1 Task Partitions

We experiment with ten different source/target task partitions as shown in Table 1. Depending on the type of the target tasks, we can divide these ten settings into several groups:

- **R→R (Random→Random)**: We first experiment with the R→R setting where both source and target tasks are randomly selected, meaning that they can cover any task type. This setting mimics the learning paradigm of humans and reflects whether cross-task generalization can help obtain a general-purpose few-shot learner.
- **X→Cls (X=Cls, Both, Non-Cls)**: The target tasks involve classification, while the source tasks can be classification, non-classification tasks or both. This setting helps us better understand the influence of the source task distribution.
- **X→Non-Cls (X=Cls, Both, Non-Cls)**: The only difference between this and the previous setting is the type of target tasks. We investigate how meta learning improves cross-task generalization when target tasks are non-classification tasks.
- **X→QA (X=QA, Non-QA)**: Compared to the previous one, this group is more fine-grained. We only select target tasks from question answering (QA) instead of all non-classification tasks. We conduct experiment on different source task types, including QA and Non-QA tasks.

- **NP→P (Non-Paraphrase Cls→Paraphrase)**: This group has the finest granularity in our setting. We choose paraphrase identification which is a sub-category of classification as the target, and non-paraphrase classification as the source. The final two groups help understand how meta learning performs in more fine-grained scenarios.

Note that we ensure that there is no overlap between the source and target tasks. Following Ye et al. (2021), we use 16 samples per class in the training (or support) and validation (or query) sets for classification tasks, and 32 samples per set for non-classification tasks. For every task, we sample the training and validation sets 5 times with different random seeds to reduce variance in few-shot evaluation and cover more diverse samples in upstream learning. We provide full details of tasks and partitions in Appendix A.1.

5.2 Methods Compared

We mainly use T5-Large (Raffel et al., 2019) as the backbone language model and compare the following methods in our work.

- **Prompt Tuning (PT) on target tasks**. It is our baseline without the upstream learning. We directly apply PT (Lester et al., 2021) to target tasks and use its performance as the basis for computing average relative gain for other methods.
- **Model-Agnostic Meta-Learning (MAML)**. We apply MAML (Finn et al., 2017) in the upstream learning (meta-training) stage. The learned meta-parameters are used to initialize prompt embeddings for learning target tasks.
- **First-order MAML (FoMAML) and Reptile**. We also investigate two first-order meta learning algorithms: FoMAML (Finn et al., 2017) and Reptile (Nichol et al., 2018). Compared to MAML, they are more memory-efficient.
- **Multi-task learning (MTL)**. We conduct multi-task learning on source tasks instead of meta learning to obtain initial parameters. This is a straight-forward yet effective method as demonstrated by Vu et al. (2022).
- **Fine-tuning on target tasks**. Fine-tuning is the dominant paradigm where the whole language model is tuned for learning target tasks. We include it to verify whether cross-task generalization can help PT outperform fine-tuning.

In addition, we conduct experiments with differ-

ent backbone models to verify MPT’s robustness.


5.3 Implementation Details

All our methods are implemented with PyTorch/Transformers library (Wolf et al., 2020). We use higher library (Grefenstette et al., 2019) for higher-order optimization in meta learning methods. The prompt length in PT is set to 100 tokens following Lester et al. (2021). We provide details of other hyperparameters in Appendix A.4.

Since it is infeasible to search for optimal hyperparameters for each of the meta- and multi-task learning methods in each of the settings, we select them based on the R→R setting. We randomly select 5 tasks that are not in the source and target sets as validation tasks for hyperparameter search. The hyperparameters with best validation performance (ARG) are used for upstream learning. We select the inner learning rate, the outer learning rate and total training steps for MAML and adopt the same three hyperparameters for FoMAML and Reptile.

6 Results and Analysis

We now address the four research questions asked before in §1 with empirical results.

 **Q1.** Can meta prompt tuning improve cross-task generalization? Is it better than multi-task learning?

The ARG of different methods *w.r.t.* PT in various settings are shown in Table 2; more detailed results on every target task are in Appendix A.2.

• **MPT can indeed help cross-task generalization.** From the results in Table 2, we observe that MPT outperforms the baseline PT in most cases with +ve ARG scores. Out of 30 different runs for three meta learning methods in ten different settings (see the 1st block of results), MPT achieves better performance than PT in 23 runs, demonstrating its effectiveness in cross-task generalization.

For the R→R setting, MAML achieves the best performance, showing that it is a good general-purpose few-shot learner. For adapting to classification tasks, MAML outperforms PT by **20.16%** if the prompt embeddings are initialized from other classification tasks. The results in a more fine-grained setting (NP→P) also indicate the ability of MAML to learn classification tasks. While Reptile performs the best (20.44%) in this setting, MAML still outperforms PT by a large margin (**11.14%**).

However, as shown in Table 2, MAML falls behind FoMAML when adapting to non-classification

tasks. Among the three meta learning methods, FoMAML achieves the best performance (**9.81%**) on non-classification target tasks in the Both→Non-Cls setting, showing effective knowledge transfer. We observe similar results in more fine-grained settings QA/Non-QA→QA, where FoMAML outperforms MAML and Reptile significantly. While Reptile is claimed empirically to be better than MAML/FoMAML (Lee et al., 2022), it falls short of MAML/FoMAML in many cases. This might be because MAML and FoMAML are more similar compared to Reptile from a gradient perspective (Nichol et al., 2018). And since the hyperparameter search is done based on MAML (§5.3), which means Reptile’s method may be suboptimal.

In addition, we can see that meta learning helps PT outperform fine-tuning in several settings including Cls→Cls (MAML, FoMAML), Both→Cls (FoMAML) and NP→P (MAML, Reptile), which demonstrates the superiority of MPT.

• **MPT does not always outperform multi-task learning (MTL).** While meta learning is specifically designed for quickly adapting to unseen target tasks, it does not always outperform MTL in PT. From Table 2, we can observe that MTL achieves better performance than MPT in many cases, especially on non-classification target tasks. We analyze the reasons as follows:

- Meta learning methods have been shown to be highly sensitive to the hyperparameters (Antoniou et al., 2019), which we could not tune exhaustively due to memory/time constraints (see Appendix A.5 for hyperparameter sensitivity analysis). As mentioned in §5.3, we select the hyperparameters of MAML using the R→R setting, and then use the same hyperparameters for all meta learning methods in all settings, which might limit the performance of MPT.
- There might be less shared structure (or features) among non-classification tasks compared to classification. The classification tasks mostly involve sentence-level classification and in some cases the task labels correlate well (*e.g.*, AG News and DBpedia). Thus, they share some common semantics in both source and target tasks. The model can learn similar patterns (inferring the label of the entire input sentence) during both meta-training and meta-testing stages, enabling better knowledge transfer. The non-classification set on the other hand can include different types of tasks such as QA and summarization; modeling


Method	R→R	Cls →Cls	Both →Cls	Non-Cls →Cls	Cls →Non-Cls	Both →Non-Cls	Non-Cls →Non-Cls	QA →QA	Non-QA →QA	NP →P
MAML	8.78 \pm 0.69	20.16 \pm 0.84	10.57 \pm 1.03	6.34 \pm 0.48	0.32 \pm 0.04	7.54 \pm 0.73	6.71 \pm 0.39	-16.59 \pm 1.36	3.26 \pm 0.24	11.14 \pm 0.93
FoMAML	1.24 \pm 0.18	18.80 \pm 1.13	17.84 \pm 1.21	7.32 \pm 0.42	6.42 \pm 0.51	9.81 \pm 0.64	3.88 \pm 0.31	16.63 \pm 1.58	9.83 \pm 0.76	-0.68 \pm 0.07
Reptile	8.42 \pm 0.46	-5.17 \pm 0.71	-4.18 \pm 0.37	2.42 \pm 0.21	-1.54 \pm 0.18	-3.38 \pm 0.49	0.78 \pm 0.07	0.77 \pm 0.09	-0.09 \pm 0.01	20.44 \pm 1.34
Multi-task learning	7.14 \pm 0.62	-5.64 \pm 0.92	5.73 \pm 0.43	4.97 \pm 0.39	8.51 \pm 1.16	13.47 \pm 0.97	19.67 \pm 1.72	25.65 \pm 1.93	17.23 \pm 1.08	-5.19 \pm 0.86
Fine-tuning	-12.61 \pm 1.57	16.02 \pm 1.44	16.02 \pm 1.44	16.02 \pm 1.44	-35.70 \pm 2.73	-35.70 \pm 2.73	-35.70 \pm 2.73	-47.37 \pm 2.97	-47.37 \pm 2.97	1.56 \pm 0.12

Table 2: **Average relative gain (ARG %) of different methods with respect to prompt tuning (PT) in various settings.** Bold indicates the best ARG score. ‘Cls’, ‘QA’, ‘P’ and ‘NP’ respectively stand for ‘classification’, ‘question answering’, ‘paraphrase’ and ‘non-paraphrase classification’.

them typically requires a Seq2Seq formulation. These tasks typically lack shared task semantics. For example, the structure of QA is context + question + answer, requiring reasoning ability. In contrast, the structure of summarization is long document + short summary, requiring summarizing ability. Although it has been shown that QA can help summarization in content selection (Arumae and Liu, 2019), it is more difficult for MPT to capture transferable knowledge as success of meta learning eventually depends on how much the tasks share (Finn, 2022).

To provide an in-depth analysis of the difference between classification and non-classification tasks, we consider from the perspective of task similarity. Following (Lin et al., 2022), the correlation between input subspaces (the norm of projected subspace onto the other subspace) for two tasks could serve as the similarity score between them. We randomly pick 5 (cls,cls) task pairs as similar tasks. For dissimilar tasks, we randomly pick 5 (QA, summarization) task pairs. The average similarity score for similar task pairs is 0.768 while the average similarity score for dissimilar task pairs is only 0.306 (see Appendix A.6 for detailed results), which verifies that classification tasks share more structure than non-classification tasks.

Given the performance gap between MPT and MTL in some settings, we believe that exploring more advanced MPT methods could be a promising research direction.

 **Q2.** What happens with more labelled data for source/target tasks (beyond few-shot settings)?

As mentioned in §5.1, we mainly explore how MPT improves cross-task generalization when both the source and target tasks are few-shot, which corresponds to the way humans learn (Lake et al., 2017). We used 16 samples per class for classification tasks, and 32 samples per dataset for non-classification tasks. To validate whether more la-

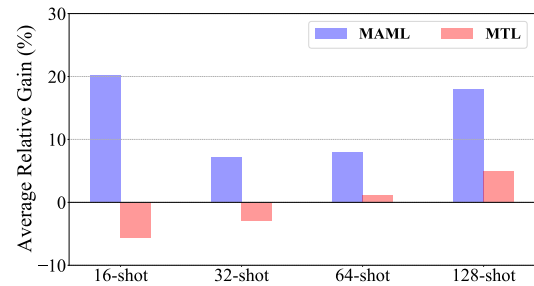



Figure 3: ARG (%) of MPT (MAML) and multi-task learning w.r.t. prompt tuning (ARG = 0) for varying data size of source tasks in the Cls→Cls setting.

belled data for source/target tasks can influence the performance of MPT, we conduct controlled experiments with {32, 64, 128} samples per class for source/target tasks in the Cls→Cls setting.

