


# CROSSFIT : A Few-shot Learning Challenge for Cross-task Generalization in NLP

Qinyuan Ye Bill Yuchen Lin Xiang Ren

University of Southern California

{qinyuany, yuchen.lin, xiangren}@usc.edu

## Abstract

Humans can learn a new language task efficiently with only few examples, by leveraging their knowledge obtained when learning prior tasks. In this paper, we explore whether and how such *cross-task generalization* ability can be acquired, and further applied to build better *few-shot learners* across diverse NLP tasks. We introduce CROSSFIT , a problem setup for studying cross-task generalization ability, which standardizes seen/unseen task partitions, data access during different learning stages, and the evaluation protocols. To instantiate different seen/unseen task partitions in CROSSFIT and facilitate in-depth analysis, we present the NLP Few-shot Gym, a repository of 160 diverse few-shot NLP tasks created from open-access NLP datasets and converted to a unified text-to-text format. Our analysis reveals that the few-shot learning ability on unseen tasks can be improved via an upstream learning stage using a set of seen tasks. We also observe that the selection of upstream learning tasks can significantly influence few-shot performance on unseen tasks, asking further analysis on task similarity and transferability.<sup>1</sup>

## 1 Introduction

Pre-trained language models fine-tuned with abundant task-specific data have become the predominant recipe for state-of-the-art results in NLP. However, these approaches are heavily dependent on large-scale labeled datasets that are expensive to create, and the resulting models still generalize poorly to out-of-distribution inputs created with small, harmless perturbations (Ribeiro et al., 2020). In retrospect, researchers have advocated for building more human-like, general linguistic intelligence that can “reuse previously acquired knowledge about a language and adapt to a new task quickly” (Yogatama et al., 2019; Linzen, 2020).

<sup>1</sup>Our code is at <https://github.com/INK-USC/CrossFit>.

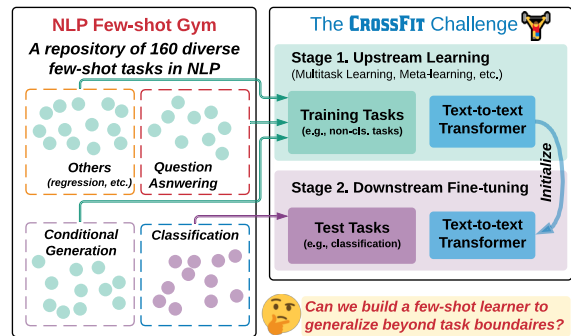


Figure 1: We present the CROSSFIT Challenge to study cross-task generalization in a diverse task distribution. To support this problem setting, we introduce the NLP Few-shot Gym, a repository of 160 diverse few-shot, text-to-text tasks in NLP.

Existing work has approached this problem via better few-shot fine-tuning, by re-formulating target tasks into cloze questions that resembles the pre-training objective (Schick and Schütze, 2020a,b), generating prompts and using demonstrations (Gao et al., 2020). Such progress primarily focus on improving *instance-level generalization*, *i.e.*, how to better generalize from few labeled instances to make predictions about new instances, *within the scope of one individual task*. From a broader perspective, human-like learning ability also benefits from *task-level generalization*, or *cross-task generalization*, *i.e.*, how to learn a new task efficiently given experiences of learning previous tasks.

Such ability has been widely studied in computer vision and robotics community (Yu et al., 2020; Triantafillou et al., 2020), but is relatively under-explored in NLP. Pruksachatkun et al. (2020) and Vu et al. (2020) study transferability between *one* intermediate task and a given target task, while it’s possible to further improve performance with *multiple* intermediate tasks. Han et al. (2018) and Bansal et al. (2020a) focus on cross-task generalization within the scope of classification tasks, whereas hu-

mans can generalize across different task formats (classification, multiple choice, generation, etc.), goals (question answering, fact checking, etc.) and domains (biomedical, social media, etc.).

Towards developing general linguistic intelligence, we present CROSSFIT, a few-shot learning challenge to acquire, evaluate and analyze cross-task generalization in a realistic setting, with standardized training pipeline, data access and evaluation protocol. The CROSSFIT challenge requires a model to first learn from a set of seen tasks in an upstream learning stage, and then perform few-shot learning on a set of unseen tasks, as illustrated in Fig. 1. In accompany, we introduce the NLP Few-shot Gym, a repository of 160 few-shot NLP tasks gathered from open-access resources, covering a wide range of capabilities and goals, and formulated into a unified text-to-text format. To analyze the capability and limitation of existing approaches to the CROSSFIT challenge, we design eight specific seen/unseen task partitions.

With the CROSSFIT Challenge and the NLP Few-shot Gym, we aim to investigate the following research questions:

- **Q1.** Can we teach cross-task generalization ability to pre-trained models with existing methods?
- **Q2.** During upstream learning, is it better to be “well-rounded” (learning from diverse tasks) or be “specialized and targeted” (learning from tasks in the same category with unseen tasks)?
- **Q3.** Does it help if we have more labelled data for seen tasks during upstream learning?

To address the above questions, we empirically analyze the performance of multi-task learning and three meta-learning algorithms (MAML (Finn et al., 2017), first-order MAML and Reptile (Nichol et al., 2018)). We observe that these approaches can indeed lead to better few-shot performance on unseen tasks. Interestingly, simple multi-task learning outperforms existing meta-learning methods in many cases, encouraging future research on identifying the reasons and developing improved meta-learning methods. For Q2, we observe that performance of individual unseen tasks varies with different selection of seen tasks, calling for more thorough investigation of the relationship between task similarity and transferability. As for Q3, we find that enlarging the size of upstream data does not necessitate better cross-task generalization abilities. We envision cross-task generalization to be an integral component towards general linguistic

intelligence, and we hope CROSSFIT serves as a useful testbed for driving related progress.

## 2 Related Work

**Few-shot Fine-tuning.** Few-shot learning refers to teaching models a new task with a small number of annotated examples. Large-scale pre-trained language models (e.g., BERT (Devlin et al., 2019)) have demonstrated great ability to learn new tasks efficiently via *fine-tuning* (Zhang et al., 2021). Schick and Schütze (2020a,b) proposed *pattern-exploiting training* (PET), which formulates text classification and NLI tasks into cloze questions (or “prompts”) that resemble masked language modeling. PET can be further improved by generating prompts *automatically* and incorporating demonstrations into the input (Gao et al., 2020); and by densifying the supervision signal with label conditioning (Tam et al., 2021). While successful, in these approaches the downstream tasks are learned in isolation. Our work aims to boost few-shot learning ability on unseen tasks via acquiring cross-task generalization ability from diverse seen tasks.

**Meta-learning in NLP.** Recent works have explored meta-learning methods for relation classification (Han et al., 2018; Gao et al., 2019), general text classification (Dou et al., 2019; Bansal et al., 2020a,b), low-resource machine translation (Gu et al., 2018), cross-lingual NLI/QA (Nooralahzadeh et al., 2020). In general, these works apply meta-learning algorithms to a set of sub-tasks; however the sub-tasks are either *synthetic* (e.g., classifying a new set of five relations is a new sub-task) or drawn from a rather *narrow* distribution (e.g., QA in one language is a sub-task). In our work, we explore a more realistic setting – learning from a set of NLP tasks with *diverse* goals: classification, question answering, conditional generation, etc. This setting is attracting attention in NLP community rapidly and is also explored in very recent work (Zhong et al., 2021; Mishra et al., 2021; Bragg et al., 2021; Wei et al., 2021).

**Unifying NLP Task Formats.** Researchers have explored unifying the formats of different tasks, in order to better enable knowledge transfer, e.g., DecaNLP (McCann et al., 2018), UFO-Entail (Yin et al., 2020) and EFL (Wang et al., 2021). Following T5 (Raffel et al., 2020), we adopt a unified text-to-text format that subsumes all text-based tasks of interest. Related to our work, UnifiedQA

(Khashabi et al., 2020) examines the feasibility of training a general cross-format QA model with multi-task learning. Our work extends from these ideas, and we significantly enlarge the task repository to 160 to broaden the coverage, in hopes to build a general-purpose few-shot learner.

### 3 The CROSSFIT Challenge

In this section, we present the CROSSFIT Challenge, a problem setting for acquiring and evaluating cross-task generalization. Ideally, a strong CROSSFIT system can capture cross-task generalization ability from a set of seen tasks and thus adapts to new unseen tasks efficiently.

#### 3.1 Preliminaries

The meaning of “task” is overloaded: “tasks” can be categorized at different granularity (e.g., text classification vs. QA, yes/no QA vs. machine reading comprehension), and from different aspects (e.g., domain, label space). Herein we take a general formulation by defining a “task” with its training and testing examples. We define a task  $T$  as a tuple of  $(\mathcal{D}_{train}, \mathcal{D}_{dev}, \mathcal{D}_{test})$ . Each set  $\mathcal{D}$  is a set of annotated examples  $\{(x_i, y_i)\}$  in text-to-text format. In few-shot setting, the size of  $\mathcal{D}_{train}$  and  $\mathcal{D}_{dev}$  are required to be small (e.g., 16 example per class for classification tasks).

Existing work mostly focuses on improving *instance-level* generalization for individual task by using *task-specific* templates. Performance on individual tasks is used as the measure of success. For the CROSSFIT Challenge, we aim to acquire *cross-task generalization* and build better *general-purpose* few-shot learners, which calls for a different problem setting with distinct training procedure and evaluation protocol.

#### 3.2 Problem Setting

**Tasks and Data.** To acquire and evaluate cross-task generalization, we first gather a large repository of few-shot tasks  $\mathcal{T}$ , and partition them into three non-overlapping sets  $\mathcal{T}_{train}, \mathcal{T}_{dev}, \mathcal{T}_{test}$ . In hopes to examine the capability and limitation of an approach in different settings, and to answer our research questions, we design multiple task partitions with different focuses. Details of the repository and partitions, or as we name them, the NLP Few-shot Gym, are deferred to §4.

**Learning Stages.** A CROSSFIT method may learn from  $\mathcal{T}_{train}$  and perform necessary tuning

with  $\mathcal{T}_{dev}$  in the upstream learning stage; it is then evaluated with few-shot tasks in  $\mathcal{T}_{test}$ :

- **Upstream learning stage.** Here, the algorithm has access to the  $\mathcal{D}_{train}$  and  $\mathcal{D}_{dev}$  for each training task in  $\mathcal{T}_{train}$ , while  $\mathcal{D}_{test}$  is unavailable. The algorithm also has access to all data in  $\mathcal{T}_{dev}$ , but for validation purpose only (i.e., it is not allowed to use  $\mathcal{T}_{dev}$  to update model weights).
- **Few-shot learning stage.** In this stage,  $\mathcal{T}_{test}$  became available. Models resulting from the upstream learning stage are required to learn from  $\mathcal{D}_{train}$  via a particular few-shot learning method (e.g., direct fine-tuning). The final few-shot learning performance is evaluated on  $\mathcal{D}_{test}$ .<sup>2</sup>

**Evaluation Metric.** Evaluating the performance of a model on a diverse collection of NLP tasks is inherently challenging, as different tasks use different metrics. It is thus not reasonable to simply aggregate performance of classification tasks (e.g., accuracy, F1) and generation tasks (e.g., ROUGE, BLEU) by taking the average.

To address this problem, we first narrow down to a collection of 7 evaluation metrics: classification F1, accuracy, QA F1, exact match (EM), Rouge-L, Matthew correlation, and Pearson correlation, which cover all tasks in our experiments. Then, we define *Average Relative Gain (ARG)*, a metric that computes relative performance changes before and after the *upstream learning stage* for each test task, and finally take the average across all test tasks.

For example, suppose we have  $\mathcal{T}_{test} = \{T_A, T_B\}$ . If an upstream learning algorithm helps improve the few-shot learning performance from 50% F1 score to 70% on task  $T_A$  (i.e., a 40% relative improvement), and from 40% accuracy to 30% on task  $T_B$  (i.e., -25% relative improvement), the final ARG on  $\mathcal{T}_{test}$  would be computed as  $\frac{40\% + (-25\%)}{2} = 7.5\%$ .

The ARG metric reflects the *overall* performance gain on all tasks in  $\mathcal{T}_{test}$ , no matter what specific metrics each task uses. We use ARG for a high-level comparison, and we still analyze the performance for each task (e.g., absolute performance metrics, performance growth with “more shots”, sensitivity to different selection of  $\mathcal{T}_{train}$ ) in our in-depth analysis.

<sup>2</sup>For clarification, the performance on the  $\mathcal{D}_{dev}$  of a task in  $\mathcal{T}_{dev}$  or  $\mathcal{T}_{test}$  will be used for tuning hyper-parameters during fine-tuning. The overall performance on  $\mathcal{T}_{dev}$  is used for tuning tuning hyper-parameters during upstream learning.

