

# LMTurk: Few-Shot Learners as Crowdsourcing Workers

Anonymous ACL submission

## Abstract

Vast efforts have been devoted to creating high-performance few-shot learners, i.e., large-scale pretrained language models (PLMs) that perform well with little downstream task training data. Training PLMs has incurred significant cost, but utilizing the few-shot learners is still challenging due to their enormous size. This work focuses on a crucial question: How to make effective use of these few-shot learners? We propose LMTurk, a novel approach that treats few-shot learners as crowdsourcing workers. The rationale is that crowdsourcing workers are in fact few-shot learners: They are shown a few illustrative examples to learn about a task and then start annotating. LMTurk employs few-shot learners built upon PLMs as workers. We show that the resulting annotations can be utilized to train models that solve the task well and are small enough to be deployable in practical scenarios. Altogether, LMTurk is an important step towards making effective use of current PLMs.

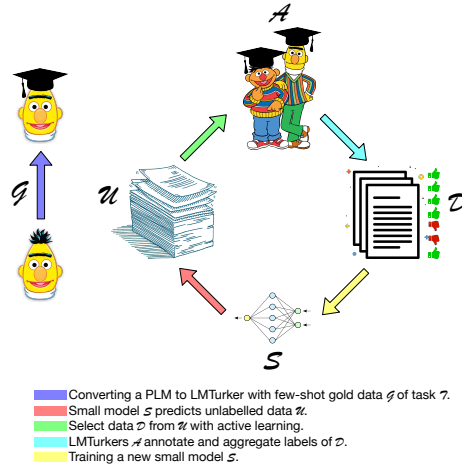


Figure 1: LMTurk overview; best viewed in color. We few-shot adapt PLMs to task  $\mathcal{T}$  (left) and then use them as crowdsourcing workers in active learning. We show that these PLM workers are effective in training a small model  $\mathcal{S}$  through a customized active learning loop (right). LMTurk is a novel way to take advantage of large-scale PLMs: It creates models small enough to be deployed in resource-limited real-world settings.

## 1 Introduction

Equipped with prolific linguistic features (Liu et al., 2019; Tenney et al., 2019; Belinkov and Glass, 2019; Rogers et al., 2020) and rich world knowledge (Petroni et al., 2019; Poerner et al., 2020; Kassner et al., 2021), large-scale pretrained language models (PLMs) have been shown to be versatile: They are now basic building blocks (Bommasani et al., 2021) of systems solving diverse NLP tasks in many languages (Wang et al., 2018, 2019; Hu et al., 2020; Xu et al., 2020; Khashabi et al., 2020; Park et al., 2021; Adelani et al., 2021).

Recent work shows that PLMs are effective *few-shot learners* (Brown et al., 2020a; Schick and Schütze, 2021b; Gao et al., 2021; Tam et al., 2021) through *priming* (Brown et al., 2020a; Tsimpoukelli et al., 2021) or *prompting* (Li and Liang, 2021; Liu et al., 2021b; Lester et al., 2021; Zhao and Schütze, 2021). Developing few-shot learn-

ers is crucial because current NLP systems require much more data than humans (Yin et al., 2020). Few-shot learners tend to perform well; however, they still fall behind systems trained with abundant data. Furthermore, the enormous size of PLMs hinders their deployment in practice. For example, it is challenging to fit the 11 billion T5-XXL (Raffel et al., 2020) model on a single regular GPU.

Our goal in this paper is to devise methods that make *more effective use of current few-shot learners*. This is crucial because an increasing number of gigantic few-shot learners are trained; how to use them effectively is thus an important question. In particular, we want an alternative to hard-to-deploy huge models. At the same time, we want to take full advantage of the PLMs’ strengths: Their versatility ensures wide applicability across tasks; their vast store of knowledge about language and the world (learned in pretraining) manifests in the data

061 efficiency of few-shot learners, reducing labor and  
062 time consumption in data annotation.