• **Source** We report the results of MAML and MTL with more labelled data for the source tasks in Fig. 3. We can observe that: (i) MPT outperforms PT (ARG = 0) and MTL in all cases, showing its robustness to data sizes. (ii) Increasing the number of samples in source tasks *does not* necessarily lead to better cross-task generalization for MPT. The best ARG is achieved for 16-shot, which justifies using few-shot source tasks. (iii) The performance of MTL improves with more data for source tasks, showing a different learning pattern from MPT.

• **Target** Table 3 shows the results for increasing the number of examples in target tasks. We can see that: (i) The performance gain of MPT is evident even at 128-shot (8.36%), demonstrating that it *does* help cross-task generalization beyond few-shot. (ii) MPT outperforms MTL by a large margin in all settings. (iii) MTL is unstable in terms of ARG scores; while it outperforms PT in 64-shot (1.96%), it falls behind PT in all other settings, indicating that MPT is a better choice when adapting to classification tasks.

 **Q3.** Does MPT help with more diverse source tasks?

Method	Shot			
	16	32	64	128
MPT (MAML)	20.16	9.10	5.64	8.36
Multi-task learning	-5.64	-14.17	1.96	-0.20

Table 3: ARG (%) of different methods when more labelled data is used in target tasks.

Method	Source task number		
	12	24	45
MPT (MAML)	8.44	12.89	20.16


Table 4: ARG (%) of MPT (MAML) when using different number of source tasks in the Cls→Cls setting.

MPT aims to learn to initialize the prompt embeddings from source tasks, which may cover different types. We hypothesize that the diversity of source tasks might influence its performance. To verify this, we analyze the influence of different source task selections on the same target tasks in two settings: varying the type and number of tasks.

• **Type of tasks.** The results of learning from different types of source tasks are reported in Table 2. The performance of MPT on non-classification target tasks improves when using more diverse source tasks, *e.g.*, from Non-Cls/Cls→Non-Cls to Both→Non-Cls. However, for adapting to classification task, the best ARG is achieved when all source tasks are classification, *i.e.*, the Cls→Cls setting. Hence, we can conclude that increasing the type diversity of source tasks *does not* necessarily improve cross-task generalization, which is consistent with the finding in (Ye et al., 2021).

• **Number of tasks.** To investigate the impact of the number of source tasks, we conduct controlled experiments on {12, 24} source tasks sampled from the original 45 source tasks in the Cls→Cls setting (see Appendix A.3 for a full list). From Table 4, we can observe that the performance of MPT keeps improving as the number of source tasks increases, showing better cross-task generalization.

It is worthwhile to note that while our work provides some insights on the choice of source tasks, more systematic studies on how to select the most suitable source tasks given a set of target tasks are needed. We hope that future analysis can provide a more comprehensive understanding of the relationship between source and target tasks.

 **Q4.** Is the performance gain of MPT consistent across different backbone language models?

Method	MAML	FoMAML	Reptile	MTL	Fine-tuning
T5-Large	11.14	-0.68	20.44	-5.19	1.56
T5-Base	9.24	4.15	7.96	1.64	7.41
T5-XLarge	14.35	2.46	10.74	5.72	-9.61
BART-Large	7.63	1.16	8.94	-2.37	2.74
GPT2-Large	3.19	-2.68	4.62	-1.43	3.75

Table 5: Average relative gain (ARG %) of all methods with different backbone models in the NP→P setting. ‘MTL’ stands for ‘multi-task learning’.

Our experiments and analysis so far use T5-Large as the backbone model. To verify whether the performance gain of MPT is consistent across different backbone models, we extend the experiments to T5-Base, T5-XLarge, BART-Large and GPT2-Large in the NP→P setting. From the results shown in Table 5, we can see that MPT still outperforms PT and MTL by a large margin when using other PLMs as the backbone model, showing its robustness to model size and type. In addition, the consistent gain of MPT with T5-XLarge could also verify the effectiveness of MPT for huge PLMs which have been shown to perform better in prompt tuning (Lester et al., 2021).

6.1 Further Analysis

Prompt tuning (PT) vs. Fine-tuning (FT). While PT shows strong few-shot learning ability, FT remains the dominant paradigm. As shown in Table 2, FT outperforms PT when adapting to classification tasks even in few-shot settings, which might be because PT has only a few tunable parameters. Though MPT is based on PT, its performance gain over FT in all cases suggests that it can learn to initialize the prompt embeddings from source tasks, enabling effective knowledge transfer.

7 Conclusion

In this paper, we have introduced meta prompt tuning (MPT), which learns to initialize the prompt embeddings for adapting to a target task. We have identified key research questions and systematically studied where and how meta learning can improve cross-task generalization in prompt tuning. We have empirically analyzed a representative set of meta learning methods in a variety of adaptation settings on a large, diverse collection of few-shot tasks. Extensive experimental results and analysis verify the effectiveness of MPT. Given the findings, in the future, we would like to explore more advanced meta learning algorithms which can consistently outperform multi-task learning.

652
653
654
655
656
657
658

659
660
661
662
663
664
665
666
667
668

669
670
671

672
673
674
675
676
677
678
679

680
681
682
683
684
685
686

687
688
689
690
691

692
693
694
695
696
697
698

699
700
701
702
703
704

705
706
707

References

Tiago A. Almeida, José María G. Hidalgo, and Akebo Yamakami. 2011. [Contributions to the study of sms spam filtering: New collection and results](#). In *Proceedings of the 11th ACM Symposium on Document Engineering*, DocEng '11, page 259–262, New York, NY, USA. Association for Computing Machinery.

Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. 2019. [MathQA: Towards interpretable math word problem solving with operation-based formalisms](#). 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 2357–2367, Minneapolis, Minnesota. Association for Computational Linguistics.

Antreas Antoniou, Harrison Edwards, and Amos Storkey. 2019. [How to train your MAML](#). In *International Conference on Learning Representations*.

Kristjan Arumae and Fei Liu. 2019. [Guiding extractive summarization with question-answering rewards](#). 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 2566–2577, Minneapolis, Minnesota. Association for Computational Linguistics.

Francesco Barbieri, Jose Camacho-Collados, Luis Espinosa Anke, and Leonardo Neves. 2020. [TweetEval: Unified benchmark and comparative evaluation for tweet classification](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1644–1650, Online. Association for Computational Linguistics.

Max Bartolo, Alastair Roberts, Johannes Welbl, Sebastian Riedel, and Pontus Stenetorp. 2020. [Beat the AI: Investigating adversarial human annotation for reading comprehension](#). *Transactions of the Association for Computational Linguistics*, 8:662–678.

Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. 2013. [Semantic parsing on Freebase from question-answer pairs](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1533–1544, Seattle, Washington, USA. Association for Computational Linguistics.

Chandra Bhagavatula, Ronan Le Bras, Chaitanya Malaviya, Keisuke Sakaguchi, Ari Holtzman, Hannah Rashkin, Doug Downey, Wen tau Yih, and Yejin Choi. 2020. [Abductive commonsense reasoning](#). In *International Conference on Learning Representations*.

Yonatan Bisk, Rowan Zellers, Ronan LeBras, Jianfeng Gao, and Yejin Choi. 2020. [PIQA: reasoning about physical commonsense in natural language](#). In *The*

Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, *The Thirty-Second Innovative Applications of Artificial Intelligence Conference*, IAAI 2020, *The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence*, EAAI 2020, New York, NY, USA, February 7-12, 2020, pages 7432–7439. AAAI Press.

Michael Boratko, Xiang Li, Tim O’Gorman, Rajarshi Das, Dan Le, and Andrew McCallum. 2020. [ProtoQA: A question answering dataset for prototypical common-sense reasoning](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1122–1136, Online. Association for Computational Linguistics.

Jan A. Botha, Manaal Faruqui, John Alex, Jason Baldridge, and Dipanjan Das. 2018. [Learning to split and rephrase from Wikipedia edit history](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 732–737, Brussels, Belgium. Association for Computational Linguistics.

Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.

Susan Carey and E. Bartlett. 1978. Acquiring a single new word. *Proceedings of the Stanford Child Language Conference*, 15:17–29.

Ankush Chatterjee, Kedhar Nath Narahari, Meghana Joshi, and Puneet Agrawal. 2019. [SemEval-2019 task 3: EmoContext contextual emotion detection in text](#). In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 39–48, Minneapolis, Minnesota, USA. Association for Computational Linguistics.

Michael Chen, Mike D’Arcy, Alisa Liu, Jared Fernandez, and Doug Downey. 2019. [CODAH: An adversarially-authored question answering dataset for common sense](#). In *Proceedings of the 3rd Workshop on Evaluating Vector Space Representations for NLP*, pages 63–69, Minneapolis, USA. Association for Computational Linguistics.

Wenhu Chen, Hongmin Wang, Jianshu Chen, Yunkai Zhang, Hong Wang, Shiyang Li, Xiyu Zhou, and William Yang Wang. 2020a. [Tabfact: A large-scale dataset for table-based fact verification](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

Xilun Chen, Asish Ghoshal, Yashar Mehdad, Luke Zettlemoyer, and Sonal Gupta. 2020b. [Low-resource domain adaptation for compositional task-oriented semantic parsing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5090–5100, Online. Association for Computational Linguistics.