## 4 NLP Few-shot Gym

Towards learning to generalize across tasks in CROSSFIT challenge, we need a resource that contains sufficient number of tasks, covering a wide range of NLP applications, and presented in a unified text-to-text format. Herein, we introduce the NLP Few-shot Gym, a repository of 160 few-shot tasks gathered from existing open-access datasets.

### 4.1 Dataset Selection

We choose to use Huggingface Datasets<sup>3</sup> (Lhoest et al., 2021) as the pool of our candidate tasks. We filter these datasets on a case-by-case basis, mainly using the following criteria: (1) We focus on English monolingual datasets. (2) We exclude datasets that require information retrieval, as they require a separate retriever. (3) We exclude sequence labeling tasks (e.g., dependency parsing, NER), which are highly dependent on tokenization, and are hard to evaluate in text-to-text format. (4) We exclude datasets dealing with extremely long documents (e.g., a scientific paper) as input, as most pre-trained models cannot process such long input sequences. We finalize our selection with 160 datasets which are detailed in Appendix A.

### 4.2 A Unified Text-to-Text Format

We follow Raffel et al. (2020) and convert all of our datasets into a unified text-to-text format. For example, the task of natural language inference (originally a sentence-pair classification problem) becomes: premise: <premise> hypothesis: <hypothesis>, and the target sequence is either the word entailment, contradiction or neutral. As for machine reading comprehension tasks, the input format is question: <question> context: <context> and the target sequence is the correct answer span. We also reference the format for QA tasks from UnifiedQA (Khashabi et al., 2020).

### 4.3 Formulating Few-shot Tasks

We mainly follow the practice in (Gao et al., 2020) for few-shot sampling. For classification and regression tasks, we include 16 training examples *per class* in  $D_{train}$ . For other types of tasks, we include 32 examples in  $D_{train}$ . In conformity with real-world situations where labeled data are scarce,

<sup>3</sup><https://huggingface.co/datasets>. It is an extensible library that provides access to 626 open-access NLP datasets (as of Feb 25th, 2021) with a unified, open-source API.



Figure 2: Task Ontology for the NLP Few-shot Gym. Full information is listed in Appendix A.

we assume a development set  $D_{dev}$  which shares the same size with  $D_{train}$ .

We sample  $D_{train}$  and  $D_{dev}$  splits from each dataset’s original train set with 5 different random seeds. This helps us reduce variance during few-shot evaluation, and also enlarges the number of few-shot tasks used for learning. Consequently, the “effective size” of our NLP Few-shot Gym is  $160 \times 5 = 800$ , while we use the number 160 throughout the paper to avoid possible confusion.

We use the original development set for each dataset as  $D_{test}$ , or withhold 20% of the dataset when the official development split is not available. The held-out test examples are sampled *once* before sampling  $D_{train}$  and  $D_{dev}$ .

### 4.4 Task Ontology and Partitions

As mentioned in §3.2, a CROSSFIT method is expected to first acquire cross-task generalization on a set of  $\mathcal{T}_{train}$  and evaluate such ability on  $\mathcal{T}_{test}$ . To comprehensively analyze to what extent a trained model can generalize, and how its behavior differs in different scenarios, we need to build different partitions of  $(\mathcal{T}_{train}, \mathcal{T}_{dev}, \mathcal{T}_{test})$ .

Towards this goal, we first manually classify the 160 tasks and form a **task ontology** with categories and sub-categories, as shown in Fig. 2. The first-level categories include classification, question answering, conditional generation, and oth-

ers.<sup>4</sup> Further, we design eight different partitions of  $(\mathcal{T}_{train}, \mathcal{T}_{dev}, \mathcal{T}_{test})$ . We illustrate four partitions in Fig. 3 and provide more details in Table 1.

Our Partition 1 randomly split all 160 few-shot tasks into the three sets, where  $|\mathcal{T}_{train}| = 120$  and  $|\mathcal{T}_{dev}| = |\mathcal{T}_{test}| = 20$ . The design of Partition 1 mimics the real-world language learning environment where the goal is to build a general-purpose few-shot learner, and a set of diverse tasks ( $\mathcal{T}_{train}$ ) are used to train the learner. Our Partition 2.1-2.3 withhold 10 classification tasks for development and 10 more for testing. The  $\mathcal{T}_{train}$  is controlled to have either 100% classification tasks, 100% non-classification tasks, or half-and-half. These three partitions help us to understand the influence brought by different task distribution in  $\mathcal{T}_{train}$ . The remaining four partitions still focus on crossing task boundaries, but in a finer granularity: seen and unseen tasks are in the same category, but not the same sub-category. For example, Partition 3.1 has 57 non-NLI classification tasks as  $\mathcal{T}_{train}$ , and 8 NLI tasks as  $\mathcal{T}_{test}$ . These partitions help us to understand whether cross-task generalization in this finer granularity is easier for model to acquire.

## 5 Methods to CROSSFIT

We mainly use BART-Base (Lewis et al., 2020) as the text-to-text transformer for our analysis in the CROSSFIT setup. We leave confirmatory experiments with T5-v1.1-Base and BART-Large model in Appendix C.

**Direct Fine-tuning on Test Tasks.** This serves as the basic baseline method for the CROSSFIT challenge, which does not make use of  $\mathcal{T}_{train}$  or  $\mathcal{T}_{dev}$ , or go through the upstream learning stage. For each task  $T \in \mathcal{T}_{test}$ , we directly fine-tune the text-to-text model with its  $\mathcal{D}_{train}$ , tune the hyperparameters with  $\mathcal{D}_{dev}$ , and assess its performance with the test set  $\mathcal{D}_{test}$ . We use the performance of direct fine-tuning as the base for computing ARG scores of other CROSSFIT approaches. We expect a model trained with upstream learning would capture cross-task generalization ability and thus have better ARG scores.

**Multi-task Learning (MTL).** A straightforward yet effective method is to combine the data<sup>5</sup> in the training tasks to learn a multi-task

<sup>4</sup>We later discuss the limitation of this design in §6-Q2

<sup>5</sup>Both  $\mathcal{D}_{train}$  and  $\mathcal{D}_{dev}$  are used, as  $\mathcal{D}_{dev}$  is used for gradient updates in meta-learning algorithm. We do so to make sure that the data access for the two methods is fair.

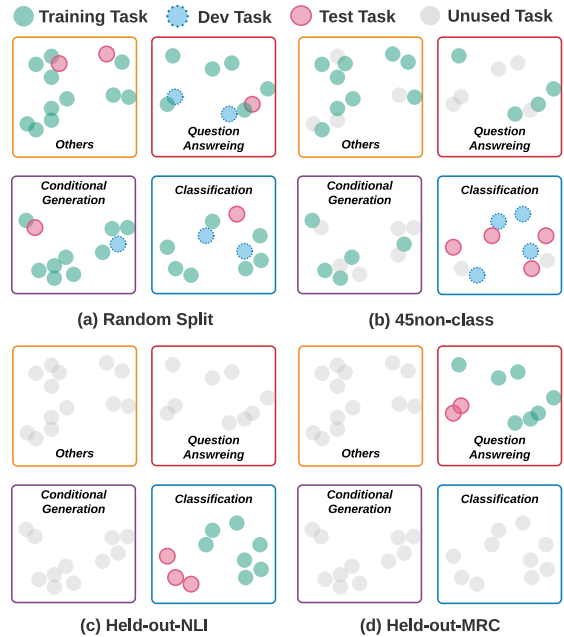


Figure 3: **Illustration for different task partitions.** We evaluate a CROSSFIT approach on different task partitions to examine its generalization ability in different scenarios. Full details in Table 1. The locations and distances in this figure are hypothetical and for illustrative purposes only.

model, before fine-tuning it on each test task. Specifically, we gather source-target examples for all tasks in  $\mathcal{T}_{train}$  and fine-tune the text-to-text model with these examples. Then we use the resulting checkpoint as initialization and perform the same procedure in “direct fine-tuning” for each test task in  $\mathcal{T}_{test}$ . The performance gain over the *direct fine-tuning* is used for computing its overall ARG score.

**Model-Agnostic Meta-learning (MAML).** Cross-task generalization ability, closely aligns with the concept of learning to learn. Hence, we use MAML (Finn et al., 2017), a representative meta-learning approach during upstream learning. The core concept of MAML is to learn a set of initialization weight, from which the model adapts fast to a new task within few gradient updates. In MAML training, we iterate through tasks in  $\mathcal{T}_{train}$  to update the model. For each train task  $(\mathcal{D}_{train}, \mathcal{D}_{dev})$ , we first sample a support batch  $\mathcal{B}_{support}$  from  $\mathcal{D}_{train}$  and a query batch  $\mathcal{B}_{query}$  from  $\mathcal{D}_{dev}$ . We use  $f_{\theta}$  to denote the text-to-text model with parameters  $\theta$ . Using  $\mathcal{B}_{support}$ , we first compute the updated parameters  $\theta'$  with gradient descent (*i.e.*, the inner loop). Due to the large size of pre-trained text-to-text models, we

No.	Shorthand	$\mathcal{T}_{train}$	$\mathcal{T}_{dev}$	$\mathcal{T}_{test}$	ARG(Multi)	ARG(MAML)	ARG(FoMAML)	ARG(Rept.)	Details
1	Random	120	20	20	35.06%	28.50%	22.69%	25.90%	Fig. 4(a)
2.1	45cls	45 cls.	10 cls.	10 cls.	11.68%	9.37%	10.28%	13.36%	Fig. 5
2.2	23cls+22non-cl	23 cls. + 22 non-cl.	10 cls.	10 cls.	11.82%	9.69%	13.75%	14.34%	
2.3	45non-cl	45 non-cl.	10 cls.	10 cls.	11.91%	9.33%	11.20%	14.14%	
3.1	Held-out-NLI	57 non-NLI cls.	/	8 NLI	16.94%	12.30%	12.33%	14.46%	Fig. 4(b)
3.2	Held-out-Para	61 non-Paraphrase cls.	/	4 Para. Iden.	18.21%	17.90%	21.57%	19.72%	Fig. 4(c)
4.1	Held-out-MRC	42 non-MRC QA	/	9 MRC	32.81%	27.28%	28.85%	28.85%	Fig. 4(d)
4.2	Held-out-MCQA	29 non-MC QA	/	22 MC QA	12.20%	4.69%	6.73%	7.67%	Fig. 4(e)

Table 1: ( $\mathcal{T}_{train}, \mathcal{T}_{dev}, \mathcal{T}_{test}$ ) partitions used in the study (full lists in Appendix B), and their ARG scores when upstream learning methods are applied. “cls.” stands for “classification”, “Para. Iden.” for “paraphrase identification”, “MRC” for “machine reading comprehension” and “MCQA” for “multiple-choice QA”.


use one gradient update in the inner loop, *i.e.*,  $\theta' = \theta - \alpha \nabla_{\theta} \mathcal{L}(f_{\theta}, \mathcal{B}_{support})$ . Then we apply the updated text-to-text model  $f_{\theta'}$  to  $\mathcal{B}_{query}$ , and do one step of meta-optimization (*i.e.*, the outer loop), with  $\theta \leftarrow \theta - \beta \nabla_{\theta} \mathcal{L}(f_{\theta'}, \mathcal{B}_{query})$ .

**First-order MAML.** First-order MAML (Finn et al., 2017) avoids second-order optimization and improves training stability using the first-order approximation by differentiating with respect to the fast weights  $\theta'$  instead of the original parameters  $\theta$  for the gradient  $\nabla_{\theta} \mathcal{L}(f_{\theta'}, \mathcal{B}_{query})$ , *i.e.*,  $\theta \leftarrow \theta - \beta \nabla_{\theta'} \mathcal{L}(f_{\theta'}, \mathcal{B}_{query})$ .

**Reptile.** Reptile (Nichol et al., 2018) is another memory-efficient, first-order meta-learning algorithm that first makes multiple gradient updates in the inner loop, then directly uses  $\theta' - \theta$  to approximate  $\nabla_{\theta} \mathcal{L}(f_{\theta'}, \mathcal{B}_{query})$ , *i.e.*,  $\theta \leftarrow \theta + \beta(\theta' - \theta)$ .

## 6 Empirical Analysis

In this section we look to interpret the results and answer our research questions. We summarize the ARG scores in Table 1 and plot the performance of each test task (for each partition) in Fig. 4-5.

 Q1. Can we teach pre-trained LMs to generalize across tasks with existing methods?