063 In this work, we propose **LMTurk**, Language  
064 **Model as mechanical Turk**. Our basic idea (see  
065 Figure 1) is that, for an NLP task  $\mathcal{T}$ , we treat few-  
066 shot learners as non-expert workers, resembling  
067 crowdsourcing workers that annotate resources for  
068 human language technology. We are inspired by the  
069 fact that we can view a crowdsourcing worker as a  
070 type of few-shot learner: A few examples demon-  
071 strating  $\mathcal{T}$  teach her enough about  $\mathcal{T}$  to conduct ef-  
072 fective annotation. For example, Snow et al. (2008)  
073 train workers with a few examples of annotating  
074 emotion; He et al. (2015) conduct short training  
075 sessions for workers before annotation; Lee et al.  
076 (2021) train workers with learning curricula.

077 Snow et al. (2008) pioneered crowdsourcing in  
078 NLP (Howe et al., 2006; Howe, 2008), motivated  
079 by the high cost of TreeBank annotation (Marcus  
080 et al., 1993; Miller et al., 1993). Crowdsourcing  
081 organizes human workers over the Web to annotate  
082 data. Workers need not be experts to be effective,  
083 resulting in reduced *per-label cost*. Active learning  
084 (Hachey et al., 2005; Felder and Brent, 2009) can  
085 be incorporated (Laws et al., 2011) to further de-  
086 crease annotation cost, by lowering the number of  
087 labels to be annotated. LMTurk treats PLM-based  
088 few-shot learners as non-expert workers that pro-  
089 duce training sets, which are then used to train a  
090 small machine learning model  $\mathcal{S}$  specialized for  
091  $\mathcal{T}$ . This scenario is analogous to active learning.  
092 We achieve two benefits: (i) low annotation cost  
093 because humans only need to annotate a few shots  
094 of data; (ii) solving practical NLP tasks with small  
095 models that are more real-world deployable.

096 LMTurk resonates with Laws et al. (2011)’s ear-  
097 lier idea of combining crowdsourcing and active  
098 learning. They consider human workers as “noisy  
099 annotators” while we explore the utilization of mod-  
100 ern NLP few-shot learners (built upon machine  
101 learning models) as workers – which have the ad-  
102 vantage of being free, instantly interactive, fast,  
103 responsive, and non-stopping.

104 Our **contributions**: (i) We propose LMTurk, a  
105 method that uses few-shot learners as crowdsourc-  
106 ing workers. Figure 1 shows the overview of LM-  
107 Turk. (ii) We vary an array of important design  
108 choices, identifying strengths and weaknesses of  
109 LMTurk. (iii) Unlike much work on active learning  
110 in a synthetic oracle setting, we develop methods  
111 for handling the varying quality of annotation that

112 does not come from an oracle. (iv) We extensively  
113 evaluate LMTurk on five datasets, showing that  
114 LMTurk can guide a small model  $\mathcal{S}$  to progres-  
115 sively improve on  $\mathcal{T}$ .  $\mathcal{S}$  can then be deployed in  
116 practical scenarios. (v) This is the first work show-  
117 ing that few-shot learners give rise to effective NLP  
118 models through crowdsourcing and active learning  
119 – with the benefits of low annotation cost and prac-  
120 tical deployability.

## 121 2 Related Work

**Few-shot learners in NLP.** Significant progress  
122 has been made in developing (Devlin et al., 2019;  
123 Peters et al., 2018; Yang et al., 2019; Brown et al.,  
124 2020b), understanding (Liu et al., 2019; Tenney  
125 et al., 2019; Belinkov and Glass, 2019; Hewitt and  
126 Liang, 2019; Hewitt and Manning, 2019; Zhao  
127 et al., 2020a; Rogers et al., 2020), and utilizing  
128 (Houlsby et al., 2019; Zhao et al., 2020b; Brown  
129 et al., 2020b; Li and Liang, 2021; Schick and  
130 Schütze, 2021a; Lester et al., 2021; Mi et al.,  
131 2021a) PLMs. Brown et al. (2020b), Schick and  
132 Schütze (2021a), and Liu et al. (2021b) show that  
133 PLMs can serve as data-efficient few-shot learners,  
134 through priming or prompting (Liu et al., 2021a).  
135 For example, GPT3 achieves near state-of-the-art  
136 performance on COPA (Roemmele et al., 2011)  
137 with only 32 annotated data.