708
709
710
711
712
713
714

715
716
717
718
719
720
721

722
723
724
725
726
727
728

729
730
731
732
733

734
735
736

737
738
739
740
741
742
743

744
745
746
747
748
749
750

751
752
753
754
755
756
757

758
759
760
761
762
763
764

878	Angela Fan, Yacine Jernite, Ethan Perez, David Grangier, Jason Weston, and Michael Auli. 2019. ELI5: Long form question answering . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 3558–3567, Florence, Italy. Association for Computational Linguistics.	934
879		935
880		936
881		937
882		938
883		939
884	Manaal Faruqui and Dipanjan Das. 2018. Identifying well-formed natural language questions . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 798–803, Brussels, Belgium. Association for Computational Linguistics.	940
885		941
886		942
887		943
888		944
889		945
890	Chelsea Finn. 2022. Deep multi-task and meta learning .	946
891	Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-agnostic meta-learning for fast adaptation of deep networks . In <i>International conference on machine learning</i> , pages 1126–1135. PMLR.	947
892		948
893		949
894		950
895	Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. Making pre-trained language models better few-shot learners . In <i>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)</i> , pages 3816–3830, Online. Association for Computational Linguistics.	951
896		952
897		953
898		954
899		955
900		956
901		957
902		958
903	Tianyu Gao, Xu Han, Ruobing Xie, Zhiyuan Liu, Fen Lin, Leyu Lin, and Maosong Sun. 2020. Neural snowball for few-shot relation learning .	959
904		960
905		961
906	Marta Garnelo, Dan Rosenbaum, Christopher Maddison, Tiago Ramalho, David Saxton, Murray Shanahan, Yee Whye Teh, Danilo Rezende, and SM Ali Eslami. 2018. Conditional neural processes . In <i>International Conference on Machine Learning</i> , pages 1704–1713. PMLR.	962
907		963
908		964
909		965
910		966
911		967
912	Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. 2019. SAMSum corpus: A human-annotated dialogue dataset for abstractive summarization . In <i>Proceedings of the 2nd Workshop on New Frontiers in Summarization</i> , pages 70–79, Hong Kong, China. Association for Computational Linguistics.	968
913		969
914		970
915		971
916		972
917		973
918		974
919	Andrew Gordon, Zornitsa Kozareva, and Melissa Roemmele. 2012. SemEval-2012 task 7: Choice of plausible alternatives: An evaluation of commonsense causal reasoning . In <i>*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)</i> , pages 394–398, Montréal, Canada. Association for Computational Linguistics.	975
920		976
921		977
922		978
923		979
924		980
925		981
926		982
927		983
928		984
929	Edward Grefenstette, Brandon Amos, Denis Yarats, Phu Mon Htut, Artem Molchanov, Franziska Meier, Douwe Kiela, Kyunghyun Cho, and Soumith Chintala. 2019. Generalized inner loop meta-learning . <i>arXiv preprint arXiv:1910.01727</i> .	985
930		986
931		987
932		988
933		989
	Jiatao Gu, Yong Wang, Yun Chen, Victor O. K. Li, and Kyunghyun Cho. 2018. Meta-learning for low-resource neural machine translation . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3622–3631, Brussels, Belgium. Association for Computational Linguistics.	990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

991	comprehension with contextual commonsense reasoning . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2391–2401, Hong Kong, China. Association for Computational Linguistics.	
992		
993		
994		
995		
996		
997		
998	Yukun Huang, Kun Qian, and Zhou Yu. 2022. Learning a better initialization for soft prompts via meta-learning . <i>arXiv preprint arXiv:2205.12471</i> .	
999		
1000		
1001	Kelvin Jiang, Dekun Wu, and Hui Jiang. 2019. FreebaseQA: A new factoid QA data set matching trivia-style question-answer pairs with Freebase . In <i>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)</i> , pages 318–323, Minneapolis, Minnesota. Association for Computational Linguistics.	
1002		
1003		
1004		
1005		
1006		
1007		
1008		
1009		
1010	Akhil Kedia, Sai Chetan Chinthakindi, and Wonho Ryu. 2021. Beyond reptile: Meta-learned dot-product maximization between gradients for improved single-task regularization . In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 407–420, Punta Cana, Dominican Republic. Association for Computational Linguistics.	
1011		
1012		
1013		
1014		
1015		
1016		
1017	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 <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 252–262, New Orleans, Louisiana. Association for Computational Linguistics.	
1018		
1019		
1020		
1021		
1022		
1023		
1024		
1025		
1026	Tushar Khot, Peter Clark, Michal Guerquin, Peter Jansen, and Ashish Sabharwal. 2020. Qasc: A dataset for question answering via sentence composition . <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 34(05):8082–8090.	
1027		
1028		
1029		
1030		
1031	Tushar Khot, Ashish Sabharwal, and Peter Clark. 2018. Scitail: A textual entailment dataset from science question answering . In <i>Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018</i> , pages 5189–5197. AAAI Press.	
1032		
1033		
1034		
1035		
1036		
1037		
1038		
1039		
1040	Byeongchang Kim, Hyunwoo Kim, and Gunhee Kim. 2019. Abstractive summarization of Reddit posts with multi-level memory networks . In <i>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)</i> , pages 2519–2531, Minneapolis, Minnesota. Association for Computational Linguistics.	
1041		
1042		
1043		
1044		
1045		
1046		
1047		
	Gregory Koch, Richard Zemel, Ruslan Salakhutdinov, et al. 2015. Siamese neural networks for one-shot image recognition . In <i>ICML deep learning workshop</i> , volume 2, page 0. Lille.	1048 1049 1050 1051
	Neema Kotonya and Francesca Toni. 2020. Explainable automated fact-checking for public health claims . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 7740–7754, Online. Association for Computational Linguistics.	1052 1053 1054 1055 1056 1057
	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: A benchmark for question answering research . <i>Transactions of the Association for Computational Linguistics</i> , 7:452–466.	1058 1059 1060 1061 1062 1063 1064 1065 1066
	Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. RACE: Large-scale Reading comprehension dataset from examinations . In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 785–794, Copenhagen, Denmark. Association for Computational Linguistics.	1067 1068 1069 1070 1071 1072 1073
	Brenden M. Lake, Tomer D. Ullman, Joshua B. Tenenbaum, and Samuel J. Gershman. 2017. Building machines that learn and think like people . <i>Behavioral and Brain Sciences</i> , 40:e253.	1074 1075 1076 1077
	Teven Le Scao and Alexander Rush. 2021. How many data points is a prompt worth? In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2627–2636, Online. Association for Computational Linguistics.	1078 1079 1080 1081 1082 1083
	Rémi Lebret, David Grangier, and Michael Auli. 2016. Neural text generation from structured data with application to the biography domain . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 1203–1213, Austin, Texas. Association for Computational Linguistics.	1084 1085 1086 1087 1088 1089
	Hung-yi Lee, Shang-Wen Li, and Ngoc Thang Vu. 2022. Meta learning for natural language processing: A survey . <i>arXiv preprint arXiv:2205.01500</i> .	1090 1091 1092
	Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, D. Kontokostas, Pablo N. Mendes, Sebastian Hellmann, M. Morsey, Patrick van Kleef, S. Auer, and C. Bizer. 2015. Dbpedia - a large-scale, multilingual knowledge base extracted from wikipedia. <i>Semantic Web</i> , 6:167–195.	1093 1094 1095 1096 1097 1098
	Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> ,	1099 1100 1101 1102

1103		pages 3045–3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	1159
1104			1160
1105			1161
1106	Hector J. Levesque, Ernest Davis, and Leora Morgenstern. 2012. The winograd schema challenge. In <i>Proceedings of the Thirteenth International Conference on Principles of Knowledge Representation and Reasoning</i> , KR’12, page 552–561. AAAI Press.		1162
1107			1163
1108			1164
1109			1165
1110			
1111	Omer Levy, Minjoon Seo, Eunsol Choi, and Luke Zettlemoyer. 2017. Zero-shot relation extraction via reading comprehension. In <i>Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017)</i> , pages 333–342, Vancouver, Canada. Association for Computational Linguistics.		1166
1112			1167
1113			1168
1114			1169
1115			1170
1116			1171
1117			
1118	Junyi Li, Tianyi Tang, Jian-Yun Nie, Ji-Rong Wen, and Xin Zhao. 2022. Learning to transfer prompts for text generation. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 3506–3518, Seattle, United States. Association for Computational Linguistics.		1172
1119			1173
1120			1174
1121			1175
1122			1176
1123			
1124	Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. In <i>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)</i> , pages 4582–4597, Online. Association for Computational Linguistics.		1177
1125			1178
1126			1179
1127			
1128			1180
1129			1181
1130			1182
1131			1183
1132			1184
1133	Xin Li and Dan Roth. 2002. Learning question classifiers. In <i>COLING 2002: The 19th International Conference on Computational Linguistics</i> .		1185
1134			
1135	Bill Yuchen Lin, Seyeon Lee, Rahul Khanna, and Xiang Ren. 2020a. Birds have four legs?! NumerSense: Probing Numerical Commonsense Knowledge of Pre-Trained Language Models. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 6862–6868, Online. Association for Computational Linguistics.		1186
1136			1187
1137			1188
1138			1189
1139			1190
1140			1191
1141			1192
1142	Bill Yuchen Lin, Wangchunshu Zhou, Ming Shen, Pei Zhou, Chandra Bhagavatula, Yejin Choi, and Xiang Ren. 2020b. CommonGen: A constrained text generation challenge for generative commonsense reasoning. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 1823–1840, Online. Association for Computational Linguistics.		1193
1143			
1144			1194
1145			1195
1146			1196
1147			1197
1148			
1149	Kevin Lin, Oyvind Taffjord, Peter Clark, and Matt Gardner. 2019. Reasoning over paragraph effects in situations. In <i>Proceedings of the 2nd Workshop on Machine Reading for Question Answering</i> , pages 58–62, Hong Kong, China. Association for Computational Linguistics.		1198
1150			1199
1151			1200
1152			1201
1153			1202
1154			
1155	Sen Lin, Li Yang, Deliang Fan, and Junshan Zhang. 2022. TRGP: Trust region gradient projection for continual learning. In <i>International Conference on Learning Representations</i> .		1203
1156			1204
1157			1205
1158			1206
			1207
			1208
			1209
			1210
			1211
			1212
			1213
			1214
			1215