**Overall Performance.** From Table 1, we observe that, on average, the tested upstream learning methods indeed improve cross-task generalization: their ARG scores are positive, meaning that they are better than *direct fine-tuning* (ARG=0%). Further, by aggregating results from all upstream learning methods and task partitions, we find that the performance on 51.47% test tasks are significantly improved ( $> 5\%$  relative improvement compared to direct fine-tuning); 35.93% tasks are relatively unaffected (between  $\pm 5\%$ ); and 12.60% tasks suffer from worse performance ( $< -5\%$ ).

**Correlated Performance Gains.** The performance gain obtained with different upstream learning methods are correlated with each other – *i.e.*, tasks that benefit from multi-task learning is likely to also benefit from meta-learning. For the *Random partition*, the Spearman Correlation between the relative improvement brought by MTL and MAML is 0.66, with  $p$  value equals to 0.0015. This suggests that different upstream learning methods, while taking different optimization objectives, capture similar inductive bias from  $\mathcal{T}_{train}$ .

**MTL is a strong baseline.** Surprisingly, the most straight-forward multi-task learning method is hard to beat. This could be counter-intuitive, as meta-learning methods are specifically designed for rapid generalization to unseen tasks, sharing the same goal with our CROSSFIT challenge. We think there are three possible reasons: (1) Due to memory constraints, we limit the number of inner-loop updates to be one, which may be insufficient. Also, meta-learning methods are highly sensitive to hyper-parameters and even random seeds (Antoniou et al., 2019), which we do not tune exhaustively for practical reasons. (2) Text-to-text transformers have much more complex architectures, while most meta-learning methods are typically applied to small feed-forward/convolutional networks. (3) The CROSSFIT challenge has a highly diverse set upstream tasks, which may introduce under-explored difficulties. That being said, we believe it is important to identify the true cause, and to develop improved meta-learning methods for the CROSSFIT challenge as future work.

**Forgetting Pre-Trained Knowledge.** A few test tasks have negative performance gain after upstream learning, including Glue-COLA (measuring linguistic acceptability) and Domain Crawl (separating domain names into tokens) in the Random



Figure 4: Experimental results for the CROSSFIT challenge with different task partitions. The details of each partition is shown in Table 1. Relative performance gain is computed based on the results of *direct fine-tuning*. Best viewed in color. Green color is used to highlight the Average Relative Gain (ARG) for each method.

Partition setting. For Glue-COLA, similar observations are reported by Pruksachatkun et al. (2020) in an intermediate-task transfer learning setting, where the authors conjecture *catastrophic forgetting* of the masked language modeling (MLM) tasks may be the cause. BART uses denoising pre-training objective, a variant of MLM. Intuitively, Domain Crawl is also one of the most similar tasks to denoising in all test tasks, which further supports this hypothesis. We thus conjecture that for test tasks that resemble pre-training objectives, upstream learning could hurt performance due to the *catastrophic forgetting* phenomena.

Understanding negative transfer (Wu et al., 2020) and selecting source tasks to avoid negative transfer (Vu et al., 2020) are also growing research topics. In this work we refrain from further investigation; however we believe combating negative transfer and thus improving CROSSFIT performance is a promising future direction.

💡 Q2. Well-rounded or specialized? Which is a better strategy of upstream learning?

“Learning to be well-rounded vs. learning to be specialized” is a common dilemma that human learners struggles with. For the CROSSFIT challenge, the former refers to learning from a set of diverse tasks in upstream learning; the latter refers to learning from a set of tasks closer to target few-shot tasks. To study this research question, we want to find out which option works better in upstream learning. Put differently, we aim to **analyze the influence of upstream task selection** for a fixed set of the downstream tasks.

**Setup.** We first conduct controlled experiments with *Partition 2.1-2.3*, where  $\mathcal{T}_{test}$  is a fixed set of classification tasks, and  $\mathcal{T}_{train}$  varies. In Partition 2.1, all tasks in  $\mathcal{T}_{train}$  are *classification* tasks (i.e., “specialized and targeted”); in Partition

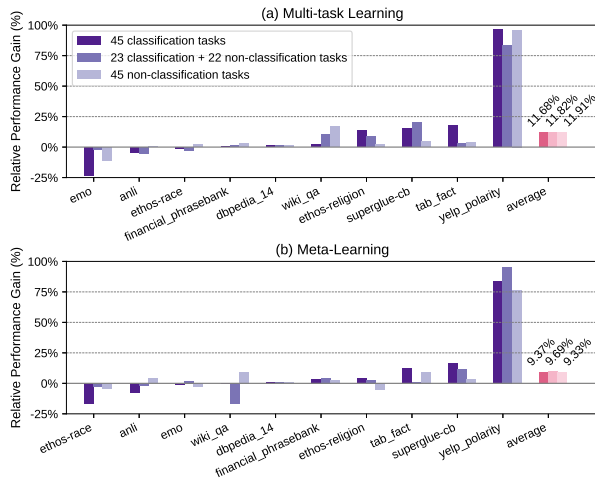


Figure 5: Comparison for the controlled experiment on Partition 2.1-2.3.  $\mathcal{T}_{test}$  is a fixed set of 10 classification tasks, while  $\mathcal{T}_{train}$  varies.

2.2, half of the tasks are *classification* tasks (i.e., “well-rounded”); in Partition 2.3, all tasks are *non-classification* tasks (i.e., “specialized in an opposite direction”, for a controlled experiment).

**Analysis and Discussion.** It is surprising at first that non-classification tasks and classification tasks are equivalently helpful in terms of ARG scores (see Fig. 5). On a second thought, this observation is encouraging as it demonstrates that acquiring cross-task generalization is feasible and promising, even when  $\mathcal{T}_{train}$  and  $\mathcal{T}_{test}$  are drastically different. It also suggests that our categorization of tasks (§4.4) may not align with how models learn transferable skills: selecting  $\mathcal{T}_{train}$  tasks that have the same format and goal as the test task may not lead to optimal transfer.

In retrospect, we acknowledge that our design of ontology and partitions based on task format and goal is flawed. This is merely one aspect of “task similarity”. However, understanding the complex relationship between tasks is another challenging and under-explored problem. We consider our ontology as a starting point, rather than a fixed final one. We use the current ontology to guide our experiment and analysis, and we hope future analysis could help build a more informative ontology.

**Case Studies.** We further look at cases where a test task appear in  $\mathcal{T}_{test}$  of multiple partitions. For example, AI2\_ARC and Race-High are in the  $\mathcal{T}_{test}$  of both Random partition and Held-out-MCQA partition. We present the results in Table 2. In general, the performance of these tasks varies when

Test Task	Partition	$\Delta_{multi}$	$\Delta_{meta}$
Glue-QNLI	Random	15.89%	11.55%
	Held-Out-NLI	10.88%	10.94%
AI2_ARC	Random	1.30%	4.22%
	Held-Out-MCQA	6.49%	-6.22%
Race-High	Random	26.71%	6.59%
	Held-Out-MCQA	7.27%	-6.28%
QuoRef	Random	25.47%	3.99%
	Held-Out-MRC	12.25%	4.64%

Table 2: Performance comparison of test task performance when different  $\mathcal{T}_{train}$  sets are used in upstream learning. See text in Q2 for in-depth analysis.

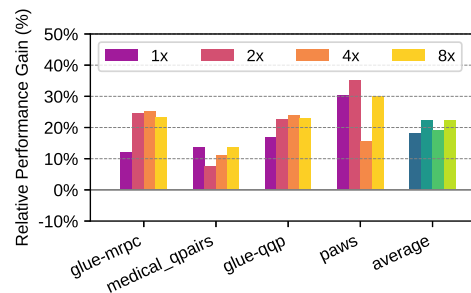


Figure 6: Controlling upstream learning data size in with Held-out-Para Partition. Enlarging the size of data during upstream learning *does not* necessitate better cross-task generalization ability.

different  $\mathcal{T}_{train}$  sets are used. However, we have not found consistent patterns of what type of  $\mathcal{T}_{train}$  lead to better performance for a specific test task.

💡 Q3. Does it help if we have more labelled data for upstream tasks?

As described in §4.3, we limit our upstream tasks to be also few-shot: classification tasks have 16 examples per class, and non-classification tasks have 32 examples. This decision is empirically determined following prior works (Schick and Schütze, 2020a,b; Gao et al., 2020) and makes our extensive analysis practical and efficient. It is possible that using more data for each upstream task can significantly improve cross-task generalization. To investigate this, we conduct a set of controlled experiments where the number of examples in upstream tasks are changed to [2, 4, 8] times of the original size. We use the Held-out-Para Partition and multi-task learning for the experiments, and present the result in Fig. 6. Surprisingly, we find that the effect from using more upstream data is inconsistent on different target tasks. The overall ARG for all sizes are close: even 8x larger up-



stream data leads to only 4% improvement in ARG. We conclude that enlarging the size of data during upstream learning *does not* necessitate better cross-task generalization ability. This also justifies our decision to keep upstream tasks few-shot.

#### Q4-Q6. Additional Analysis

Due to space limit, we summarize our other findings below and defer the details to Appendix C.

**Few-Shot → More-Shot (Q4).** In practice, users may continue to collect data over time. We wonder if cross-task generalization ability is still helpful for medium/high-resource target tasks. We find that the performance gain from upstream learning is still evident when 1024 shots are available. The performance gap diminishes with millions of training examples.

**Using Different Base Models (Q5).** We extend our analysis on BART-base (139M) to larger pre-trained text-to-text Transformers: BART-Large (406M) and T5-v1.1-Base (248M). Generally, the performance grows with models sizes with only few exceptions, which suggests that upstream learning methods we use are model-agnostic, and can be applied to larger models to further improve few-shot performance.

**Integration with PET Training (Q6).** Pattern-exploiting training (PET) (Schick and Schütze, 2020a,b) was originally proposed for classification tasks and *encoder* language models. We test a few variants of PET training with BART-Base and try applying PET training after upstream learning. In general we observe deteriorated performance compared to direct fine-tuning. We hypothesize that PET methods are not directly applicable to *encoder-decoder* language models used in our study.

## 7 Conclusion and Future Work

In this paper, we study the problem of building better few-shot learners via acquiring cross-task generalization ability from diverse NLP tasks. Towards our goal, we introduce the CROSSFIT Challenge, an task setup that standardizes the training pipeline, data access and evaluation protocol. We also present the NLP Few-shot Gym, a repository of 160 diverse few-shot NLP tasks, to support CROSSFIT learning in different scenarios. We empirically demonstrated that cross-task generalization can be acquired via multi-task learning and

meta-learning; confirmed that the selection of seen tasks would influence the few-shot performance on unseen tasks.

We have highlighted several unexpected or undesired observations in our analysis, for which we invite future work in understanding and combating related issues. In addition, we envision the CROSSFIT Challenge and the NLP Few-shot Gym to serve as the testbed for many interesting “meta-problems”, such as (1) learning to generate prompt for diverse task formats and further improve learning efficiency (Shin et al., 2020; Gao et al., 2020); (2) learning to select appropriate source tasks to learn from during upstream learning (Zamir et al., 2018; Standley et al., 2020), potentially with task2vec methods (Achille et al., 2019; Vu et al., 2020); (3) applying task augmentation strategies to prevent over-fitting (Murty et al., 2021); (4) learning to accumulate knowledge and avoid catastrophic forgetting in an continual learning setup (Jin et al., 2021); (5) decomposing complex tasks into atomic tasks and exploring cross-task generalization through the lens of compositionality (Andreas et al., 2016; Khot et al., 2021).

## Acknowledgments

We thank authors and crowd-workers of all datasets used in our study. We thank huggingface datasets team for making datasets more accessible. We thank anonymous reviewers and members of USC INK Lab for their valuable feedback. This work is supported in part by the Office of the Director of National Intelligence (ODNI), Intelligence Advanced Research Projects Activity (IARPA), via Contract No. 2019-19051600007; the DARPA MCS program under Contract No. N660011924033; the Defense Advanced Research Projects Agency with award W911NF-19-20271; NSF IIS 2048211.