138 However, little to no work discusses or explores  
139 the actual *practical utility* of these few-shot learn-  
140 ers. We aim to develop effective methods of utiliz-  
141 ing them in practical scenarios.

142 **Crowdsourcing** has a long history in human  
143 language technology (Alonso et al., 2008; Callison-  
144 Burch, 2009; Trautmann et al., 2020); specialized  
145 workshops were organized (Callison-Burch and  
146 Dredze, 2010; Paun and Hovy, 2019). It has numer-  
147 ous applications (Yuen et al., 2011), but we focus  
148 on its application as voting systems. To reduce *per-*  
149 *label* cost, crowdsourcing organizes non-expert hu-  
150 man workers distributed across the Web for annota-  
151 tion, instead of employing linguistic experts (Jami-  
152 son and Gurevych, 2015; Bhardwaj et al., 2019;  
153 Nangia et al., 2021). Snow et al. (2008) show  
154 that averaging ten crowdsourced labels matches  
155 an expert-level label for recognizing textual entail-  
156 ment (Dagan et al., 2006). Paun et al. (2018) show  
157 that incorporating structure in annotation models is  
158 important. Measuring label disagreements is also  
159 crucial (Dumitrache et al., 2021).

160 LMTurk utilizes NLP few-shot learners as non-  
161

expert workers. The few-shot training data can be viewed as the examples shown to humans before annotating. The process is free, fast, responsive, and non-stopping.

**Active learning** (AL; Cohn et al. (1996); Settles (2009)) strives to reduce *the number of examples* to be annotated via identifying informative examples with acquisition functions. Settles and Craven (2008) evaluate AL algorithms for sequence labeling. Zhang et al. (2017); Shen et al. (2017); Sidhant and Lipton (2018) apply AL to deep neural networks. Simpson and Gurevych (2018) devise a scalable Bayesian preference learning method for identifying convincing arguments. Lee et al. (2020) propose to consider user feedback in AL systems. Ein-Dor et al. (2020) explore AL for BERT. Schröder and Niekler (2020) review text classification with AL. Liang et al. (2020); Margatina et al. (2021) integrate contrastive learning into AL. Zhang and Plank (2021) identify examples with datamap (Swayamdipta et al., 2020).

We incorporate AL in LMTurk to reduce the amount of examples to be annotated by PLMs, reducing the computational cost of running several inference passes. This contributes to a more environmentally friendly (Strubell et al., 2019; Schwartz et al., 2020; Patterson et al., 2021) scenario.

Perhaps closest to our work, Yoo et al. (2021) conduct data augmentation via priming GPT3 and Wang et al. (2021) mix human- and GPT3-annotated data, focusing on cost analysis. GPT3 is not free.<sup>1</sup> Also, strategies of priming GPT3 may not generalize well to other PLMs.<sup>2</sup> In this work, we prompt publicly available free PLMs. This also makes the process more flexible; e.g., the PLM can be updated with gradient descent.

### 3 LMTurk

#### 3.1 Training few-shot learners

We first adapt a PLM to task  $\mathcal{T}$  with a few-shot human-labeled gold dataset  $\mathcal{G} = \{\mathcal{G}_{train}; \mathcal{G}_{dev}\}$  of  $\mathcal{T}$ . This procedure mimics one of the initial but crucial steps in crowdsourcing: A few example annotations are shown to the workers, demonstrating  $\mathcal{T}$ ; workers learn about the task and start annotating (Snow et al., 2008; He et al., 2015; Roit et al., 2020; Trautmann et al., 2020; Lee et al., 2021)

<sup>1</sup><https://beta.openai.com/pricing>

<sup>2</sup>For example, priming strategies have to adapt to GPT3’s maximum sequence length. However, maximum sequence length – as a hyperparameter – could vary across PLMs.