1216	Julian J. McAuley and Jure Leskovec. 2013. Hidden factors and hidden topics: understanding rating dimensions with review text . In <i>Seventh ACM Conference on Recommender Systems, RecSys '13, Hong Kong, China, October 12-16, 2013</i> , pages 165–172. ACM.	Alex Nichol, Joshua Achiam, and John Schulman. 2018. On first-order meta-learning algorithms . <i>arXiv preprint arXiv:1803.02999</i> .	1271
1217			1272
1218			1273
1219			
1220		Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2020. Adversarial NLI: A new benchmark for natural language understanding . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4885–4901, Online. Association for Computational Linguistics.	1274
1221			1275
1222	Clara H. McCreery, Namit Katariya, Anitha Kannan, Manish Chablani, and Xavier Amatriain. 2020. Effective transfer learning for identifying similar questions: Matching user questions to COVID-19 faqs . In <i>KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020</i> , pages 3458–3465. ACM.		1276
1223			1277
1224			1278
1225			1279
1226			1280
1227		A. Othman and M. Jemni. 2012. English-asl gloss parallel corpus 2012: Aslg-pc12. In <i>5th Workshop on the Representation and Processing of Sign Languages: Interactions between Corpus and Lexicon LREC</i> .	1281
1228			1282
1229	Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? a new dataset for open book question answering . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2381–2391, Brussels, Belgium. Association for Computational Linguistics.		1283
1230			1284
1231		Bo Pang and Lillian Lee. 2005. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales . In <i>Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL'05)</i> , pages 115–124, Ann Arbor, Michigan. Association for Computational Linguistics.	1285
1232			1286
1233			1287
1234			1288
1235			1289
1236	Sewon Min, Mike Lewis, Luke Zettlemoyer, and Han-naneh Hajishirzi. 2022. MetaICL: Learning to learn in context . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2791–2809, Seattle, United States. Association for Computational Linguistics.		1290
1237			1291
1238		Dimitris Pappas, Petros Stavropoulos, Ion Androustopoulos, and Ryan McDonald. 2020. BioMRC: A dataset for biomedical machine reading comprehension . In <i>Proceedings of the 19th SIGBioMed Workshop on Biomedical Language Processing</i> , pages 140–149, Online. Association for Computational Linguistics.	1292
1239			1293
1240			1294
1241			1295
1242			1296
1243	Nikhil Mishra, Mostafa Rohaninejad, Xi Chen, and Pieter Abbeel. 2018. A simple neural attentive meta-learner . In <i>International Conference on Learning Representations</i> .		1297
1244			1298
1245		Fabio Petroni, Patrick Lewis, Aleksandra Piktus, Tim Rocktäschel, Yuxiang Wu, Alexander H. Miller, and Sebastian Riedel. 2020. How context affects language models' factual predictions . In <i>Automated Knowledge Base Construction</i> .	1299
1246			1300
1247	Ioannis Mollas, Zoe Chrysopoulou, Stamatis Karlos, and Grigorios Tsoumakas. 2020. Ethos: an online hate speech detection dataset . <i>ArXiv preprint</i> , abs/2006.08328.		1301
1248			1302
1249			1303
1250		Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2463–2473, Hong Kong, China. Association for Computational Linguistics.	1304
1251	Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R. Bowman. 2020. CrowS-pairs: A challenge dataset for measuring social biases in masked language models . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1953–1967, Online. Association for Computational Linguistics.		1305
1252			1306
1253			1307
1254			1308
1255			1309
1256			1310
1257			1311
1258	Courtney Napoles, Matthew Gormley, and Benjamin Van Durme. 2012. Annotated Gigaword . In <i>Proceedings of the Joint Workshop on Automatic Knowledge Base Construction and Web-scale Knowledge Extraction (AKBC-WEKEX)</i> , pages 95–100, Montréal, Canada. Association for Computational Linguistics.	Mohammad Taher Pilehvar and Jose Camacho-Collados. 2019. WiC: the word-in-context dataset for evaluating context-sensitive meaning representations . In <i>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)</i> , pages 1267–1273, Minneapolis, Minnesota. Association for Computational Linguistics.	1313
1259			1314
1260			1315
1261			1316
1262			1317
1263			1318
1264	Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1797–1807, Brussels, Belgium. Association for Computational Linguistics.		1319
1265			1320
1266			1321
1267		Amir Pouran Ben Veyseh, Franck Dernoncourt, Quan Hung Tran, and Thien Huu Nguyen. 2020. What does this acronym mean? introducing a new dataset for acronym identification and disambiguation . In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 3285–	1322
1268			1323
1269			1324
1270			1325
			1326
			1327

1328	3301, Barcelona, Spain (Online). International Committee on Computational Linguistics.	
1329		
1330	Chengwei Qin and Shafiq Joty. 2022a. Continual few-shot relation learning via embedding space regularization and data augmentation . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2776–2789, Dublin, Ireland. Association for Computational Linguistics.	
1331		
1332		
1333		
1334		
1335		
1336		
1337	Chengwei Qin and Shafiq Joty. 2022b. LFPT5: A unified framework for lifelong few-shot language learning based on prompt tuning of t5 . In <i>International Conference on Learning Representations</i> .	
1338		
1339		
1340		
1341	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer . <i>arXiv preprint arXiv:1910.10683</i> .	
1342		
1343		
1344		
1345		
1346	Nazneen Fatema Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. Explain yourself! leveraging language models for commonsense reasoning . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4932–4942, Florence, Italy. Association for Computational Linguistics.	
1347		
1348		
1349		
1350		
1351		
1352		
1353	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 2383–2392, Austin, Texas. Association for Computational Linguistics.	
1354		
1355		
1356		
1357		
1358		
1359	Anna Rogers, Olga Kovaleva, Matthew Downey, and Anna Rumshisky. 2020. Getting closer to ai complete question answering: A set of prerequisite real tasks . <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 34(05):8722–8731.	
1360		
1361		
1362		
1363		
1364	Amrita Saha, Rahul Aralikkatte, Mitesh M. Khapra, and Karthik Sankaranarayanan. 2018. DuoRC: Towards complex language understanding with paraphrased reading comprehension . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1683–1693, Melbourne, Australia. Association for Computational Linguistics.	
1365		
1366		
1367		
1368		
1369		
1370		
1371		
1372	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavathula, and Yejin Choi. 2020. Winogrande: An adversarial winograd schema challenge at scale . <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 34(05):8732–8740.	
1373		
1374		
1375		
1376		
1377	Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. 2016. Meta-learning with memory-augmented neural networks . In <i>International conference on machine learning</i> , pages 1842–1850. PMLR.	
1378		
1379		
1380		
1381		
	Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. 2019. Social IQa: Commonsense reasoning about social interactions . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 4463–4473, Hong Kong, China. Association for Computational Linguistics.	1382
		1383
		1384
		1385
		1386
		1387
		1388
		1389
		1390
	Elvis Saravia, Hsien-Chi Toby Liu, Yen-Hao Huang, Junlin Wu, and Yi-Shin Chen. 2018. CARER: Contextualized affect representations for emotion recognition . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3687–3697, Brussels, Belgium. Association for Computational Linguistics.	1391
		1392
		1393
		1394
		1395
		1396
		1397
	Timo Schick and Hinrich Schütze. 2021. Exploiting cloze-questions for few-shot text classification and natural language inference . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 255–269, Online. Association for Computational Linguistics.	1398
		1399
		1400
		1401
		1402
		1403
		1404
	Jurgen Schmidhuber. 1987. Evolutionary principles in self-referential learning. on learning now to learn: The meta-meta-meta...-hook . Diploma thesis, Technische Universität München, Germany, 14 May.	1405
		1406
		1407
		1408
	Emily Sheng and David Uthus. 2020. Investigating societal biases in a poetry composition system . In <i>Proceedings of the Second Workshop on Gender Bias in Natural Language Processing</i> , pages 93–106, Barcelona, Spain (Online). Association for Computational Linguistics.	1409
		1410
		1411
		1412
		1413
		1414
	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 <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 4222–4235, Online. Association for Computational Linguistics.	1415
		1416
		1417
		1418
		1419
		1420
		1421
	Damien Sileo, Tim Van De Cruys, Camille Pradel, and Philippe Muller. 2019. Mining discourse markers for unsupervised sentence representation learning . In <i>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)</i> , pages 3477–3486, Minneapolis, Minnesota. Association for Computational Linguistics.	1422
		1423
		1424
		1425
		1426
		1427
		1428
		1429
		1430
	Jake Snell, Kevin Swersky, and Richard Zemel. 2017. Prototypical networks for few-shot learning . <i>Advances in neural information processing systems</i> , 30.	1431
		1432
		1433
	Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank . In <i>Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing</i> , pages	1434
		1435
		1436
		1437
		1438
		1439

1440	1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.		
1441			
1442	Kai Sun, Dian Yu, Jianshu Chen, Dong Yu, Yejin Choi, and Claire Cardie. 2019. DREAM: A challenge data set and models for dialogue-based reading comprehension . <i>Transactions of the Association for Computational Linguistics</i> , 7:217–231.		
1443			
1444			
1445			
1446			
1447	Oyvind Tafjord, Peter Clark, Matt Gardner, Wen-tau Yih, and Ashish Sabharwal. 2019a. Quarel: A dataset and models for answering questions about qualitative relationships . <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 33(01):7063–7071.		
1448			
1449			
1450			
1451			
1452	Oyvind Tafjord, Matt Gardner, Kevin Lin, and Peter Clark. 2019b. QuaRTz: An open-domain dataset of qualitative relationship questions . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 5941–5946, Hong Kong, China. Association for Computational Linguistics.		
1453			
1454			
1455			
1456			
1457			
1458			
1459			
1460	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. CommonsenseQA: A question answering challenge targeting commonsense knowledge . In <i>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)</i> , pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.		
1461			
1462			
1463			
1464			
1465			
1466			
1467			
1468			
1469	Niket Tandon, Bhavana Dalvi, Keisuke Sakaguchi, Peter Clark, and Antoine Bosselut. 2019. WIQA: A dataset for “what if...” reasoning over procedural text . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 6076–6085, Hong Kong, China. Association for Computational Linguistics.		
1470			
1471			
1472			
1473			
1474			
1475			
1476			
1477			
1478	James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a large-scale dataset for fact extraction and VERification . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 809–819, New Orleans, Louisiana. Association for Computational Linguistics.		
1479			
1480			
1481			
1482			
1483			
1484			
1485			
1486			
1487	Eleni Triantafyllou, Richard S. Zemel, and Raquel Urtasun. 2017. Few-shot learning through an information retrieval lens .		
1488			
1489			
1490	Sowmya Vajjala and Ivana Lučić. 2018. OneStopEnglish corpus: A new corpus for automatic readability assessment and text simplification . In <i>Proceedings of the Thirteenth Workshop on Innovative Use of NLP for Building Educational Applications</i> , pages 297–304, New Orleans, Louisiana. Association for Computational Linguistics.		
1491			
1492			
1493			
1494			
1495			
1496			
		Oriol Vinyals, Charles Blundell, Timothy Lillicrap, koray kavukcuoglu, and Daan Wierstra. 2016. Matching networks for one shot learning . In <i>Advances in Neural Information Processing Systems</i> , volume 29. Curran Associates, Inc.	1497
			1498
			1499
			1500
			1501
		Tu Vu, Brian Lester, Noah Constant, Rami Al-Rfou’, and Daniel Cer. 2022. SPoT: Better frozen model adaptation through soft prompt transfer . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5039–5059, Dublin, Ireland. Association for Computational Linguistics.	1502
			1503
			1504
			1505
			1506
			1507
			1508
		William Yang Wang. 2017. “liar, liar pants on fire”: A new benchmark dataset for fake news detection . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 422–426, Vancouver, Canada. Association for Computational Linguistics.	1509
			1510
			1511
			1512
			1513
			1514
		Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mohananey, Wei Peng, Sheng-Fu Wang, and Samuel R. Bowman. 2020. BLiMP: The benchmark of linguistic minimal pairs for English . <i>Transactions of the Association for Computational Linguistics</i> , 8:377–392.	1515
			1516
			1517
			1518
			1519
			1520
		Alex Warstadt, Amanpreet Singh, and Samuel R. Bowman. 2019. Neural network acceptability judgments . <i>Transactions of the Association for Computational Linguistics</i> , 7:625–641.	1521
			1522
			1523
			1524
		Johannes Welbl, Nelson F. Liu, and Matt Gardner. 2017. Crowdsourcing multiple choice science questions . In <i>Proceedings of the 3rd Workshop on Noisy User-generated Text</i> , pages 94–106, Copenhagen, Denmark. Association for Computational Linguistics.	1525
			1526
			1527
			1528
			1529
		Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 1112–1122, New Orleans, Louisiana. Association for Computational Linguistics.	1530
			1531
			1532
			1533
			1534
			1535
			1536
			1537
			1538
		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 <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	1539
			1540
			1541
			1542
			1543
			1544
			1545
			1546
			1547
			1548
			1549
			1550
		Tomer Wolfson, Mor Geva, Ankit Gupta, Matt Gardner, Yoav Goldberg, Daniel Deutch, and Jonathan	1551
			1552