## References

- A. Achille, Michael Lam, Rahul Tewari, A. Ravichandran, Subhansu Maji, Charless C. Fowlkes, Stefano Soatto, and P. Perona. 2019. Task2vec: Task embedding for meta-learning. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6429–6438.
- 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.
- Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. 2016. [Learning to compose neural networks for question answering](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1545–1554, San Diego, California. Association for Computational Linguistics.
- Antreas Antoniou, Harrison Edwards, and Amos Storkey. 2019. [How to train your MAML](#). In *International Conference on Learning Representations*.
- Trapit Bansal, Rishikesh Jha, and Andrew McCallum. 2020a. [Learning to few-shot learn across diverse natural language classification tasks](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5108–5123, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Trapit Bansal, Rishikesh Jha, Tsendsuren Munkhdalai, and Andrew McCallum. 2020b. [Self-supervised meta-learning for few-shot natural language classification tasks](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 522–534, Online. Association for Computational Linguistics.
- Roy Bar-Haim, Ido Dagan, Bill Dolan, Lisa Ferro, Danilo Giampiccolo, Bernardo Magnini, and Idan Szpektor. 2006. The second pascal recognizing textual entailment challenge. In *Proceedings of the second PASCAL challenges workshop on recognising textual entailment*, volume 6, pages 6–4. Venice.
- 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.
- Luisa Bentivogli, Peter Clark, Ido Dagan, and Danilo Giampiccolo. 2009. The fifth pascal recognizing textual entailment challenge. In *TAC*.
- 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 Le Bras, Jianfeng Gao, and Yejin Choi. 2020. [Piqa: Reasoning about physical commonsense in natural language](#). In *Thirty-Fourth AAAI Conference on Artificial Intelligence*.
- 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.
- Jonathan Bragg, Arman Cohan, Kyle Lo, and Iz Beltagy. 2021. [Flex: Unifying evaluation for few-shot nlp](#).
- 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*.
- 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.
- Anthony Chen, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2020a. [MOCHA: A dataset for training and evaluating generative reading comprehension metrics](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6521–6532, Online. 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. 2020b. [Tabfact: A large-scale dataset for table-based fact verification](#). In *International Conference on Learning Representations*.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. [BoolQ: Exploring the surprising difficulty of natural yes/no questions](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2924–2936, Minneapolis, Minnesota. Association for Computational Linguistics.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *ArXiv*, abs/1803.05457.
- Arman Cohan, Waleed Ammar, Madeleine van Zuylen, and Field Cady. 2019. [Structural scaffolds for citation intent classification in scientific publications](#). 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 3586–3596, Minneapolis, Minnesota. Association for Computational Linguistics.
- Ido Dagan, Oren Glickman, and Bernardo Magnini. 2005. The pascal recognising textual entailment challenge. In *Machine Learning Challenges Workshop*, pages 177–190. Springer.
- Pradeep Dasigi, Nelson F. Liu, Ana Marasović, Noah A. Smith, and Matt Gardner. 2019. [Quoref: A reading comprehension dataset with questions requiring coreferential reasoning](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5925–5932, Hong Kong, China. Association for Computational Linguistics.
- Thomas Davidson, Dana Warmusley, Michael Macy, and Ingmar Weber. 2017. Automated hate speech detection and the problem of offensive language. In *Proceedings of the 11th International AAAI Conference on Web and Social Media, ICWSM ’17*, pages 512–515.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- T. Diggelmann, Jordan L. Boyd-Graber, Jannis Bulian, Massimiliano Ciaramita, and Markus Leippold. 2020. Climate-fever: A dataset for verification of real-world climate claims. *ArXiv*, abs/2012.00614.
- Emily Dinan, Stephen Roller, Kurt Shuster, Angela Fan, Michael Auli, and Jason Weston. 2019. [Wizard of wikipedia: Knowledge-powered conversational agents](#). In *International Conference on Learning Representations*.
- William B. Dolan and Chris Brockett. 2005. [Automatically constructing a corpus of sentential paraphrases](#). In *Proceedings of the Third International Workshop on Paraphrasing (IWP2005)*.
- Zi-Yi Dou, Keyi Yu, and Antonios Anastasopoulos. 2019. [Investigating meta-learning algorithms for low-resource natural language understanding tasks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1192–1197, Hong Kong, China. Association for Computational Linguistics.
- Matthew Dunn, Levent Sagun, Mike Higgins, V. U. Güney, Volkan Cirik, and Kyunghyun Cho. 2017. Searchqa: A new q&a dataset augmented with context from a search engine. *ArXiv*, abs/1704.05179.
- Ondřej Dušek, David M. Howcroft, and Verena Rieser. 2019. [Semantic noise matters for neural natural language generation](#). In *Proc. of the 12th International Conference on Natural Language Generation*, pages 421–426, Tokyo, Japan. Association for Computational Linguistics.
- Ondřej Dušek, Jekaterina Novikova, and Verena Rieser. 2020. [Evaluating the State-of-the-Art of End-to-End Natural Language Generation: The E2E NLG Challenge](#). *Computer Speech & Language*, 59:123–156.
- Hady Elsahar, Pavlos Vougiouklis, Arslan Remaci, Christophe Gravier, Jonathon Hare, Frederique Laforest, and Elena Simperl. 2018. [T-REx: A large scale alignment of natural language with knowledge base triples](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC-2018)*, Miyazaki, Japan. European Languages Resources Association (ELRA).
- Alexander Fabbri, Irene Li, Tianwei She, Suyi Li, and Dragomir Radev. 2019. [Multi-news: A large-scale multi-document summarization dataset and abstractive hierarchical model](#). In *Proceedings of the 57th*

- Annual Meeting of the Association for Computational Linguistics*, pages 1074–1084, Florence, Italy. Association for Computational Linguistics.
- Angela Fan, Yacine Jernite, Ethan Perez, David Grangier, Jason Weston, and Michael Auli. 2019. [ELI5: Long form question answering](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3558–3567, Florence, Italy. Association for Computational Linguistics.
- Manaal Faruqui and Dipanjan Das. 2018. [Identifying well-formed natural language questions](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 798–803, Brussels, Belgium. Association for Computational Linguistics.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. [Model-agnostic meta-learning for fast adaptation of deep networks](#). In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1126–1135. PMLR.
- Tianyu Gao, A. Fisch, and Danqi Chen. 2020. Making pre-trained language models better few-shot learners. *ArXiv*, abs/2012.15723.
- Tianyu Gao, Xu Han, Hao Zhu, Zhiyuan Liu, Peng Li, Maosong Sun, and Jie Zhou. 2019. [FewRel 2.0: Towards more challenging few-shot relation classification](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6250–6255, Hong Kong, China. Association for Computational Linguistics.
- Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and Bill Dolan. 2007. The third pascal recognizing textual entailment challenge. In *Proceedings of the ACL-PASCAL workshop on textual entailment and paraphrasing*, pages 1–9. Association for Computational Linguistics.
- Ona de Gibert, Naiara Perez, Aitor García-Pablos, and Montse Cuadros. 2018. [Hate Speech Dataset from a White Supremacy Forum](#). In *Proceedings of the 2nd Workshop on Abusive Language Online (ALW2)*, pages 11–20, Brussels, Belgium. Association for Computational Linguistics.
- Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. 2019. [SAMSum corpus: A human-annotated dialogue dataset for abstractive summarization](#). In *Proceedings of the 2nd Workshop on New Frontiers in Summarization*, pages 70–79, Hong Kong, China. Association for Computational Linguistics.
- Andrew Gordon, Zornitsa Kozareva, and Melissa Roemmele. 2012. [SemEval-2012 task 7: Choice of plausible alternatives: An evaluation of commonsense causal reasoning](#). In *\*SEM 2012: The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pages 394–398, Montréal, Canada. Association for Computational Linguistics.
- Jiatao Gu, Yong Wang, Yun Chen, Victor O. K. Li, and Kyunghyun Cho. 2018. [Meta-learning for low-resource neural machine translation](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3622–3631, Brussels, Belgium. Association for Computational Linguistics.
- Harsha Gurulingappa, Abdul Mateen Rajput, Angus Roberts, Juliane Fluck, Martin Hofmann-Apitius, and Luca Toldo. 2012. [Development of a benchmark corpus to support the automatic extraction of drug-related adverse effects from medical case reports](#). *Journal of Biomedical Informatics*, 45(5):885–892. Text Mining and Natural Language Processing in Pharmacogenomics.
- Xu Han, Hao Zhu, Pengfei Yu, Ziyun Wang, Yuan Yao, Zhiyuan Liu, and Maosong Sun. 2018. [FewRel: A large-scale supervised few-shot relation classification dataset with state-of-the-art evaluation](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4803–4809, Brussels, Belgium. Association for Computational Linguistics.
- Luheng He, Mike Lewis, and Luke Zettlemoyer. 2015. [Question-answer driven semantic role labeling: Using natural language to annotate natural language](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 643–653, Lisbon, Portugal. Association for Computational Linguistics.
- Johannes Hoffart, Mohamed Amir Yosef, Ilaria Bordino, Hagen Fürstenauf, Manfred Pinkal, Marc Spaniol, Bilyana Taneva, Stefan Thater, and Gerhard Weikum. 2011. [Robust disambiguation of named entities in text](#). In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, pages 782–792, Edinburgh, Scotland, UK. Association for Computational Linguistics.
- Eduard Hovy, Laurie Gerber, Ulf Hermjakob, Chin-Yew Lin, and Deepak Ravichandran. 2001. [Toward semantics-based answer pinpointing](#). In *Proceedings of the First International Conference on Human Language Technology Research*.
- Lifu Huang, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. [Cosmos QA: Machine reading comprehension with contextual commonsense reasoning](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2391–2401, Hong Kong, China. Association for Computational Linguistics.

- Chao Jiang, Mounica Maddela, Wuwei Lan, Yang Zhong, and Wei Xu. 2020. [Neural CRF model for sentence alignment in text simplification](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7943–7960, Online. Association for Computational Linguistics.
- Kelvin Jiang, Dekun Wu, and Hui Jiang. 2019. [FreebaseQA: A new factoid QA data set matching trivia-style question-answer pairs with Freebase](#). 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 318–323, Minneapolis, Minnesota. Association for Computational Linguistics.
- Xisen Jin, Mohammad Rostami, and Xiang Ren. 2021. Lifelong learning of few-shot learners across nlp tasks. *ArXiv*, abs/2104.08808.
- Daniel Khashabi, Snigdha Chaturvedi, Michael Roth, Shyam Upadhyay, and Dan Roth. 2018. [Looking beyond the surface: A challenge set for reading comprehension over multiple sentences](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 252–262, New Orleans, Louisiana. Association for Computational Linguistics.
- Daniel Khashabi, Sewon Min, Tushar Khot, Ashish Sabharwal, Oyvind Tafjord, Peter Clark, and Hananeh Hajishirzi. 2020. [UNIFIEDQA: Crossing format boundaries with a single QA system](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1896–1907, Online. Association for Computational Linguistics.
- Tushar Khot, Peter Clark, Michal Guerquin, Peter Jansen, and Ashish Sabharwal. 2020. [Qasc: A dataset for question answering via sentence composition](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):8082–8090.
- Tushar Khot, Daniel Khashabi, Kyle Richardson, Peter Clark, and Ashish Sabharwal. 2021. [Text modular networks: Learning to decompose tasks in the language of existing models](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1264–1279, Online. Association for Computational Linguistics.
- Tushar Khot, Ashish Sabharwal, and Peter Clark. 2018. SciTail: A textual entailment dataset from science question answering. In *AAAI*.
- Byeongchang Kim, Hyunwoo Kim, and Gunhee Kim. 2019. [Abstractive summarization of Reddit posts with multi-level memory networks](#). 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 2519–2531, Minneapolis, Minnesota. Association for Computational Linguistics.
- Neema Kotonya and Francesca Toni. 2020. [Explainable automated fact-checking for public health claims](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7740–7754, Online. Association for Computational Linguistics.
- 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](#). *Transactions of the Association for Computational Linguistics*, 7:453–466.
- Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. [RACE: Large-scale Reading comprehension dataset from examinations](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 785–794, Copenhagen, Denmark. Association for Computational Linguistics.
- Rémi Lebret, David Grangier, and Michael Auli. 2016. [Neural text generation from structured data with application to the biography domain](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1203–1213, Austin, Texas. Association for Computational Linguistics.
- 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. *Semantic Web*, 6:167–195.
- Hector J. Levesque, Ernest Davis, and Leora Morgenstern. 2012. The winograd schema challenge. In *Proceedings of the Thirteenth International Conference on Principles of Knowledge Representation and Reasoning, KR’12*, page 552–561. AAAI Press.
- Omer Levy, Minjoon Seo, Eunsol Choi, and Luke Zettlemoyer. 2017. [Zero-shot relation extraction via reading comprehension](#). In *Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017)*, pages 333–342, Vancouver, Canada. Association for Computational Linguistics.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational*

- Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Quentin Lhoest, Albert Villanova del Moral, Yacine Jernite, A. Thakur, Patrick von Platen, Suraj Patil, Julien Chaumond, Mariama Drame, Julien Plu, Lewis Tunstall, Joe Davison, Mario vSavsko, Guntan Chhablani, Bhavitvya Malik, Simon Brandeis, Teven Le Scao, Victor Sanh, Canwen Xu, Nicolas Patry, Angelina McMillan-Major, Philipp Schmid, Sylvain Gugger, Clement Delangue, Th’eo Matussiere, Lysandre Debut, Stas Bekman, Pierric Cistac, Thibault Goehringer, Victor Mustar, Franccois Lagunas, Alexander M. Rush, and Thomas Wolf. 2021. Datasets: A community library for natural language processing.
- Xin Li and Dan Roth. 2002. [Learning question classifiers](#). In *COLING 2002: The 19th International Conference on Computational Linguistics*.
- 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 *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6862–6868, Online. Association for Computational Linguistics.
- 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 *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1823–1840, Online. Association for Computational Linguistics.
- Kevin Lin, Oyvind Tafjord, Peter Clark, and Matt Gardner. 2019. [Reasoning over paragraph effects in situations](#). In *Proceedings of the 2nd Workshop on Machine Reading for Question Answering*, pages 58–62, Hong Kong, China. Association for Computational Linguistics.
- Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017. [Program induction by rationale generation: Learning to solve and explain algebraic word problems](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 158–167, Vancouver, Canada. Association for Computational Linguistics.
- Tal Linzen. 2020. [How can we accelerate progress towards human-like linguistic generalization?](#) In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5210–5217, Online. Association for Computational Linguistics.
- Annie Louis, Dan Roth, and Filip Radlinski. 2020. [“I’d rather just go to bed”: Understanding indirect answers](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7411–7425, Online. Association for Computational Linguistics.
- Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011. [Learning word vectors for sentiment analysis](#). In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon, USA. Association for Computational Linguistics.
- Pekka Malo, Ankur Sinha, Pekka Korhonen, Jyrki Wallenius, and Pyry Takala. 2014. [Good debt or bad debt: Detecting semantic orientations in economic texts](#). *J. Assoc. Inf. Sci. Technol.*, 65(4):782–796.
- Irene Manotas, Ngoc Phuoc An Vo, and Vadim Sheinin. 2020. [LiMiT: The literal motion in text dataset](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 991–1000, Online. Association for Computational Linguistics.
- Marco Marelli, Stefano Menini, Marco Baroni, Luisa Bentivogli, Raffaella Bernardi, and Roberto Zamparelli. 2014. [A SICK cure for the evaluation of compositional distributional semantic models](#). In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC-2014)*, pages 216–223, Reykjavik, Iceland. European Languages Resources Association (ELRA).
- Marie-Catherine de Marneffe, Mandy Simons, and Judith Tonhauser. 2019. [The commitmentbank: Investigating projection in naturally occurring discourse](#). *Proceedings of Sinn und Bedeutung*, 23(2):107–124.
- Binny Mathew, Punyajoy Saha, Seid Muhie Yimam, Chris Biemann, Pawan Goyal, and Animesh Mukherjee. 2020. [Hatexplain: A benchmark dataset for explainable hate speech detection](#). *arXiv preprint arXiv:2012.10289*.
- Julian McAuley and J. Leskovec. 2013. [Hidden factors and hidden topics: understanding rating dimensions with review text](#). *Proceedings of the 7th ACM conference on Recommender systems*.
- Bryan McCann, N. Keskar, Caiming Xiong, and R. Socher. 2018. [The natural language decathlon: Multitask learning as question answering](#). *ArXiv*, abs/1806.08730.
- 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 *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD ’20*, page 3458–3465, New York, NY, USA. Association for Computing Machinery.
- 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 *Proceedings of the 2018 Conference on*

- Empirical Methods in Natural Language Processing*, pages 2381–2391, Brussels, Belgium. Association for Computational Linguistics.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2021. Cross-task generalization via natural language crowdsourcing instructions. *arXiv preprint arXiv:2104.08773*.
- Ioannis Mollas, Zoe Chrysopoulou, Stamatis Karlos, and Grigorios Tsoumakas. 2020. Ethos: an online hate speech detection dataset. *ArXiv*, abs/2006.08328.
- Shikhar Murty, T. Hashimoto, and Christopher D. Manning. 2021. Dreca: A general task augmentation strategy for few-shot natural language inference.
- 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 *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1953–1967, Online. Association for Computational Linguistics.
- Courtney Napoles, Matthew Gormley, and Benjamin Van Durme. 2012. Annotated Gigaword. In *Proceedings of the Joint Workshop on Automatic Knowledge Base Construction and Web-scale Knowledge Extraction (AKBC-WEKEX)*, pages 95–100, Montréal, Canada. Association for Computational Linguistics.
- 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 *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1797–1807, Brussels, Belgium. Association for Computational Linguistics.
- Alex Nichol, Joshua Achiam, and John Schulman. 2018. On first-order meta-learning algorithms. *ArXiv*, abs/1803.02999.
- Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2020. Adversarial NLI: A new benchmark for natural language understanding. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4885–4901, Online. Association for Computational Linguistics.
- Farhad Nooralahzadeh, Giannis Bekoulis, Johannes Bjerva, and Isabelle Augenstein. 2020. Zero-shot cross-lingual transfer with meta learning. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4547–4562, Online. Association for Computational Linguistics.
- A. Othman and M. Jemni. 2012. English-asl gloss parallel corpus 2012: Aslg-pc12.
- Bo Pang and Lillian Lee. 2005. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL’05)*, pages 115–124, Ann Arbor, Michigan. Association for Computational Linguistics.
- Dimitris Pappas, Petros Stavropoulos, Ion Androustopoulos, and Ryan McDonald. 2020. BioMRC: A dataset for biomedical machine reading comprehension. In *Proceedings of the 19th SIGBioMed Workshop on Biomedical Language Processing*, pages 140–149, Online. Association for Computational Linguistics.
- 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 *Automated Knowledge Base Construction*.
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2463–2473, Hong Kong, China. 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 *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1267–1273, Minneapolis, Minnesota. Association for Computational Linguistics.
- 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 *Proceedings of the 28th International Conference on Computational Linguistics*, pages 3285–3301, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Yada Pruksachatkun, Jason Phang, Haokun Liu, Phu Mon Htut, Xiaoyi Zhang, Richard Yuanzhe Pang, Clara Vania, Katharina Kann, and Samuel R. Bowman. 2020. Intermediate-task transfer learning with pretrained language models: When and why does it work? In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5231–5247, Online. Association for Computational Linguistics.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring

- the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67.
- Altaf Rahman and Vincent Ng. 2012. Resolving complex cases of definite pronouns: The Winograd schema challenge. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 777–789, Jeju Island, Korea. Association for Computational Linguistics.
- Nazneen Fatema Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. Explain yourself! leveraging language models for commonsense reasoning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4932–4942, Florence, Italy. Association for Computational Linguistics.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.
- Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2019. Towards empathetic open-domain conversation models: A new benchmark and dataset. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5370–5381, Florence, Italy. Association for Computational Linguistics.
- Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. Beyond accuracy: Behavioral testing of NLP models with CheckList. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4902–4912, Online. Association for Computational Linguistics.
- Anna Rogers, Olga Kovaleva, Matthew Downey, and Anna Rumshisky. 2020. Getting closer to ai complete question answering: A set of prerequisite real tasks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):8722–8731.
- Amrita Saha, Rahul Aralikkatte, Mitesh M. Khapra, and Karthik Sankaranarayanan. 2018. DuoRC: Towards complex language understanding with paraphrased reading comprehension. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1683–1693, Melbourne, Australia. Association for Computational Linguistics.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2020. Winogrande: An adversarial winograd schema challenge at scale. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):8732–8740.
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. 2019. Social IQa: Commonsense reasoning about social interactions. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4463–4473, Hong Kong, China. Association for Computational Linguistics.
- Elvis Saravia, Hsien-Chi Toby Liu, Yen-Hao Huang, Junlin Wu, and Yi-Shin Chen. 2018. CARER: Contextualized affect representations for emotion recognition. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3687–3697, Brussels, Belgium. Association for Computational Linguistics.
- Timo Schick and Hinrich Schütze. 2020a. Exploiting cloze questions for few-shot text classification and natural language inference. *Computing Research Repository*, arXiv:2001.07676.
- Timo Schick and Hinrich Schütze. 2020b. It’s not just size that matters: Small language models are also few-shot learners. *Computing Research Repository*, arXiv:2009.07118.
- Emily Sheng and David Uthus. 2020. Investigating societal biases in a poetry composition system. In *Proceedings of the Second Workshop on Gender Bias in Natural Language Processing*, pages 93–106, Barcelona, Spain (Online). Association for Computational Linguistics.
- Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. 2020. AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4222–4235, Online. Association for Computational Linguistics.
- Damien Sileo, Tim Van De Cruys, Camille Pradel, and Philippe Muller. 2019. Mining discourse markers for unsupervised sentence representation learning. 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 3477–3486, Minneapolis, Minnesota. Association for Computational Linguistics.
- 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 *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.
- Trevor Scott Standley, A. Zamir, Dawn Chen, L. Guibas, Jitendra Malik, and S. Savarese. 2020.



- Which tasks should be learned together in multi-task learning? In *ICML*.
- 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**. *Transactions of the Association for Computational Linguistics*, 7:217–231.
- Oyvind Tafjord, Peter Clark, Matt Gardner, Wen-tau Yih, and Ashish Sabharwal. 2019a. **Quarel: A dataset and models for answering questions about qualitative relationships**. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):7063–7071.
- Oyvind Tafjord, Matt Gardner, Kevin Lin, and Peter Clark. 2019b. **QuaRTz: An open-domain dataset of qualitative relationship questions**. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5941–5946, Hong Kong, China. Association for Computational Linguistics.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. **CommonsenseQA: A question answering challenge targeting commonsense knowledge**. 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 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.
- Derek Tam, R. R. Menon, M. Bansal, Shashank Srivastava, and Colin Raffel. 2021. Improving and simplifying pattern exploiting training. *ArXiv*, abs/2103.11955.
- Niket Tandon, Bhavana Dalvi, Keisuke Sakaguchi, Peter Clark, and Antoine Bosselut. 2019. **WIQA: A dataset for “what if...” reasoning over procedural text**. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6076–6085, Hong Kong, China. Association for Computational Linguistics.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. **FEVER: a large-scale dataset for fact extraction and VERification**. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 809–819, New Orleans, Louisiana. Association for Computational Linguistics.
- Eleni Triantafyllou, Tyler Zhu, Vincent Dumoulin, Pascal Lamblin, Utku Evci, Kelvin Xu, Ross Goroshin, Carles Gelada, Kevin Swersky, Pierre-Antoine Manzagol, and Hugo Larochelle. 2020. **Meta-dataset: A dataset of datasets for learning to learn from few examples**. In *International Conference on Learning Representations*.
- Sowmya Vajjala and Ivana Lučić. 2018. **OneStopEnglish corpus: A new corpus for automatic readability assessment and text simplification**. In *Proceedings of the Thirteenth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 297–304, New Orleans, Louisiana. Association for Computational Linguistics.
- Tu Vu, Tong Wang, Tsendsuren Munkhdalai, Alessandro Sordani, Adam Trischler, Andrew Mattarella-Micke, Subhansu Maji, and Mohit Iyyer. 2020. **Exploring and predicting transferability across NLP tasks**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7882–7926, Online. Association for Computational Linguistics.
- Sinong Wang, Han Fang, Madian Khabsa, Hanzi Mao, and Hao Ma. 2021. Entailment as few-shot learner. *arXiv preprint arXiv:2104.14690*.
- Sinong Wang, Madian Khabsa, and Hao Ma. 2020. **To pretrain or not to pretrain: Examining the benefits of pretraining on resource rich tasks**. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2209–2213, Online. Association for Computational Linguistics.
- William Yang Wang. 2017. **“liar, liar pants on fire”: A new benchmark dataset for fake news detection**. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 422–426, Vancouver, Canada. Association for Computational Linguistics.
- 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**. *Transactions of the Association for Computational Linguistics*, 8:377–392.
- Alex Warstadt, Amanpreet Singh, and Samuel R. Bowman. 2019. **Neural network acceptability judgments**. *Transactions of the Association for Computational Linguistics*, 7:625–641.
- Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2021. Finetuned language models are zero-shot learners.
- Johannes Welbl, Nelson F. Liu, and Matt Gardner. 2017. **Crowdsourcing multiple choice science questions**. In *Proceedings of the 3rd Workshop on Noisy User-generated Text*, pages 94–106, Copenhagen, Denmark. Association for Computational Linguistics.
- Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. **A broad-coverage challenge corpus for sentence understanding through inference**. In *Proceedings of the 2018 Conference of the North American*

- Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122, New Orleans, Louisiana. Association for Computational Linguistics.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.
- Tomer Wolfson, Mor Geva, Ankit Gupta, Matt Gardner, Yoav Goldberg, Daniel Deutch, and Jonathan Berant. 2020. [Break it down: A question understanding benchmark](#). *Transactions of the Association for Computational Linguistics*, 8:183–198.
- Sen Wu, Hongyang R. Zhang, and Christopher Ré. 2020. [Understanding and improving information transfer in multi-task learning](#). In *International Conference on Learning Representations*.
- Wenhan Xiong, Jiawei Wu, Hong Wang, Vivek Kulkarini, Mo Yu, Shiyu Chang, Xiaoxiao Guo, and William Yang Wang. 2019. [TWEETQA: A social media focused question answering dataset](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5020–5031, Florence, Italy. Association for Computational Linguistics.
- Yi Yang, Wen-tau Yih, and Christopher Meek. 2015. [WikiQA: A challenge dataset for open-domain question answering](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 2013–2018, Lisbon, Portugal. Association for Computational Linguistics.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. [HotpotQA: A dataset for diverse, explainable multi-hop question answering](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2369–2380, Brussels, Belgium. Association for Computational Linguistics.
- Wenpeng Yin, Nazneen Fatema Rajani, Dragomir Radev, Richard Socher, and Caiming Xiong. 2020. [Universal natural language processing with limited annotations: Try few-shot textual entailment as a start](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8229–8239, Online. Association for Computational Linguistics.
- Dani Yogatama, Cyprien de Masson d’Autume, Jerome Connor, Tomás Kociský, Mike Chrzanowski, Lingpeng Kong, A. Lazaridou, Wang Ling, L. Yu, Chris Dyer, and P. Blunsom. 2019. [Learning and evaluating general linguistic intelligence](#). *ArXiv*, abs/1901.11373.
- Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir Radev. 2018. [Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-SQL task](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3911–3921, Brussels, Belgium. Association for Computational Linguistics.
- Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey Levine. 2020. [Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning](#). In *Proceedings of the Conference on Robot Learning*, volume 100 of *Proceedings of Machine Learning Research*, pages 1094–1100. PMLR.
- Amir R. Zamir, Alexander Sax, William B. Shen, Leonidas J. Guibas, Jitendra Malik, and Silvio Savarese. 2018. [Taskonomy: Disentangling task transfer learning](#). In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE.
- Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. 2018. [SWAG: A large-scale adversarial dataset for grounded commonsense inference](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 93–104, Brussels, Belgium. Association for Computational Linguistics.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. [HellaSwag: Can a machine really finish your sentence?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, Florence, Italy. Association for Computational Linguistics.
- Hao Zhang, Jae Ro, and Richard Sproat. 2020. [Semi-supervised URL segmentation with recurrent neural networks pre-trained on knowledge graph entities](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4667–4675, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Rui Zhang and Joel Tetreault. 2019. [This email could save your life: Introducing the task of email subject line generation](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 446–456, Florence, Italy. Association for Computational Linguistics.
- Sheng Zhang, X. Liu, J. Liu, Jianfeng Gao, Kevin Duh, and Benjamin Van Durme. 2018. [Record: Bridging](#)

the gap between human and machine commonsense reading comprehension. *ArXiv*, abs/1810.12885.

Tianyi Zhang, Felix Wu, Arzoo Katiyar, Kilian Q Weinberger, and Yoav Artzi. 2021. [Revisiting few-sample {bert} fine-tuning](#). In *International Conference on Learning Representations*.

Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 1, NIPS'15*, page 649–657, Cambridge, MA, USA. MIT Press.

Yuan Zhang, Jason Baldridge, and Luheng He. 2019. [PAWS: Paraphrase adversaries from word scrambling](#). 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 1298–1308, Minneapolis, Minnesota. Association for Computational Linguistics.

Ruiqi Zhong, Kristy Lee, Zheng Zhang, and D. Klein. 2021. Adapting language models for zero-shot learning by meta-tuning on dataset and prompt collections. *ArXiv*, abs/2104.04670.

Victor Zhong, Caiming Xiong, and Richard Socher. 2017. Seq2sql: Generating structured queries from natural language using reinforcement learning. *CoRR*, abs/1709.00103.

Ben Zhou, Daniel Khashabi, Qiang Ning, and Dan Roth. 2019. [“going on a vacation” takes longer than “going for a walk”](#): A study of temporal commonsense understanding. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3363–3369, Hong Kong, China. Association for Computational Linguistics.

## A Selected Tasks in NLP Few-shot Gym

Table 3: Tasks in NLP Few-shot Gym.

Task Name	Ontology	Reference
acronym_identification	other	Pouran Ben Veysseh et al. 2020
ade_corpus_v2-classification	cls/other	Gurulingappa et al. 2012
ade_corpus_v2-dosage	other/slot filling	Gurulingappa et al. 2012
ade_corpus_v2-effect	other/slot filling	Gurulingappa et al. 2012
adversarialqa	qa/machine reading comprehension	Bartolo et al. 2020
aeslc	cg/summarization	Zhang and Tetreault 2019
ag_news	cls/topic	Gulli (link)
ai2_arc	qa/multiple-choice qa	Clark et al. 2018
amazon_polarity	cls/sentiment analysis	McAuley and Leskovec 2013
anli	cls/nli	Nie et al. 2020
app_reviews	other/regression	Missing
aqua_rat	qa/multiple-choice qa	Ling et al. 2017
art (abductive nli)	other	Bhagavatula et al. 2020
aslg_pc12	other	Othman and Jemni 2012
biomrc	qa/machine reading comprehension	Pappas et al. 2020
blimp-anaphor_gender_agreement	other/linguistic phenomenon	Warstadt et al. 2020
blimp-anaphor_number_agreement	other/linguistic phenomenon	Warstadt et al. 2020
blimp-determiner_noun_agreement_with_adj_irregular_1	other/linguistic phenomenon	Warstadt et al. 2020
blimp-ellipsis_n_bar_1	other/linguistic phenomenon	Warstadt et al. 2020
blimp-ellipsis_n_bar_2	other/linguistic phenomenon	Warstadt et al. 2020
blimp-existential_there_quantifiers_1	other/linguistic phenomenon	Warstadt et al. 2020
blimp-irregular_past_participle_adjectives	other/linguistic phenomenon	Warstadt et al. 2020
blimp-sentential_negation_npi_licensor_present	other/linguistic phenomenon	Warstadt et al. 2020
blimp-sentential_negation_npi_scope	other/linguistic phenomenon	Warstadt et al. 2020
blimp-wh_questions_object_gap	other/linguistic phenomenon	Warstadt et al. 2020
boolq	qa/binary	Clark et al. 2019
break-QDMR	other	Wolfson et al. 2020
break-QDMR-high-level	other	Wolfson et al. 2020
circa	cls/other	Louis et al. 2020
climate_fever	cls/fact checking	Diggelmann et al. 2020
codah	qa/multiple-choice qa	Chen et al. 2019
common_gen	other	Lin et al. 2020b
commonsense_qa	qa/multiple-choice qa	Talmor et al. 2019
cos_e	other/generate explanation	Rajani et al. 2019
cosmos_qa	qa/multiple-choice qa	Huang et al. 2019
crawl_domain	other	Zhang et al. 2020
crows_pairs	other	Nangia et al. 2020
dbpedia_14	cls/topic	Lehmann et al. 2015
definite_pronoun_resolution	other	Rahman and Ng 2012
discovery	cls/other	Sileo et al. 2019
dream	qa/multiple-choice qa	Sun et al. 2019
duorc	qa/machine reading comprehension	Saha et al. 2018
e2e_nlg_cleaned	other	Dušek et al. 2020, 2019
eli5-askh	qa/long-form qa	Fan et al. 2019
eli5-asks	qa/long-form qa	Fan et al. 2019
eli5-eli5	qa/long-form qa	Fan et al. 2019
emo	cls/emotion	Chatterjee et al. 2019
emotion	cls/emotion	Saravia et al. 2018
empathetic_dialogues	cg/dialogue	Rashkin et al. 2019
ethos-directed_vs_generalized	cls/hate speech detection	Mollas et al. 2020
ethos-disability	cls/hate speech detection	Mollas et al. 2020
ethos-gender	cls/hate speech detection	Mollas et al. 2020
ethos-national_origin	cls/hate speech detection	Mollas et al. 2020
ethos-race	cls/hate speech detection	Mollas et al. 2020
ethos-religion	cls/hate speech detection	Mollas et al. 2020
ethos-sexual_orientation	cls/hate speech detection	Mollas et al. 2020
financial_phrasebank	cls/sentiment analysis	Malo et al. 2014
freebase_qa	qa/closed-book qa	Jiang et al. 2019
gigaword	cg/summarization	Napoles et al. 2012
glue-cola	cls/other	Warstadt et al. 2019
glue-mnli	cls/nli	Williams et al. 2018
glue-mrpc	cls/paraphrase	Dolan and Brockett 2005
glue-qnli	cls/nli	Rajpurkar et al. 2016
glue-qqp	cls/paraphrase	(link)
glue-rtc	cls/nli	Dagan et al. 2005; Bar-Haim et al. 2006
glue-sst2	cls/sentiment analysis	Giampiccolo et al. 2007; Bentivogli et al. 2009
glue-wnli	cls/nli	Socher et al. 2013
google_wellformed_query	cls/other	Levesque et al. 2012
hate_speech18	cls/hate speech detection	Faruqui and Das 2018
hate_speech_offensive	cls/hate speech detection	de Gibert et al. 2018
hatexplain	cls/hate speech detection	Davidson et al. 2017
health_fact	cls/hate speech detection	Mathew et al. 2020
hellaswag	cls/fact checking	Kotonya and Toni 2020
hotpot_qa	qa/multiple-choice qa	Zellers et al. 2019
imdb	qa/machine reading comprehension	Yang et al. 2018
jeopardy	cls/sentiment analysis	Maas et al. 2011
kilt_ay2	qa/closed-book qa	(link)
	other/entity linking	Hoffart et al. 2011