We achieve this adaptation through P-Tuning (Liu et al., 2021b). Taking movie review classification as an example, the goal is to associate a binary label  $y$  from  $\{-1, +1\}$  to an input sentence  $\mathbf{x} = (x_1, \dots, x_n)$  where  $x_i$  refers to a token. Unlike finetuning and its variants (Devlin et al., 2019; Houlsby et al., 2019; Zhao et al., 2020b) that train a classifier head, P-Tuning reformulates a sentence into a cloze-style query; the PLM is then requested to respond to the query with an answer selected from a list of candidates. Concretely, an input pair

$(\mathbf{x}, y) = (\text{“watching it leaves you giddy.”}, -1)$

is reformulated to:

“[v] watching it leaves you giddy. It is [MASK].”

in which the underlined tokens are prompting words that give the model a hint about  $\mathcal{T}$ . “[v]” – whose trainable embedding vector is randomly initialized – is a prompting token injecting extra free parameters. The PLM is then requested to pick a word from {“bad”, “good”} to fill in the position of “[MASK]”. A mapping {“bad”  $\rightarrow$  -1, “good”  $\rightarrow$  +1} is used to transform the selected answer to a label such that standard evaluation measures like accuracy can be computed. Prompting has been shown to effectively adapt a PLM to  $\mathcal{T}$  with only a few annotations; see (Liu et al., 2021a) for a comprehensive review of prompting. We refer to a PLM adapted to  $\mathcal{T}$  as an **LMTurker**  $A$ .

We select prompting words and mappings based on the small development set  $\mathcal{G}_{dev}$ . §4.2 provides details on prompting and datasets.

#### 3.2 Aggregating annotations

Individual workers are subject to biases (Snow et al., 2008); therefore, crowdsourcing often collects labels from several workers (Yuen et al., 2011) for an example  $\mathbf{x}$  and then aggregates them for quality control (Alonso et al., 2008). It is straightforward to obtain a group of LMTurkers  $\mathcal{A} = \{A_1, \dots, A_k\}$ , by adapting the PLM to  $\mathcal{T}$  with  $k$  different prompts. A querying sentence  $\mathbf{x}$  is then annotated by every LMTurker, resulting in a list of labels  $\mathbf{y} = [y_1, \dots, y_k]$ . We evaluate different methods aggregating  $\mathbf{y}$  to a single label  $\hat{y}$ .

**BestWorker.** Among the  $k$  LMTurkers, we pick the one performing best on the dev set  $\mathcal{G}_{dev}$ .

**MajorityVoting.** We select the most frequent label in  $\mathbf{y} = [y_1, \dots, y_k]$  as  $\hat{y}$ .

To estimate an LMTurker’s confidence on label  $y_i$ , we compare the logits<sup>3</sup> computed by the PLM:

$$y_i = \arg \max(\text{logit}(y^1), \dots, \text{logit}(y^N)),$$

where  $N$  refers to the label set size, e.g.,  $N=2$  for  $y$  from  $\{-1, +1\}$ . We evaluate several methods of aggregating annotations according to PLM logits:

**LogitVoting.** We average the logits from all  $k$  LMTurkers  $\{A_1, \dots, A_k\}$  to compute  $\hat{y}$ :

$$\hat{y} = \arg \max(\frac{1}{k} \sum_{i=1}^k \text{logit}(y_i^1), \dots, \frac{1}{k} \sum_{i=1}^k \text{logit}(y_i^N)),$$

**WeightedLogitVoting.** We use LMTurkers’ performance on  $\mathcal{G}_{dev}$  to weight their logits and then aggregate the predictions:

$$\hat{y} = \arg \max(\sum_{i=1}^k w_i \text{logit}(y_i^1), \dots, \sum_{i=1}^k w_i \text{logit}(y_i^N))$$

$$w_i = f(A_i, \mathcal{G}_{dev}) / \sum_{i=1}^k f(A_i, \mathcal{G}_{dev})$$

where  $f(A_i, \mathcal{G}_{dev})$  is the performance of the  $i$ th LMTurker  $A_i$  on  $\mathcal{G}_{dev}$ .