1553	Berant. 2020. Break it down: A question understanding benchmark . <i>Transactions of the Association for Computational Linguistics</i> , 8:183–198.	1610
1554		1611
1555		1612
1556	Wenhan Xiong, Jiawei Wu, Hong Wang, Vivek Kulka-	
1557	rni, Mo Yu, Shiyu Chang, Xiaoxiao Guo, and	
1558	William Yang Wang. 2019. TWEETQA: A social	1614
1559	media focused question answering dataset . In <i>Pro-</i>	1615
1560	<i>ceedings of the 57th Annual Meeting of the Asso-</i>	1616
1561	<i>ciation for Computational Linguistics</i> , pages 5020–	1617
1562	5031, Florence, Italy. Association for Computational	1618
1563	Linguistics.	1619
1564	Yi Yang, Wen-tau Yih, and Christopher Meek. 2015.	1620
1565	WikiQA: A challenge dataset for open-domain ques-	1621
1566	tion answering . In <i>Proceedings of the 2015 Con-</i>	1622
1567	<i>ference on Empirical Methods in Natural Language</i>	1623
1568	<i>Processing</i> , pages 2013–2018, Lisbon, Portugal. As-	1624
1569	sociation for Computational Linguistics.	
1570	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio,	
1571	William Cohen, Ruslan Salakhutdinov, and Christo-	
1572	pher D. Manning. 2018. HotpotQA: A dataset for	1625
1573	diverse, explainable multi-hop question answering .	1626
1574	In <i>Proceedings of the 2018 Conference on Empiri-</i>	1627
1575	<i>cal Methods in Natural Language Processing</i> , pages	1628
1576	2369–2380, Brussels, Belgium. Association for Com-	1629
1577	putational Linguistics.	1630
1578	Qinyuan Ye, Bill Yuchen Lin, and Xiang Ren. 2021.	
1579	CrossFit: A few-shot learning challenge for cross-	1631
1580	task generalization in NLP . In <i>Proceedings of the</i>	1632
1581	<i>2021 Conference on Empirical Methods in Natural</i>	1633
1582	<i>Language Processing</i> , pages 7163–7189, Online and	1634
1583	Punta Cana, Dominican Republic. Association for	1635
1584	Computational Linguistics.	1636
1585	Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga,	
1586	Dongxu Wang, Zifan Li, James Ma, Irene Li, Qing-	
1587	ning Yao, Shanelle Roman, Zilin Zhang, and Dragomir	
1588	Radev. 2018. Spider: A large-scale human-labeled	1637
1589	dataset for complex and cross-domain semantic pars-	1638
1590	ing and text-to-SQL task . In <i>Proceedings of the 2018</i>	
1591	<i>Conference on Empirical Methods in Natural Lan-</i>	1639
1592	<i>guage Processing</i> , pages 3911–3921, Brussels, Bel-	1640
1593	gium. Association for Computational Linguistics.	1641
1594	Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin	
1595	Choi. 2018. SWAG: A large-scale adversarial dataset	1642
1596	for grounded commonsense inference . In <i>Proceed-</i>	1643
1597	<i>ings of the 2018 Conference on Empirical Methods in</i>	1644
1598	<i>Natural Language Processing</i> , pages 93–104, Brus-	1645
1599	sels, Belgium. Association for Computational Lin-	
1600	guistics.	
1601	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali	
1602	Farhadi, and Yejin Choi. 2019. HellaSwag: Can a ma-	1646
1603	chine really finish your sentence? In <i>Proceedings of</i>	1647
1604	<i>the 57th Annual Meeting of the Association for Com-</i>	1648
1605	<i>putational Linguistics</i> , pages 4791–4800, Florence,	1649
1606	Italy. Association for Computational Linguistics.	1650
1607	Hao Zhang, Jae Ro, and Richard Sproat. 2020. Semi-	1651
1608	supervised URL segmentation with recurrent neu-	1652
1609	ral networks pre-trained on knowledge graph enti-	1653
	ties . In <i>Proceedings of the 28th International Con-</i>	1654
	<i>ference on Computational Linguistics</i> , pages 4667–	1655
	4675, Barcelona, Spain (Online). International Com-	1656
	mittee on Computational Linguistics.	1657
	Rui Zhang and Joel Tetreault. 2019. This email could	1658
	save your life: Introducing the task of email subject	1659
	line generation . In <i>Proceedings of the 57th Annual</i>	1660
	<i>Meeting of the Association for Computational Lin-</i>	1661
	<i>guistics</i> , pages 446–456, Florence, Italy. Association	1662
	for Computational Linguistics.	1663
	Sheng Zhang, X. Liu, J. Liu, Jianfeng Gao, Kevin	
	Duh, and Benjamin Van Durme. 2018. Record:	1664
	Bridging the gap between human and machine com-	1665
	monsense reading comprehension . <i>ArXiv preprint</i> ,	
	abs/1810.12885.	
	Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015.	
	Character-level convolutional networks for text clas-	1666
	sification . In <i>Advances in Neural Information Pro-</i>	1667
	<i>cessing Systems 28: Annual Conference on Neural In-</i>	1668
	<i>formation Processing Systems 2015, December 7-12,</i>	1669
	<i>2015, Montreal, Quebec, Canada</i> , pages 649–657.	
	Yuan Zhang, Jason Baldridge, and Luheng He. 2019.	
	PAWS: Paraphrase adversaries from word scrambling .	1670
	In <i>Proceedings of the 2019 Conference of the North</i>	1671
	<i>American Chapter of the Association for Computa-</i>	1672
	<i>tional Linguistics: Human Language Technologies,</i>	1673
	<i>Volume 1 (Long and Short Papers)</i> , pages 1298–1308,	1674
	Minneapolis, Minnesota. Association for Computa-	1675
	tional Linguistics.	1676
	Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and	
	Sameer Singh. 2021. Calibrate before use: Improv-	1677
	ing few-shot performance of language models . In	1678
	<i>Proceedings of the 38th International Conference</i>	1679
	<i>on Machine Learning</i> , volume 139 of <i>Proceedings</i>	1680
	<i>of Machine Learning Research</i> , pages 12697–12706.	1681
	PMLR.	1682
	Ruiqi Zhong, Kristy Lee, Zheng Zhang, and Dan Klein.	
	2021. Adapting language models for zero-shot learn-	1683
	ing by meta-tuning on dataset and prompt collections .	1684
	In <i>Findings of the Association for Computational</i>	1685
	<i>Linguistics: EMNLP 2021</i> , pages 2856–2878, Punta	1686
	Cana, Dominican Republic. Association for Compu-	1687
	tational Linguistics.	1688
	Victor Zhong, Caiming Xiong, and Richard Socher.	
	2017. Seq2sql: Generating structured queries	1689
	from natural language usin . <i>ArXiv preprint</i> ,	1690
	abs/1709.00103.	1691
	Ben Zhou, Daniel Khoshabi, Qiang Ning, and Dan Roth.	
	2019. “going on a vacation” takes longer than “go-	1692
	ing for a walk”: A study of temporal commonsense	1693
	understanding . In <i>Proceedings of the 2019 Confer-</i>	1694
	<i>ence on Empirical Methods in Natural Language Pro-</i>	1695
	<i>cessing and the 9th International Joint Conference</i>	1696
	<i>on Natural Language Processing (EMNLP-IJCNLP)</i> ,	1697
	pages 3363–3369, Hong Kong, China. Association	1698
	for Computational Linguistics.	1699

Beier Zhu, Yulei Niu, Yucheng Han, Yue Wu, and Hanwang Zhang. 2022. [Prompt-aligned gradient for prompt tuning](#). *arXiv preprint arXiv:2205.14865*.

A Appendix

A.1 Task List

We report the full list of tasks used in ten different settings in [Table 8](#). All tasks are taken from CROSSFIT ([Ye et al., 2021](#)).

A.2 Detailed Results on Every Target Task

We mainly report average relative gain (ARG) in our experiments (§6). In this section, we show detailed results on each target task in [Fig. 4](#) ~ [Fig. 13](#).

A.3 Details of Sampled Tasks

We sample $\{12, 24\}$ tasks from the original 45 source tasks in the Cls→Cls setting to investigate the influence of the number of source tasks. The details of sampled tasks are shown in [Table 9](#).

A.4 Details of Hyperparameters

For meta-training, we set the inner and outer learning rates to $3e-5$ and $5e-1$, respectively. We use 5000 for total training steps. We set the inner batch size to 2, 4 and 4, and inner update steps to 1, 1 and 10 for MAML, FoMAML and Reptile, respectively. For multi-task learning, we set the learning rate, batch size and number of epochs to $5e-1$, 4 and 20, respectively. For MAML, we select the inner learning rate from $\{2e-5, 3e-5, 5e-5\}$, the outer learning rate from $\{2e-1, 3e-1, 5e-1\}$, and total training steps from $\{2500, 5000, 10000\}$. We adopt the same three hyperparameters for FoMAML and Reptile. The search range for the inner update steps of Reptile is $\{2, 4, 6, 8, 10\}$. For multi-task learning, we select the learning rate from $\{2e-1, 3e-1, 5e-1\}$, the batch size from $\{2, 4, 6, 8\}$, and the number of epochs from $\{5, 10, 20\}$.

For downstream learning, we mainly follow the settings in [Ye et al. \(2021\)](#). For prompt tuning, we select the learning rate from $\{5e-1, 4e-1, 3e-1, 2e-1\}$ based on the validation performance. For fine-tuning, the search range for the learning rate is $\{5e-4, 3e-4, 2e-4, 1e-4\}$. We set the batch size, total training steps and evaluation interval to 8, 3000 and 50, respectively.