Continued on next page

Task Name	Ontology	Reference
kilt_fever	cls/fact checking	Thorne et al. 2018
kilt_hotpotqa	qa/closed-book qa	Yang et al. 2018
kilt_nq	qa/closed-book qa	Kwiatkowski et al. 2019
kilt_trex	qa/closed-book qa	Elsahar et al. 2018
kilt_wow	cg/dialogue	Dinan et al. 2019
kilt_zsre	qa/closed-book qa	Levy et al. 2017
lama-conceptnet	qa/closed-book qa	Petroni et al. 2019, 2020
lama-google_re	qa/closed-book qa	Petroni et al. 2019, 2020
lama-squad	qa/closed-book qa	Petroni et al. 2019, 2020
lama-trex	qa/closed-book qa	Petroni et al. 2019, 2020
liar	cls/fact checking	Wang 2017
limit	other	Manotas et al. 2020
math_qa	qa/multiple-choice qa	Amini et al. 2019
mc_taco	qa/binary	Zhou et al. 2019
medical_questions_pairs	cls/paraphrase	McCreery et al. 2020
mocha	other/regression	Chen et al. 2020a
multi_news	cg/summarization	Fabbri et al. 2019
numer_sense	qa/closed-book qa	Lin et al. 2020a
onestop_english	cls/other	Vajjala and Lučić 2018
openbookqa	qa/multiple-choice qa	Mihaylov et al. 2018
paws	cls/paraphrase	Zhang et al. 2019
piqa	other	Bisk et al. 2020
poem_sentiment	cls/sentiment analysis	Sheng and Uthus 2020
proto_qa	other	Boratko et al. 2020
qa_srl	other	He et al. 2015
qasc	qa/multiple-choice qa	Khot et al. 2020
quail	qa/multiple-choice qa	Rogers et al. 2020
quarel	qa/multiple-choice qa	Tafjord et al. 2019a
quartz-no_knowledge	qa/multiple-choice qa	Tafjord et al. 2019b
quartz-with_knowledge	qa/multiple-choice qa	Tafjord et al. 2019b
quoref	qa/machine reading comprehension	Dasigi et al. 2019
race-high	qa/multiple-choice qa	Lai et al. 2017
race-middle	qa/multiple-choice qa	Lai et al. 2017
reddit_tifu-title	cg/summarization	Kim et al. 2019
reddit_tifu-tldr	cg/summarization	Kim et al. 2019
ropes	qa/machine reading comprehension	Lin et al. 2019
rotten_tomatoes	cls/sentiment analysis	Pang and Lee 2005
samsun	cg/summarization	Gliwa et al. 2019
scicite	cls/other	Cohan et al. 2019
sciq	qa/multiple-choice qa	Welbl et al. 2017
scitail	cls/nli	Khot et al. 2018
search_qa	qa/closed-book qa	Dunn et al. 2017
sick	cls/nli	Marelli et al. 2014
sms_spam	cls/other	Almeida et al. 2011
social_i_qa	qa/multiple-choice qa	Sap et al. 2019
spider	cg/other	Yu et al. 2018
squad-no_context	qa/closed-book qa	Rajpurkar et al. 2016
squad-with_context	qa/machine reading comprehension	Rajpurkar et al. 2016
superglue-cb	cls/nli	de Marneffe et al. 2019
superglue-copa	qa/multiple-choice qa	Gordon et al. 2012
superglue-multirc	qa/multiple-choice qa	Khashabi et al. 2018
superglue-record	qa/machine reading comprehension	Zhang et al. 2018
superglue-rte	cls/nli	Dagan et al. 2005; Bar-Haim et al. 2006 Giampiccolo et al. 2007; Bentivogli et al. 2009
superglue-wic	cls/other	Pilehvar and Camacho-Collados 2019
superglue-wsc	cls/other	Levesque et al. 2012
swag	qa/multiple-choice qa	Zellers et al. 2018
tab_fact	cls/fact checking	Chen et al. 2020b
trec	cls/other	Li and Roth 2002; Hovy et al. 2001
trec-finegrained	cls/other	Li and Roth 2002; Hovy et al. 2001
tweet_eval-emoji	cls/emotion	Barbieri et al. 2020
tweet_eval-emotion	cls/emotion	Barbieri et al. 2020
tweet_eval-hate	cls/emotion	Barbieri et al. 2020
tweet_eval-irony	cls/emotion	Barbieri et al. 2020
tweet_eval-offensive	cls/emotion	Barbieri et al. 2020
tweet_eval-sentiment	cls/emotion	Barbieri et al. 2020
tweet_eval-stance_abortion	cls/emotion	Barbieri et al. 2020
tweet_eval-stance_atheism	cls/emotion	Barbieri et al. 2020
tweet_eval-stance_climate	cls/emotion	Barbieri et al. 2020
tweet_eval-stance_feminist	cls/emotion	Barbieri et al. 2020
tweet_eval-stance_hillary	cls/emotion	Barbieri et al. 2020
tweet_qa	qa/machine reading comprehension	Xiong et al. 2019
web_questions	qa/closed-book qa	Berant et al. 2013
wiki_auto	cls/other	Jiang et al. 2020
wiki_bio	cg/other	Lebret et al. 2016
wiki_qa	cls/other	Yang et al. 2015
wiki_split	cg/other	Botha et al. 2018
wikisql	cg/other	Zhong et al. 2017
wino_grande	qa/multiple-choice qa	Sakaguchi et al. 2020
wiqa	qa/multiple-choice qa	Tandon et al. 2019
xsum	cg/summarization	Narayan et al. 2018
yahoo_answers_topics	cls/topic	(link)
yelp_polarity	cls/sentiment analysis	Zhang et al. 2015; (link)
yelp_review_full	other/regression	Zhang et al. 2015; (link)

## B Details about Task Partition

### B.1 Partition 1. Random

```
1 {
2   "train": ['glue-mrpc', 'math_qa', 'quarel', 'e2e_nlg_cleaned', 'tweet_eval-stance_atheism', 'lama-squad',
3     'tab_fact', 'aqua_rat', 'tweet_eval-emoji', 'glue-wnli', 'codah', 'tweet_eval-offensive', 'wiki_qa',
4     'blimp-ellipsis_n_bar_1', 'openbookqa', 'sms_spam', 'acronym_identification', 'blimp-determiner_noun_agreement_with_adj_irregular_1',
5     'ethos-national_origin', 'spider', 'definite_pronoun_resolution', 'hellaswag', 'superglue-wsc', 'numer_sense', 'ade_corpus_v2-dosage',
6     'blimp-ellipsis_n_bar_2', 'kilt_ay2', 'squad-no_context', 'google_wellformed_query', 'xsum', 'wiqa',
7     'tweet_eval-stance_abortion', 'reddit_tifu_tldr', 'ade_corpus_v2-effect', 'qa_srl', 'ethos-religion',
8     'commonsense_qa', 'jeopardy', 'biomrc', 'superglue-multirc', 'ethos-race', 'eli5-askh', 'glue-qqp',
9     'paws', 'ethos-directed_vs_generalized', 'glue-sst2', 'mocha', 'tweet_eval-hate', 'glue-rte',
10    'blimp-anaphor_number_agreement', 'lama-conceptnet', 'hate_speech_offensive', 'superglue-wic',
11    'boolq', 'kilt_hotpotqa', 'quartz-no_knowledge', 'aslg_pc12', 'sick', 'tweet_eval-stance_climate',
12    'tweet_eval-sentiment', 'crows_pairs', 'glue-mnli', 'medical_questions_pairs', 'break-QDMR-high-level',
13    'qasc', 'imdb', 'ethos-gender', 'trec-finegrained', 'adversarialqa', 'onestop_english', 'web_questions',
14    'duorc', 'yelp_review_full', 'swag', 'proto_qa', 'scitail', 'tweet_eval-stance_feminist', 'limit',
15    'common_gen', 'scicite', 'blimp-irregular_past_participle_adjectives', 'social_i_qa', 'anli',
16    'kilt_zsre', 'cosmos_qa', 'superglue-record', 'squad-with_context', 'emotion', 'blimp-existential_there_quantifiers_1',
17    'race-middle', 'kilt_wow', 'sciq', 'wino_grande', 'rotten_tomatoes', 'superglue-cb', 'poem_sentiment',
18    'ropes', 'reddit_tifu_title', 'piqa', 'climate_fever', 'lama-google_re', 'search_qa', 'wiki_auto',
19    'mc_taco', 'blimp-wh_questions_object_gap', 'hotpot_qa', 'emo', 'kilt_nq', 'kilt_trex', 'quartz-with_knowledge',
20    'dbpedia_14', 'yahoo_answers_topics', 'app_reviews', 'superglue-copa', 'blimp-anaphor_gender_agreement',
21    'hate_speech18', 'gigaword', 'multi_news', 'aeslc', 'quail'],
22  "dev": ['cos_e', 'kilt_fever', 'eli5-asks', 'trec', 'eli5-eli5', 'art', 'empathetic_dialogues', 'tweet_qa',
23    'wikisql', 'lama-trex', 'tweet_eval-stance_hillary', 'discovery', 'tweet_eval-emotion', 'liar',
24    'wiki_bio', 'dream', 'ade_corpus_v2-classification', 'health_fact', 'samsun', 'financial_phrasebank'],
25  "test": ['quoref', 'wiki_split', 'ethos-disability', 'yelp_polarity', 'superglue-rte', 'glue-cola',
26    'ethos-sexual_orientation', 'blimp-sentential_negation_npi_scope', 'ai2_arc', 'amazon_polarity', 'race-high',
27    'blimp-sentential_negation_npi_licensor_present', 'tweet_eval-irony', 'break-QDMR', 'crawl_domain',
28    'freebase_qa', 'glue-qnli', 'hatexplain', 'ag_news', 'circa'],
29 }
```

### B.2 Partition 2.1. 45cls

```
1 {
2   "train": ["superglue-rte", "tweet_eval-sentiment", "discovery", "glue-rte", "superglue-wsc", "scicite",
3     "glue-mrpc", "tweet_eval-stance_hillary", "tweet_eval-offensive", "emotion", "hatexplain", "glue-cola",
4     "sick", "paws", "ethos-sexual_orientation", "glue-qqp", "tweet_eval-emotion", "sms_spam", "health_fact",
5     "glue-mnli", "imdb", "ethos-disability", "glue-wnli", "scitail", "trec-finegrained", "yahoo_answers_topics",
6     "liar", "glue-sst2", "tweet_eval-stance_abortion", "circa", "tweet_eval-stance_climate", "glue-qnli",
7     "tweet_eval-emoji", "ethos-directed_vs_generalized", "ade_corpus_v2-classification", "wiki_auto",
8     "hate_speech_offensive", "superglue-wic", "google_wellformed_query", "tweet_eval-irony", "ethos-gender",
9     "onestop_english", "trec", "rotten_tomatoes", "kilt_fever"],
10  "dev": ["tweet_eval-stance_feminist", "ethos-national_origin", "tweet_eval-hate", "ag_news", "amazon_polarity",
11    "hate_speech18", "poem_sentiment", "climate_fever", "medical_questions_pairs", "tweet_eval-stance_atheism"],
12  "test": ["superglue-cb", "dbpedia_14", "wiki_qa", "emo", "yelp_polarity", "ethos-religion", "financial_phrasebank",
13    "tab_fact", "anli", "ethos-race"],
14 }
```

### B.3 Partition 2.2. 23cls+22non-cls

```
1 {
2   "train": ["ade_corpus_v2-dosage", "biomrc", "blimp-ellipsis_n_bar_2", "blimp-sentential_negation_npi_scope",
3     "commonsense_qa", "crows_pairs", "duorc", "hellaswag", "kilt_zsre", "lama-google_re", "lama-squad",
4     "math_qa", "numer_sense", "openbookqa", "piqa", "proto_qa", "quartz-no_knowledge", "race-high",
5     "reddit_tifu_tldr", "ropes", "sciq", "wiki_bio", "discovery", "emotion", "ethos-disability",
6     "ethos-sexual_orientation", "glue-cola", "glue-mnli", "glue-mrpc", "glue-qqp", "glue-rte",
7     "glue-wnli", "hatexplain", "health_fact", "imdb", "paws", "scicite", "sick", "sms_spam",
8     "superglue-rte", "superglue-wsc", "tweet_eval-emotion", "tweet_eval-offensive", "tweet_eval-sentiment",
9     "tweet_eval-stance_hillary"],
10  "dev": ["tweet_eval-stance_feminist", "ethos-national_origin", "tweet_eval-hate", "ag_news", "amazon_polarity",
11    "hate_speech18", "poem_sentiment", "climate_fever", "medical_questions_pairs", "tweet_eval-stance_atheism"],
12  "test": ["superglue-cb", "dbpedia_14", "wiki_qa", "emo", "yelp_polarity", "ethos-religion", "financial_phrasebank",
13    "tab_fact", "anli", "ethos-race"]
14 }
```

### B.4 Partition 2.3. 45non-cls

```
1 {
```

```

2  "train": ["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", "hellaswag", "gigaword", "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"],
3  "dev": ["tweet_eval-stance_feminist", "ethos-national_origin", "tweet_eval-hate", "ag_news", "
    amazon_polarity", "hate_speech18", "poem_sentiment", "climate_fever", "medical_questions_pairs", "
    tweet_eval-stance_atheism"],
4  "test": ["superglue-cb", "dbpedia_14", "wiki_qa", "emo", "yelp_polarity", "ethos-religion", "
    financial_phrasebank", "tab_fact", "anli", "ethos-race"]
5  }

```

## B.5 Partition 3.1. Held-out-NLI

```

1  {
2  "train": ["ade_corpus_v2-classification", "ag_news", "amazon_polarity", "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-mrpc", "glue-qqp", "glue-sst2", "
    google_wellformed_query", "hate_speech18", "hate_speech_offensive", "hatexplain", "health_fact", "
    imdb", "kilt_fever", "liar", "medical_questions_pairs", "onestop_english", "paws", "poem_sentiment",
    "rotten_tomatoes", "scicite", "sick", "sms_spam", "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_auto", "wiki_qa", "yahoo_answers_topics", "yelp_polarity"
3  ],
4  "dev": [],
5  "test": ["anli", "glue-mnli", "glue-qnli", "glue-rte", "glue-wnli", "scitail", "sick", "superglue-cb"]
6  }

```

## B.6 Partition 3.2. Held-out-Para

```

1  {
2  "train": ["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_auto", "wiki_qa", "yahoo_answers_topics", "
    yelp_polarity"],
3  "dev": [],
4  "test": ["glue-mrpc", "glue-qqp", "medical_questions_pairs", "paws"]
5  }

```

## B.7 Partition 4.1. Held-out-MRC

```

1  {
2  "train": ["ai2_arc", "aqua_rat", "boolq", "codah", "commonsense_qa", "cosmos_qa", "dream", "eli5-askh",
    "eli5-asks", "eli5-eli5", "freebase_qa", "hellaswag", "jeopardy", "kilt_hotpotqa", "kilt_nq", "
    kilt_trex", "kilt_zsre", "lama-conceptnet", "lama-google_re", "lama-squad", "lama-trex", "math_qa",
    "mc_taco", "numer_sense", "openbookqa", "qasc", "quail", "quarel", "quartz-no_knowledge", "quartz-
    with_knowledge", "race-high", "race-middle", "sciq", "search_qa", "social_i_qa", "squad-no_context",
    "superglue-copa", "superglue-multirc", "swag", "web_questions", "wino_grande", "wiqa"
3  ],
4  "dev": [],
5  "test": ["adversarialqa", "biomrc", "duorc", "hotpot_qa", "quoref", "ropes", "squad-with_context", "
    superglue-record", "tweet_qa"],
6  }