We collect and aggregate annotations from five LMTurkers, i.e., we use  $k=5$  in our experiments.

### 3.3 Training a small model $\mathcal{S}$

After adapting LMTurkers to  $\mathcal{T}$  through prompting with the few-shot gold dataset  $\mathcal{G}$ , we next train a small model  $\mathcal{S}$  specialized to solve  $\mathcal{T}$ . Though large PLMs are versatile and strong performers, training and inference are faster and more efficient for small models: They are more deployable in resource-restricted scenarios, e.g., on edge devices (Jiao et al., 2020).

We mimic pool-based active learning (AL; Settles (2009)) to train  $\mathcal{S}$ . The motivation is to avoid frequent querying of LMTurkers  $\mathcal{A}$  because energy and time consumption of PLM inference is costly when the number of queries and  $|\mathcal{A}|$  are large.

Concretely, pool-based AL assumes a large collection of unlabeled data  $\mathcal{U} = \{\mathbf{x}_1, \dots, \mathbf{x}_M\}$  for  $\mathcal{T}$ .  $\mathcal{S}$  is first trained with  $\mathcal{G} = \{\mathcal{G}_{train}; \mathcal{G}_{dev}\}$ . After that, a group of examples  $\mathcal{B}$  from  $\mathcal{U}$  is sampled (c.f. §3.3.1), which LMTurkers annotate. Next, the annotated and aggregated examples  $\mathcal{B}'$  are concatenated with  $\mathcal{G}$  to train  $\mathcal{S}$ . The procedure is repeated iteratively, such that the training data for  $\mathcal{S}$  keeps expanding. We denote as  $\mathcal{S}^j$  the model trained after the  $j$ th iteration. Note that  $\mathcal{S}$  is trained from scratch in each iteration (Cohn et al., 1994).

<sup>3</sup>Calibration can be conducted to further improve the estimation (Guo et al., 2017). We leave this to future work.

#### 3.3.1 AL acquisition function

At the beginning of the  $j$ th iteration, a straightforward strategy of sampling  $\mathcal{B}$  from  $\mathcal{U}$  is **random sampling**. AL promises to select a more informative  $\mathcal{B}$  such that the trained  $\mathcal{S}^j$  performs better, under the same budget. These strategies – or *acquisition functions* – rely on  $\mathcal{S}^{j-1}$ , i.e.,  $\mathcal{S}$  from the previous iteration:  $\mathcal{S}^{j-1}$  is employed to infer  $\mathcal{U}$  to obtain labels and logits  $\mathcal{P}^{j-1} = \{(y_1, \mathbf{c}_1), \dots, (y_M, \mathbf{c}_M)\}$ ; each  $\mathbf{c}_i$  contains the logits of the  $N$  labels;  $y_i = \arg \max(\mathbf{c}_i)$ . We explore two common AL acquisition functions: Entropy (Roy and McCallum, 2001) and LeastConfident (Lewis and Gale, 1994).

**Entropy** selects from  $\mathcal{P}^{j-1}$  examples with the largest prediction entropy, computed on  $\mathbf{c}$ . Large entropy of an example  $\mathbf{x}$  implies that  $\mathcal{S}^{j-1}$  is unsure about which label to select;  $\mathbf{x}$  is then a query made to LMTurkers to obtain its annotation  $\hat{y}$ .  $(\mathbf{x}, \hat{y})$  is subsequently added to  $\mathcal{G}_{train}$  for training  $\mathcal{S}^j$ .

**LeastConfident** selects from  $\mathcal{P}^{j-1}$  examples for which the maximum logit in  $\mathbf{c}$  is the smallest. Selected examples are then annotated and added to  $\mathcal{G}_{train}$  for training  $\mathcal{S}^j$ .