A.5 Hyperparameter Sensitivity Analysis

As mentioned in [Appendix A.4](#), for MAML, we select the inner learning rate from $\{2e-5, 3e-5, 5e-5\}$, the outer learning rate from $\{2e-1, 3e-1, 5e-1\}$, and total training steps from $\{2500, 5000, 10000\}$ in the R→R setting. The best validation performance (10.14% ARG) is achieved with $\{3e-5, 5e-1, 5000\}$, while the worst validation ARG is -16.21% when using $\{5e-5, 2e-1, 2500\}$. We can see that MPT is quite sensitive to hyperparameters. It performs even worse than PT with inappropriate hyperparameters.

A.6 Task Similarity Analysis

As discussed in §6, we use the correlation between input subspaces for two tasks as the similarity score between them. Detailed results of randomly picked similar and dissimilar task pairs are shown in [Table 6](#).

A.7 Case Study

To take a closer look at the influence of different source task types on a particular target task, we further conduct a case study where we ensure that the task under consideration appears in the target task partitions.² Results are shown in [Table 7](#); for example, the first block indicates that Amazon_Polarity appears as a target task in both R→R and Cls→Cls settings. We can observe that there is no consistent conclusion on how we should choose the source tasks for a specific target task, which is consistent with our view in Q3.

A.8 Limitations

Although comprehensive, our study of MPT in this work has couple of limitations:

- As mentioned in §5.3, because of infeasibility to search for optimal hyperparameters for each of the meta learning methods in each of the ten settings, we choose to use the R→R setting as our main representative setting. This could be one of the reasons for MPT underperforming MTL in some non-classification tasks (noted in §6-Q1).
- We mainly focus on how upstream meta learning can improve the performance on target tasks. However, meta learning also enables faster convergence. We leave how it could help reduce the convergence time of PT as future work.

²As before, we ensure it does not appear in the source.

	Task Pair Index					Average
	1	2	3	4	5	
Similar	0.772	0.695	0.754	0.819	0.802	0.768
Dissimilar	0.326	0.311	0.283	0.315	0.297	0.306

Table 6: Similarity scores of randomly picked similar and dissimilar task pairs.

Target Task	Partition	Δ_{MPT}	Δ_{MTL}
Amazon_Polarity	R→R	3.10	2.25
	Cls→Cls	7.40	10.45
AI2_ARC	R→R	12.54	5.55
	Both→Non-Cls	8.17	6.69
Samsun	R→R	1.97	6.77
	Both→Non-Cls	2.50	5.71
Superglue-Copa	Both→Non-Cls	1.20	10.00
	QA→QA	-3.20	4.80

Table 7: Relative gain in % for MPT and MTL when the same target task appears in different partitions.

1757 Aside from that, meta prompt tuning (MPT) as a
1758 method has a limitation that it is Memory-intensive.
1759 Optimization-based meta learning methods, especially MAML, are memory-intensive, which limits
1760 the tuning of the inner batch size and inner update steps (Appendix A.4). One potential solution
1761 is to build more memory-efficient meta learning
1762 libraries.
1763
1764

<p>Partition: Random Source glue-mrpc, math_qa, quarel, e2e_nlg_cleaned, tweet_eval-stance_atheism, lama-squad, tab_fact, aqua_rat, tweet_eval-emoji, glue-wnli, codah, tweet_eval-offensive, wiki_qa, blimp-ellipsis_n_bar_1, openbookqa, sms_spam, acronym_identification, blimp-determiner_noun_agreement_with_adj_irregular_1, ethos-national_origin, spider, hellaswag, superglue-wsc, numer_sense, ade_corpus_v2-dosage, blimp-ellipsis_n_bar_2, kilt_ay2, squad-no_context, google_wellformed_query, xsum, wiqa, tweet_eval-stance_abortion, reddit_tifu-tldr, ade_corpus_v2-effect, qa_srl, ethos-religion, commonsense_qa, biomrc, superglue-multirc, ethos-race, eli5-askh, glue-qqp, paws, ethos-directed_vs_generalized, glue-sst2, tweet_eval-hate, glue-rte, blimp-anaphor_number_agreement, lama-conceptnet, hate_speech_offensive, superglue-wic, boolq, kilt_hotpotqa, quartz-no_knowledge, aslg_pc12, sick, tweet_eval-stance_climate, tweet_eval-sentiment, crows_pairs, glue-mnli, medical_questions_pairs, break-QDMR-high-level, qasc, imdb, ethos-gender, trec-finegrained, adversarialqa, onestop_english, web_questions, duorc, swag, proto_qa, scitail, tweet_eval-stance_feminist, limit, common_gen, scicite, blimp-irregular_past_participle_adjectives, social_i_qa, anli, kilt_zsre, cosmos_qa, superglue-record, squad-with_context, emotion, blimp-existential_there_quantifiers_1, race-middle, kilt_wow, sciq, wino_grande, rotten_tomatoes, superglue-cb, poem_sentiment, ropes, reddit_tifu-title, piqa, climate_fever, lama-google_re, search_qa, mc_taco, blimp-wh_questions_object_gap, hotpot_qa, emo, kilt_nq, kilt_trex, quartz-with_knowledge, dbpedia_14, yahoo_answers_topics, superglue-copa, blimp-anaphor_gender_agreement, hate_speech18, gigaword, multi_news, aesc, quail</p>
<p>Partition: Random Target quoref, wiki_split, ethos-disability, yelp_polarity, superglue-rte, glue-cola, ethos-sexual_orientation, blimp-sentential_negation_npi_scope, ai2_arc, amazon_polarity, race-high, blimp-sentential_negation_npi_licensor_present, tweet_eval-irony, crawl_domain, freebase_qa, glue-qnli, hatexplain, ag_news, circa, samsun</p>
<p>Partition: Classification Source superglue-rte, tweet_eval-sentiment, discovery, glue-rte, superglue-wsc, scicite, glue-mrpc, tweet_eval-stance_hillary, tweet_eval-offensive, emotion, hatexplain, glue-cola, sick, paws, ethos-sexual_orientation, glue-qqp, tweet_eval-emotion, sms_spam, health_fact, glue-mnli, imdb, ethos-disability, glue-wnli, scitail, trec-finegrained, yahoo_answers_topics, liar, glue-sst2, tweet_eval-stance_abortion, circa, tweet_eval-stance_climate, glue-qnli, tweet_eval-emoji, ethos-directed_vs_generalized, ade_corpus_v2-classification, ag_news, hate_speech_offensive, superglue-wic, google_wellformed_query, tweet_eval-irony, ethos-gender, onestop_english, trec, rotten_tomatoes, kilt_fever</p>
<p>Partition: Non-Classification Source ade_corpus_v2-dosage, art, biomrc, blimp-anaphor_number_agreement, blimp-ellipsis_n_bar_2, blimp-sentential_negation_npi_licensor_present, blimp-sentential_negation_npi_scope, break-QDMR-high-level, commonsense_qa, crows_pairs, dream, duorc, eli5-asks, eli5-eli5, freebase_qa, gigaword, hellaswag, hotpot_qa, kilt_ay2, kilt_hotpotqa, kilt_trex, kilt_zsre, lama-conceptnet, lama-google_re, lama-squad, math_qa, numer_sense, openbookqa, piqa, proto_qa, qa_srl, quarel, quartz-no_knowledge, race-high, reddit_tifu-title, reddit_tifu-tldr, ropes, sciq, social_i_qa, spider, superglue-multirc, wiki_bio, wikisql, xsum, yelp_review_full</p>
<p>Partition: Both (Classification + Non-Classification) Source ade_corpus_v2-dosage, biomrc, blimp-ellipsis_n_bar_2, blimp-sentential_negation_npi_scope, commonsense_qa, crows_pairs, duorc, hellaswag, kilt_zsre, lama-google_re, lama-squad, math_qa, numer_sense, openbookqa, piqa, proto_qa, quartz-no_knowledge, race-high, reddit_tifu-tldr, ropes, sciq, wiki_bio, discovery, emotion, ethos-disability, ethos-sexual_orientation, glue-cola, glue-mnli, glue-mrpc, glue-qqp, glue-rte, glue-wnli, hatexplain, health_fact, imdb, paws, scicite, sick, sms_spam, superglue-rte, superglue-wsc, tweet_eval-emotion, tweet_eval-offensive, tweet_eval-sentiment, tweet_eval-stance_hillary</p>
<p>Partition: Classification Target superglue-cb,dbpedia_14, wiki_qa, emo, yelp_polarity, ethos-religion, amazon_polarity, tab_fact, anli, ethos-race</p>
<p>Partition: Non-Classification Target multi_news, superglue-copa, quail, blimp-anaphor_gender_agreement, common_gen, acronym_identification, quoref, wiki_split, ai2_arc, break-QDMR, crawl_domain, samsun</p>
<p>Partition: QA Source biomrc, boolq, freebase_qa, hotpot_qa, kilt_hotpotqa, kilt_nq, kilt_trex, kilt_zsre, lama-conceptnet, lama-google_re, lama-squad, lama-trex, mc_taco, numer_sense, quoref, ropes, search_qa, squad-no_context, superglue-multirc, superglue-record, tweet_qa, web_questions</p>
<p>Partition: Non-QA Source hate_speech_offensive, google_wellformed_query, circa, glue-sst2, scitail, emo, ag_news, art, paws, kilt_ay2, glue-qnli, ade_corpus_v2-classification, hatexplain, emotion, glue-qqp, kilt_fever, dbpedia_14, glue-mnli, discovery, gigaword, amazon_polarity, tab_fact, tweet_eval-emoji, tweet_eval-offensive, tweet_eval-sentiment, imdb, liar, anli, wikisql, xsum, yahoo_answers_topics, yelp_polarity, yelp_review_full</p>
<p>Partition: QA Target ai2_arc, codah, cosmos_qa, dream, hellaswag, qasc, quail, quarel, quartz-no_knowledge, quartz-with_knowledge, sciq, superglue-copa, swag, wino_grande, wiqa</p>
<p>Partition: Non-Paraphrase Classification Source ade_corpus_v2-classification, ag_news, amazon_polarity, anli, circa, climate_fever, dbpedia_14, discovery, emo, emotion, ethos-directed_vs_generalized, ethos-disability, ethos-gender, ethos-national_origin, ethos-race, ethos-religion, ethos-sexual_orientation, financial_phrasebank, glue-cola, glue-mnli, glue-qnli, glue-rte, glue-sst2, glue-wnli, google_wellformed_query, hate_speech18, hate_speech_offensive, hatexplain, health_fact, imdb, kilt_fever, liar, onestop_english, poem_sentiment, rotten_tomatoes, scicite, scitail, sick, sms_spam, superglue-cb, superglue-rte, superglue-wic, superglue-wsc, tab_fact, trec, trec-finegrained, tweet_eval-emoji, tweet_eval-emotion, tweet_eval-hate, tweet_eval-irony, tweet_eval-offensive, tweet_eval-sentiment, tweet_eval-stance_abortion, tweet_eval-stance_atheism, tweet_eval-stance_climate, tweet_eval-stance_feminist, tweet_eval-stance_hillary, wiki_qa, yahoo_answers_topics, yelp_polarity</p>
<p>Partition: Paraphrase Target glue-mrpc, glue-qqp, medical_questions_pairs, paws</p>

Table 8: Full datasets for all settings described in Section 5.1. We provide references for all datasets in Table 10.