```

## B.8 Partition 4.2. Held-out-MCQA

```

1  {
2  "train": ["adversarialqa", "biomrc", "boolq", "duorc", "eli5-askh", "eli5-asks", "eli5-eli5", "
    freebase_qa", "hotpot_qa", "jeopardy", "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", "squad-with_context", "superglue-multirc", "superglue-
    record", "tweet_qa", "web_questions"
3  ],
4  "dev": [],

```

```

5 "test": ["ai2_arc", "aqua_rat", "codah", "commonsense_qa", "cosmos_qa", "dream", "hellaswag", "math_qa",
6 "openbookqa", "qasc", "quail", "quarel", "quartz-no_knowledge", "quartz-with_knowledge", "race-
high", "race-middle", "sciq", "social_i_qa", "superglue-copa", "swag", "wino_grande", "wiqa"]
}

```

## B.9 Partition 5. Held-out-GLUE

To examine whether combining our methods with template-based training (Schick and Schütze, 2020a,b; Gao et al., 2020) results in even better few-shot performance, we add another partition that uses all non-GLUE classification tasks as  $\mathcal{T}_{train}$ , and all GLUE tasks as  $\mathcal{T}_{test}$ .

```


1 {
2 "train": ["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", "google_wellformed_query", "hate_speech18", "hate_speech_offensive", "
hatexplain", "health_fact", "imdb", "kilt_fever", "liar", "medical_questions_pairs", "onestop_english", "paws", "poem_sentiment", "rotten_tomatoes", "scicite", "scitail", "sick", "
sms_spam", "superglue-cb", "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_auto",
, "wiki_qa", "yahoo_answers_topics", "yelp_polarity"],
3 "dev": [],
4 "test": ["glue-cola", "glue-mnli", "glue-mrpc", "glue-qnli", "glue-qqp", "glue-rte", "glue-sst2", "glue-
wnli"]
5 }

```

Continued on next page.




## C Additional Results and Analysis

 Q4. Does the improved cross-task generalization ability go beyond few-shot settings?

In real-world applications, annotated data usually grow for a few-shot task over time. Is upstream learning still helpful when a target task has more shots? To study this question, we study CommonsenseQA (in *Held-out-Multiple-Choice Partition*), ROPES (in *Held-out-MRC Partition*), and MNLI (in *Held-out-NLI Partition*) as target tasks in medium and high-resource scenarios. We take their corresponding checkpoints after upstream learning and conduct experiments in medium and high-resource scenarios. That is, we randomly sample  $\{32, 64, \dots, 4096\}$  examples from the three datasets, and use them as  $\mathcal{D}_{train}$ . Then, we sample a  $\mathcal{D}_{dev}$  with the same size as  $\mathcal{D}_{train}$ , or has the size of 1024 if  $|\mathcal{D}_{train}| > 1024$ . We also try fine-tuning with the full dataset.<sup>6</sup> The performance of these settings is shown in Fig. 7.

From Fig. 7, we see that the benefits brought by upstream learning methods extend into medium resource cases with up to 2048 training examples. For CommonsenseQA, checkpoints from upstream learning outperform direct fine-tuning significantly, even with the full dataset. This finding encourages the use of upstream learning before task-specific fine-tuning when the target task has limited annotation. On the other hand, for resource-rich tasks (e.g., MNLI), the improvement brought by upstream learning diminishes. This aligns with the findings of (Wang et al., 2020) who discuss the benefits of pre-training on resource-rich tasks.


 Q5. Can we further improve few-shot performance by using different/larger pre-trained models?

We have been mainly using BART-Base (139M parameters) as the main network, while it is possible to further push the limits of few-shot learning by using scaling up to larger models or using different model architectures. Previous work has shown that scaling up model size leads to better performance (Raffel et al., 2020; Brown et al., 2020). Moreover, since meta-learning algorithms are naturally unstable, it is important to verify whether they

<sup>6</sup>We do five random samples of 1024 examples as  $\mathcal{D}_{dev}$  and use the remaining examples in the original train set as  $\mathcal{D}_{train}$ . We use the original dev set for testing.

function as expected with larger models. In Q5, we experiment with T5-v1.1-Base (248M)<sup>7</sup> and BART-Large (406M) model with Held-out-Para Partition to verify these assumptions. We only consider first-order methods, as second-order optimization with these larger models is impossible with our available computation.

Our results are plotted in Fig. 8. In Fig. 8(a) we compare the few-shot performance of direct fine-tuning on these three pre-trained models. On average, few-shot performance grows with models size, with a few exceptions such as QQP+T5-v1.1-Base and MRPC+Bart-Large. In Fig. 8(b-c) we plot the effect brought by upstream learning method for larger models. Except for FoMAML+T5-v1.1-Base<sup>8</sup>, upstream learning methods consistently improves few-shot performance on  $\mathcal{T}_{test}$ , which verifies that upstream learning methods we use are model-agnostic, and can be applied to larger models to further improve few-shot performance.

 Q6. Can we use pattern-exploiting training to replace direct fine-tuning to achieve even better performance?

Pattern-exploiting training (PET) is a novel method that formulate a target task into cloze-style questions (Schick and Schütze, 2020a,b; Gao et al., 2020). This approach narrows the gap between the masked language modeling objective during pre-training and downstream task fine-tuning, and therefore leads to more efficient transfer. PET is demonstrated to be effective with encoder models (e.g., RoBERTa), however, whether it is applicable to text-to-text models with auto-regressive decoders is underexplored to the best of our knowledge. In Q6, we study whether applying PET-style methods to text-to-text models is feasible, and whether combining the two methods further pushes the few-shot performance.

To align with the experiment settings in (Schick and Schütze, 2020a,b; Gao et al., 2020), we introduce a new task partition “Held-out-GLUE”, which uses non-GLUE classification tasks as  $\mathcal{T}_{train}$ , and GLUE tasks as  $\mathcal{T}_{test}$ . We use the top 3 patterns in (Gao et al., 2020) for each GLUE task, and use the

<sup>7</sup>T5-Base was trained on a mixture of downstream tasks during its pre-training; such practice strays from the purpose of our study. Therefore, we use T5-v1.1-Base model, which is trained with the C4 Corpus only.

<sup>8</sup>We observe instability in training loss during FoMAML training for T5-v1.1-Base.

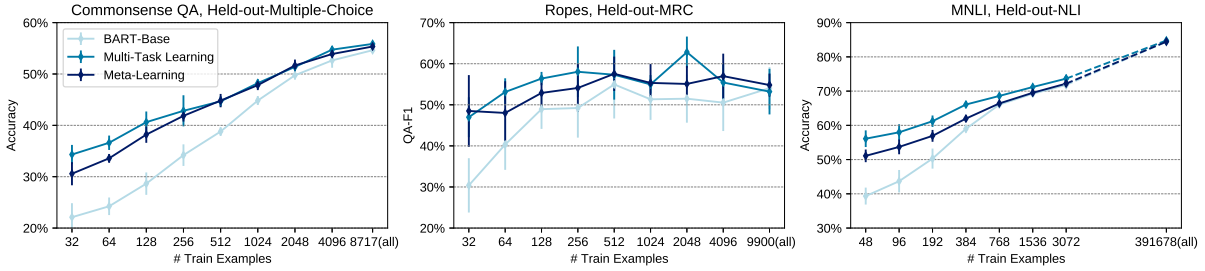


Figure 7: Performance comparisons in medium and high-resource scenarios. Benefits brought by upstream learning lasts in medium-resource scenarios.

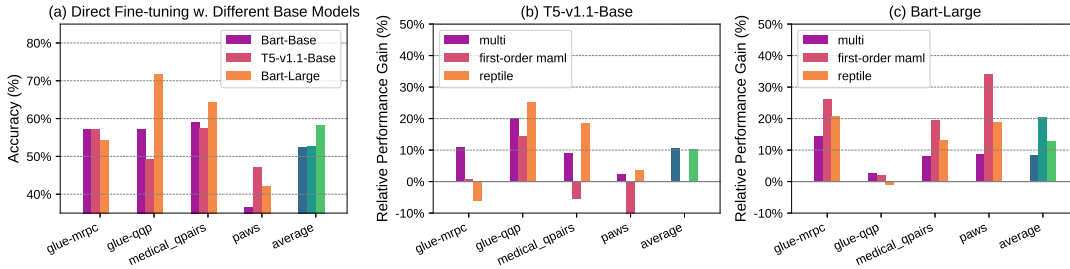


Figure 8: Extending upstream learning to larger pre-trained text-to-text models. (a) Absolute performance with direct fine-tuning with different pre-trained models. (b-c) Relative performance gain using upstream learning.

ensemble of the three models to produce the final prediction.

Since pattern-exploiting training is originally designed for encoder models (e.g., BERT/RoBERTa), we first tried two of its variants that adapts it to our auto-regressive transformer models. The first variant generates complete sentence, e.g., generate “The movie is great. A wonderful piece” from “The movie is great. A <mask> piece” for sentiment classification. The second variant generates only the word “wonderful”, from “The movie is great. A <mask> piece”. Though the first variant is more similar to the denoising pre-training objective of BART, we find the second variant to have better performance.

We then launch pattern-exploiting training using variant two with the original BART-Base models. We observe negative performance on average (leftmost blue bar in Fig. 9). Performance is improved with CoLA and MRPC, but not with the remaining GLUE tasks. We further launch experiments with/without pattern-exploiting training, with our upstream learning checkpoints. Still pattern-exploiting training leads to deteriorated performance on average.

We stop further investigation since this is out of the scope of our study. Still we believe it is important to identify the reasons and develop pattern-exploiting methods for auto-regressive models.

## D Reproducibility

**Implementation.** All our experiments are implemented with Huggingface Transformers<sup>9</sup> (Wolf et al., 2020). For higher-order optimization in the meta-learning approach optimization, we use higher library<sup>10</sup>. Our code has been uploaded in supplementary materials, and is also open-sourced at <https://github.com/INK-USC/CrossFit>.

**Hyper-parameters.** We mainly follow the practice in (Gao et al., 2020). During few-shot fine-tuning, we select the learning rate from  $\{1e-5, 2e-5, 5e-5\}$ , and the batch size from  $\{2, 4, 8\}$ , based on  $D_{dev}$  performance. We set the total number of updates to be 1000, number of warmup updates to be 100. We evaluate the model on  $D_{dev}$  every 100 steps.

**Infrastructure and Runtime.** Upstream learning are done with one single Quadro RTX 8000 (48GB). Upstream learning jobs finishes within 3 hours on average. Fine-tuning experiments are all done with one single GPU, with either NVIDIA Quadro GP100, NVIDIA Quadro RTX 8000, NVIDIA Quadro RTX 6000, NVIDIA GeForce RTX 1080 Ti, or NVIDIA GeForce RTX 2080 Ti, based on availability. Fine-tuning on one few-shot

<sup>9</sup><https://github.com/huggingface/transformers>

<sup>10</sup><https://github.com/facebookresearch/higher>

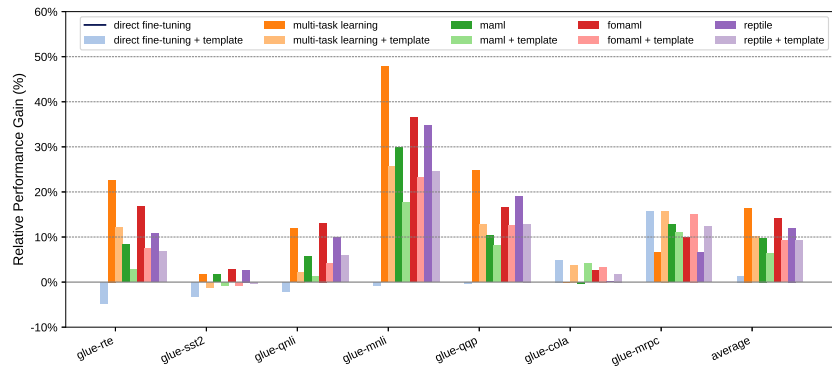


Figure 9: Combining upstream learning with pattern-exploiting training.

task (with hyperparameter tuning for all 5 random samples) takes approximately 4 hours on average.

**Number of Parameters.** BART-Base model contains 139 million parameters. T5-v1.1-Base model contains 246 million parameters. BART-Large model contains 406 million parameters.