Our AL setup is fairly standard, both in terms of acquisition functions and iterative enlargement by new sampled data  $\mathcal{B}$  at iteration  $j$  labeled by  $\mathcal{S}^{j-1}$ .

#### 3.3.2 Considering annotation quality

As in any realistic AL scenario, annotations are not perfect: LMTurkers do not score perfectly on  $\mathcal{T}$ . So *annotation quality of LMTurkers needs to be taken into consideration before training  $\mathcal{S}^j$* . Denoting the training data of  $\mathcal{S}^j$  as  $\mathcal{D}^j$ , we explore a strategy of processing  $\mathcal{D}^j$ , based on LMTurker logits  $\mathbf{l}$ .

**InstanceTresholding.** We preserve examples  $(\mathbf{x}, \hat{y}, \mathbf{l}) \in \mathcal{D}^j$  for which entropy computed on  $\mathbf{l}$  is smallest.  $\mathcal{G}^{train}$  is always preserved because it is human-labeled gold data. Note that this is different from the strategy of sampling  $\mathcal{B}$ , where we select from  $\mathcal{P}^{j-1}$  examples to which  $\mathcal{S}^{j-1}$  is most unsure (computed with  $\mathbf{c}$ ). We evaluate<sup>4</sup> the effectiveness of processing  $\mathcal{D}^j$  before training  $\mathcal{S}^j$  in §5.6.

### 3.4 Summary of LMTurk

LMTurk can be viewed as intermediate between self training (Yarowsky, 1995; Lee et al., 2013; Mi et al., 2021b) and AL. Unlike self training, *external* models provide labels to  $\mathcal{S}$ . Different from the

<sup>4</sup>Motivated by Wang et al. (2017), we also investigate the effectiveness of weighting training examples. However, we do not observe noticeable improvements of task performance. We list more details in §E.

artificial setup used in many AL experiments, the provided labels *do not have oracle quality*; so  $\mathcal{S}$  must use the annotations more carefully. We next conduct experiments investigating the effectiveness of LMTurk.

## 4 Datasets and Setup

### 4.1 Dataset

We evaluate LMTurk on five datasets: Binary (SST2) and fine-grained (five classes) sentiment classification (SST5) with the Stanford Sentiment TreeBank (Socher et al., 2013); news article topic classification with the AG’s News Corpus (AGNews; Zhang et al. (2015)); recognizing textual entailment (RTE; Dagan et al. (2006)); assessing linguistic acceptability (CoLA; Warstadt et al. (2019)). Appendix §A reports dataset statistics. SST2/SST5 and AGNews are widely used in crowdsourcing and AL (Laws et al., 2011; Ein-Dor et al., 2020; Margatina et al., 2021; Zhang and Plank, 2021). RTE and CoLA assess the models’ ability to understand textual entailment and linguistic phenomena – as opposed to text categorization. We report Matthew’s correlation coefficient for CoLA and accuracy for the others (Wang et al., 2018).

**Few-shot datasets.** Recall LMTurk uses a small human-annotated dataset  $\mathcal{G} = \{\mathcal{G}_{train}; \mathcal{G}_{dev}\}$ . Denoting  $n$  as the number of shots *per class*, we sample  $\mathcal{G}_{train}^n$  and  $\mathcal{G}_{dev}^n$  for each of  $n \in \{8, 16, 32\}$ . For SST2, RTE, and CoLA, we use the train and dev sets of GLUE (Wang et al., 2018);  $\mathcal{G}_{train}^n$  and  $\mathcal{G}_{dev}^n$  are sampled from the train set; the dev set is used as the test set. For SST5 and AGNews, we use the official datasets;  $\mathcal{G}_{train}^n$  ( $\mathcal{G}_{dev}^n$ ) is sampled from the train (dev) set; we report performance on the test set. We repeat the sampling process with three random seeds.