<p>12 source tasks superglue-rte, tweet_eval-sentiment, discovery, glue-rte, hatexplain, glue-cola, health_fact, glue-mnli, imdb, ethos-disability, glue-wnli, scitail</p>
<p>24 source tasks superglue-rte, tweet_eval-sentiment, discovery, glue-rte, superglue-wsc, scicite, hatexplain, glue-cola, tweet_eval-emotion, sms_spam, health_fact, glue-mnli, imdb, ethos-disability, glue-wnli, scitail, glue-sst2, tweet_eval-stance_abortion, glue-qnli, ethos-directed_vs_generalized, ag_news, hate_speech_offensive, ethos-gender, kilt_fever</p>

Table 9: Details of sampled {12, 24} tasks for investigating the impact of the number of source tasks.

Task Name	Reference
eli5-eli5	Fan et al. 2019
ethos-race	Mollas et al. 2020
tweet_qa	Xiong et al. 2019
tweet_eval-stance_hillary	Barbieri et al. 2020
piqa	Bisk et al. 2020
acronym_identification	Pouran Ben Veysch et al. 2020
wiki_split	Botha et al. 2018
scitail	Khot et al. 2018
emotion	Saravia et al. 2018
medical_questions_pairs	McCreery et al. 2020
blimp-anaphor_gender_agreement	Warstadt et al. 2020
sciq	Welbl et al. 2017
paws	Zhang et al. 2019
yelp_review_full	Zhang et al. 2015; (link)
freebase_qa	Jiang et al. 2019
anli	Nie et al. 2020
quartz-with_knowledge	Tafjord et al. 2019b
hatexplain	Mathew et al. 2020
yahoo_answers_topics	(link)
search_qa	Dunn et al. 2017
tweet_eval-stance_feminist	Barbieri et al. 2020
codah	Chen et al. 2019
lama-squad	Petroni et al. 2019, 2020
superglue-record	Zhang et al. 2018
spider	Yu et al. 2018
mc_taco	Zhou et al. 2019
glue-mrpc	Dolan and Brockett 2005
kilt_fever	Thorne et al. 2018
eli5-asks_qa	Fan et al. 2019
imdb	Maas et al. 2011
tweet_eval-stance_abortion	Barbieri et al. 2020
aqua_rat	Ling et al. 2017
duorc	Saha et al. 2018
lama-trex	Petroni et al. 2019, 2020
tweet_eval-stance_atheism	Barbieri et al. 2020
ropes	Lin et al. 2019
squad-no_context	Rajpurkar et al. 2016
superglue-rte	Dagan et al. 2005
qasc	Khot et al. 2020
hate_speech_offensive	Davidson et al. 2017
trec-finegrained	Li and Roth 2002; Hovy et al. 2001
glue-wnli	Levesque et al. 2012
yelp_polarity	Zhang et al. 2015; (link)
kilt_hotpotqa	Yang et al. 2018
glue-sst2	Socher et al. 2013
xsum	Narayan et al. 2018
tweet_eval-offensive	Barbieri et al. 2020
aeslc	Zhang and Tetreault 2019
emo	Chatterjee et al. 2019
hellaswag	Zellers et al. 2019
social_i_qa	Sap et al. 2019
kilt_wow	Dinan et al. 2019
scicite	Cohan et al. 2019
superglue-wsc	Levesque et al. 2012
hate_speech18	de Gibert et al. 2018
adversarialqa	Bartolo et al. 2020
break-QDMR	Wolfson et al. 2020
dream	Sun et al. 2019
circa	Louis et al. 2020
wiki_qa	Yang et al. 2015
ethos-directed_vs_generalized	Mollas et al. 2020
wiqa	Tandon et al. 2019
poem_sentiment	Sheng and Uthus 2020
kilt_ay2	Hoffart et al. 2011
cosmos_qa	Huang et al. 2019
reddit_tifu-title	Kim et al. 2019
superglue-cb	de Marneffe et al. 2019
kilt_nq	Kwiatkowski et al. 2019
quarel	Tafjord et al. 2019a
race-high	Lai et al. 2017
wino_grande	Sakaguchi et al. 2020
break-QDMR-high-level	Wolfson et al. 2020
tweet_eval-irony	Barbieri et al. 2020
liar	Wang 2017
openbookqa	Mihaylov et al. 2018
superglue-multirc	Khashabi et al. 2018
race-middle	Lai et al. 2017
quoref	Dasigi et al. 2019
cos_e	Rajani et al. 2019
reddit_tifu-tldr	Kim et al. 2019
ai2_arc	Clark et al. 2018
quail	Rogers et al. 2020
crawl_domain	Zhang et al. 2020
glue-cola	Warstadt et al. 2019

Task Name	Reference
art	Bhagavatula et al. 2020
rotten_tomatoes	Pang and Lee 2005
tweet_eval-emoji	Barbieri et al. 2020
numer_sense	Lin et al. 2020a
blimp-existential_there_quantifiers_1	Warstadt et al. 2020
eli5-askh_qa	Fan et al. 2019
ethos-national_origin	Mollas et al. 2020
boolq	Clark et al. 2019
qa_srl	He et al. 2015
sms_spam	Almeida et al. 2011
samsun	Gliwa et al. 2019
ade_corpus_v2-classification	Gurulingappa et al. 2012
superglue-wic	Pilehvar and Camacho-Collados 2019
ade_corpus_v2-dosage	Gurulingappa et al. 2012
tweet_eval-stance_climate	Barbieri et al. 2020
e2e_nlg_cleaned	Dušek et al. 2020, 2019
aslg_pc12	Othman and Jemni 2012
ag_news	Gulli (link)
math_qa	Amini et al. 2019
commonsense_qa	Talmor et al. 2019
web_questions	Berant et al. 2013
biomrc	Pappas et al. 2020
swag	Zellers et al. 2018
blimp-determiner_noun_agreement_with_adj_irregular_1	Warstadt et al. 2020
glue-mnli	Williams et al. 2018
squad-with_context	Rajpurkar et al. 2016
blimp-ellipsis_n_bar_2	Warstadt et al. 2020
financial_phrasebank	Malo et al. 2014
sick	Marelli et al. 2014
ethos-religion	Mollas et al. 2020
hotpot_qa	Yang et al. 2018
tweet_eval-emotion	Barbieri et al. 2020
dbpedia_14	Lehmann et al. 2015
ethos-gender	Mollas et al. 2020
tweet_eval-hate	Barbieri et al. 2020
ethos-sexual_orientation	Mollas et al. 2020
health_fact	Kotonya and Toni 2020
common_gen	Lin et al. 2020b
crows_pairs	Nangia et al. 2020
ade_corpus_v2-effect	Gurulingappa et al. 2012
blimp-sentential_negation_npi_scope	Warstadt et al. 2020
lama-conceptnet	Petroni et al. 2019, 2020
glue-qnli	Rajpurkar et al. 2016
quartz-no_knowledge	Tafjord et al. 2019b
google_wellformed_query	Faruqui and Das 2018
kilt_trex	Elsahar et al. 2018
blimp-ellipsis_n_bar_1	Warstadt et al. 2020
trec	Li and Roth 2002; Hovy et al. 2001
superglue-copa	Gordon et al. 2012
ethos-disability	Mollas et al. 2020
lama-google_re	Petroni et al. 2019, 2020
discovery	Sileo et al. 2019
blimp-anaphor_number_agreement	Warstadt et al. 2020
climate_fever	Diggelmann et al. 2020
blimp-irregular_past_participle_adjectives	Warstadt et al. 2020
tab_fact	Chen et al. 2020a
gigaword	Napoles et al. 2012
glue-rte	Dagan et al. 2005
tweet_eval-sentiment	Barbieri et al. 2020
limit	Manotas et al. 2020
wikisql	Zhong et al. 2017
glue-qqp	(link)
onestop_english	Vajjala and Lučić 2018
amazon_polarity	McAuley and Leskovec 2013
blimp-wh_questions_object_gap	Warstadt et al. 2020
multi_news	Fabbri et al. 2019
proto_qa	Boratko et al. 2020
wiki_bio	Lebret et al. 2016
kilt_zsre	Levy et al. 2017
blimp-sentential_negation_npi_licensor_present	Warstadt et al. 2020

Table 10: References for all datasets.