### 4.2 Training setup

Brown et al. (2020b) show that large model size is necessary for strong few-shot performance. We use ALBERT-XXLarge-v2 (Lan et al., 2020) – of size 223M parameters – as our large PLM, which is adapted to be an LMTurker  $A$  of  $\mathcal{T}$  with  $\mathcal{G}$ . With parameter reuse, ALBERT-XXLarge-v2 outperforms larger models like the 334M BERT-large (Devlin et al., 2019). In contrast,  $\mathcal{S}$  must be small to be deployable in practical scenarios. We use TinyBERT-General-4L-312D (Jiao et al., 2020), which has 14.5M parameters.

We train – with prompting – the large PLM with

	Schick and Schütze (2021a,b)	Gao et al. (2021)	Ours
SST2	n/a	93.0±0.6	93.08±0.62
SST5	n/a	49.5±1.7	46.70±0.93
RTE	69.8	71.1±5.3	70.88±1.70
AGN.	86.3±0.0	n/a	87.71±0.07
CoLA	n/a	21.8±15.9	19.71±1.89

Table 1: LMTurkers achieve comparable few-shot performance with the literature. We refer to *PET* results in Schick and Schütze (2021a,b) and results of *Prompt-based FT (auto) + demonstrations* in Gao et al. (2021).

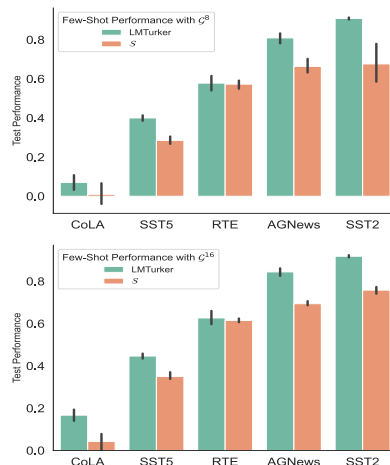


Figure 2: *Few-shot* test set performance of LMTurkers and  $\mathcal{S}$ . We use the few-shot gold datasets  $\mathcal{G}^8$  (top) and  $\mathcal{G}^{16}$  (bottom).  $\mathcal{G}^{32}$  results present similar trend; they are shown in Appendix §D.

$\mathcal{G}$  for 100 batch steps using batch size 16, AdamW (Loshchilov and Hutter, 2019) and learning rate 5e-4 with linear decay. We prompt the large PLM five times to obtain five LMTurkers; Appendix §C shows prompting details. At each iteration, we fine-tune  $\mathcal{S}$  for 20 epochs using batch size 32, Adam (Kingma and Ba, 2015) and learning rate 5e-5. Each experiment is run with three different random seeds. We use PyTorch (Paszke et al., 2019) and HuggingFace (Wolf et al., 2020).

## 5 Experiment

### 5.1 Few-shot performance (non-iterative)

We compare few-shot performance of LMTurkers and the small model  $\mathcal{S}$  when *only  $\mathcal{G}$  is used*. LMTurker performance is comparable to prior work as shown in Table 1.

Figure 2 compares performance of LMTurkers and  $\mathcal{S}$ . Appendix §B Table 3 reports numeric values. LMTurkers perform clearly better than  $\mathcal{S}$  on CoLA, SST5, AGNews, and SST2; e.g., for SST2, for train/dev size 16, LMTurker accuracy is 93.08% vs. 75.83% for  $\mathcal{S}$ . LMTurkers’ superiority over  $\mathcal{S}$

on RTE is modest. As an inference task, RTE is more challenging than classification (e.g., AGNews). We hypothesize that current few-shot learners require more data than  $\mathcal{G}^{32}$  to process difficult tasks better than  $\mathcal{S}$ . Scaling up to even larger PLMs is also a promising direction (Lester et al., 2021).

Overall, LMTurkers outperform  $\mathcal{S}$  with clear margins, evidencing that their annotations can serve as supervisions for training  $\mathcal{S}$ . We next conduct iterative training to improve performance of  $\mathcal{S}$  on  $\mathcal{T}$  with supervisions from LMTurkers.