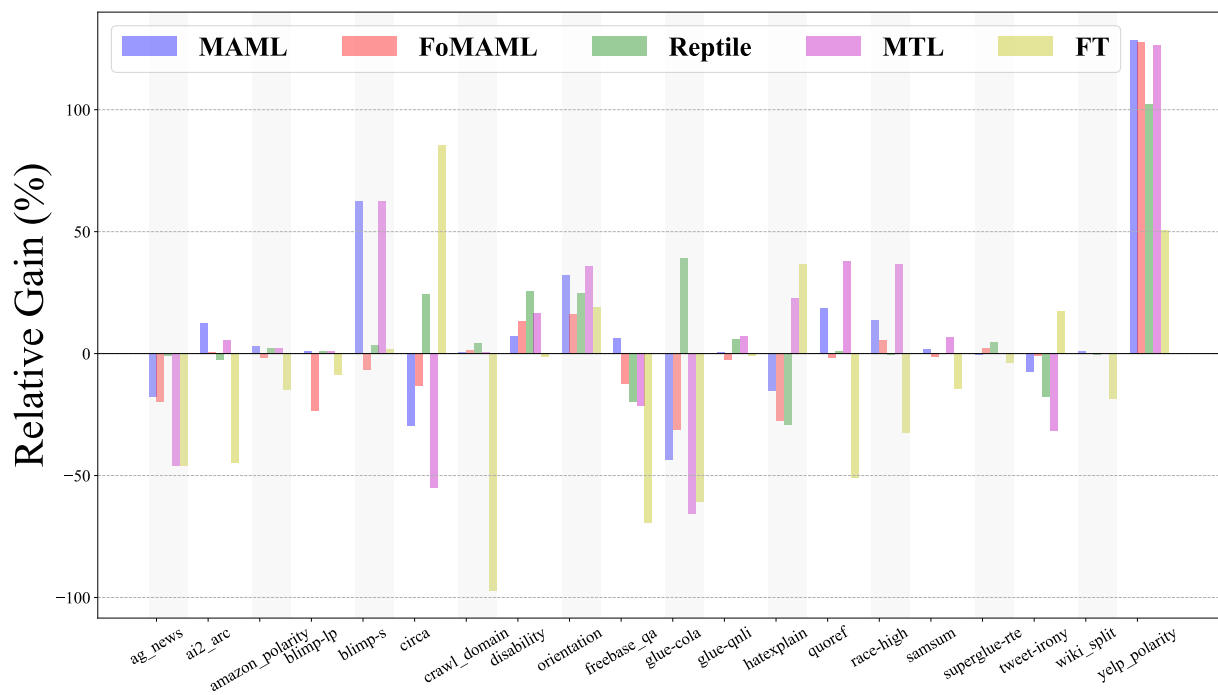


Figure 4: Random to Random

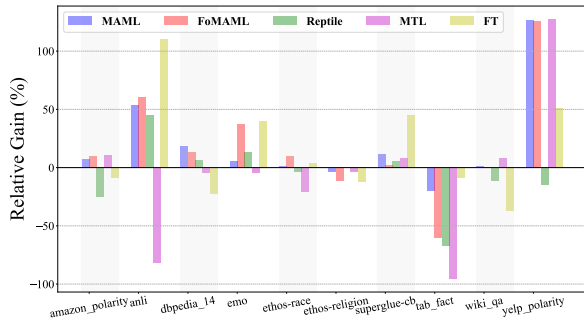


Figure 5: Classification to Classification

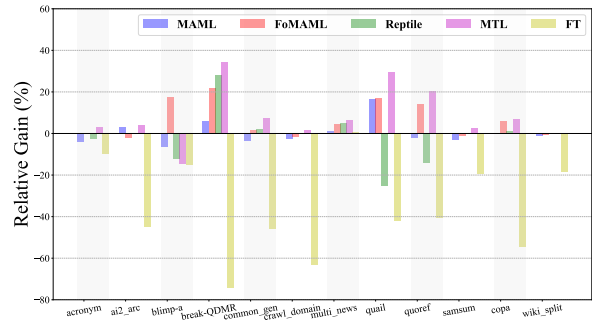


Figure 9: Classification to Non-Classification

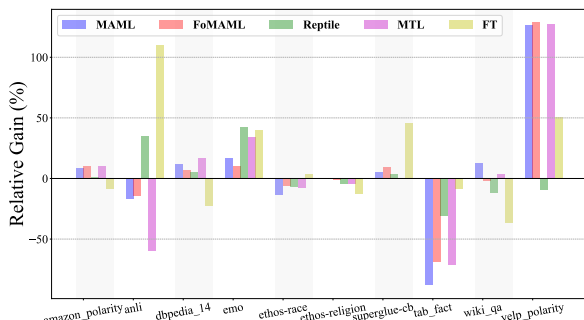


Figure 6: Non-Classification to Classification

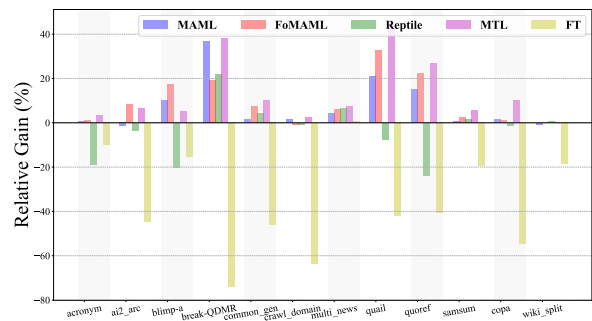


Figure 10: Both (Classification + Non-Classification) to Non-Classification

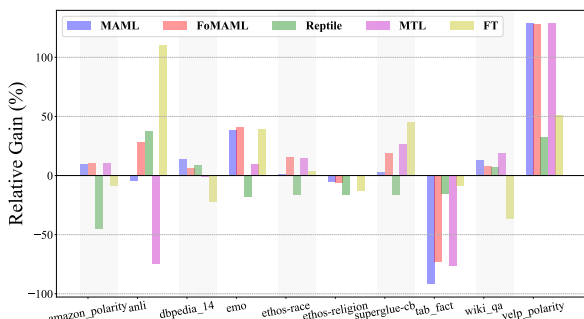


Figure 7: Both (Classification + Non-Classification) to Classification

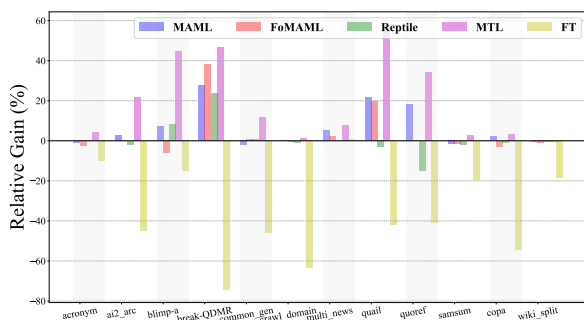


Figure 8: Non-Classification to Non-Classification

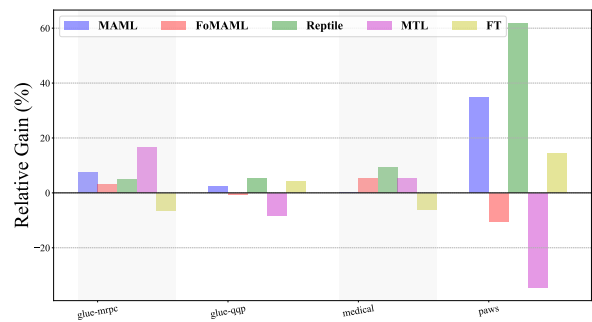


Figure 11: Non-Paraphrase Classification to Paraphrase

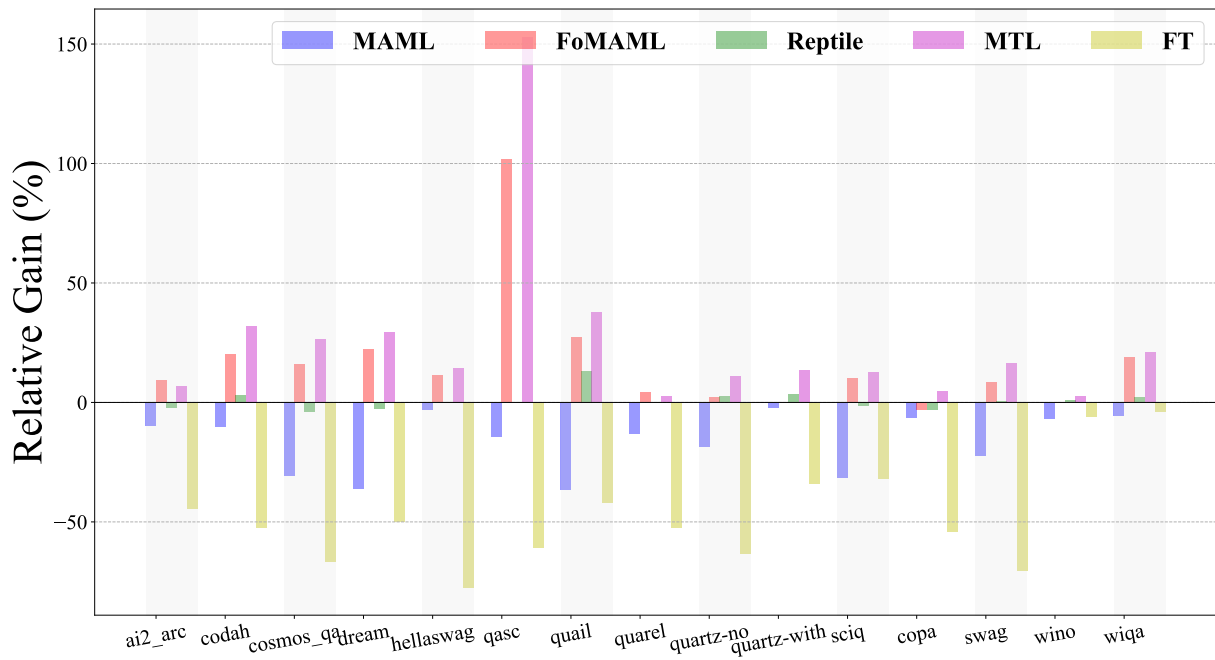


Figure 12: QA to QA

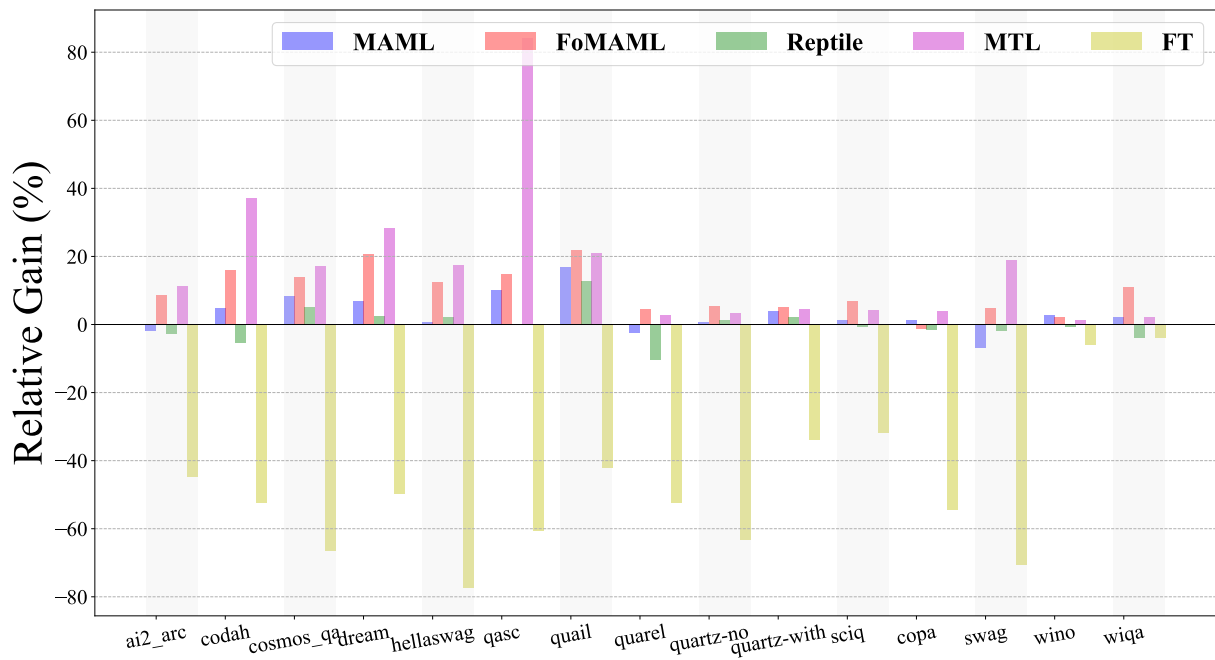


Figure 13: Non-QA to QA