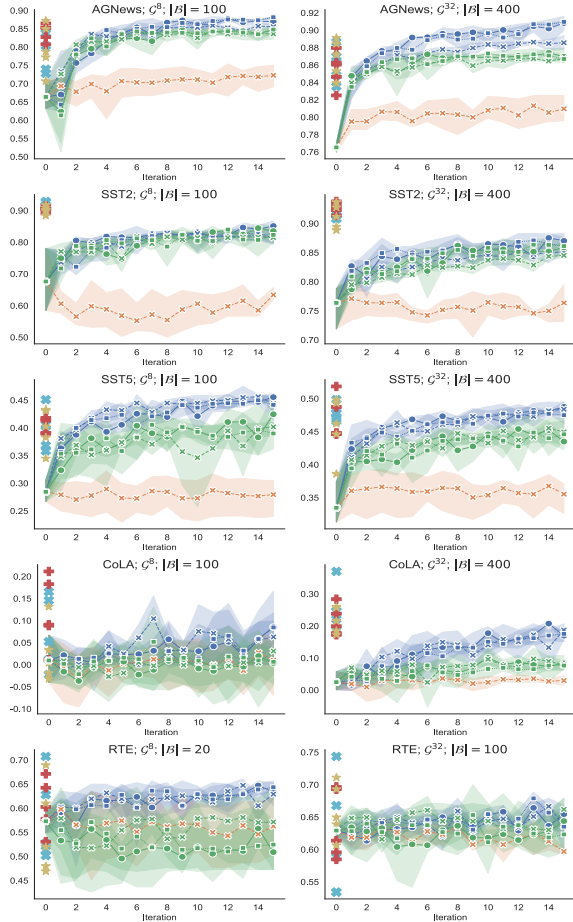


Figure 3: Improving  $\mathcal{S}$  with active learning (blue), self training (orange), and LMTurk (green). Free markers at step zero show LMTurker performances; colors distinguish random seeds. Three acquisition functions are: Entropy ( $\bullet$ ), LeastConfident ( $\blacksquare$ ), random sampling ( $\times$ ). At iteration  $j$ , each experiment is repeated three times; we show mean and standard deviation. Appendix Figure 10 visualizes more results.

## 5.2 Iterative training

We investigate the effectiveness of LMTurk by simulating scenarios analogous to active learning. Concretely, we compare three schemes of annotating the sampled data  $\mathcal{B}$  at each annotation iteration  $j$ :

- **Active learning (AL)**. We use  $\mathcal{B}$ 's gold labels to show how  $\mathcal{S}$  performs with expert annotations. Gold labels are ideal, but costly.
- **Self training (ST)**.  $\mathcal{S}^{j-1}$  (the model trained in the previous iteration) annotates  $\mathcal{B}$  (Yarowsky, 1995; Lee et al., 2013). ST trades supervision quality for annotation cost; no extra cost is introduced. Because there is no external supervision, ST is expected to be a baseline.
- **LMTurk**. We query the LMTurkers to annotate  $\mathcal{B}$ . LMTurkers are machine learning models, so there is no human labor. Based on the findings in Figure 2, LMTurker supervisions are expected to have better quality than those of ST. Yet LMTurk could fall behind AL because LMTurker labels are not gold labels.

When sampling  $\mathcal{B}$  from  $\mathcal{U}$  at each iteration  $j$ , we consider the strategies described in §3.3. We employ Random for all three schemes and Entropy/LeastConfident for AL/LMTurk. The latter two rely on  $\mathcal{S}^{j-1}$ . Regarding the number of sampled examples, we experiment with  $|\mathcal{B}|=100$  and  $|\mathcal{B}|=400$  for SST2, SST5, AGNews, CoLA. Due to RTE's small size, we use  $|\mathcal{B}|=20$  and  $|\mathcal{B}|=100$ . We run for 15 iterations of improving  $\mathcal{S}$ . To aggregate annotations from LMTurkers, we use Majority Voting (§3.2), which is widely used in crowdsourcing. See §5.3 for a comparison of aggregation methods.

Figure 3 compares AL, ST and LMTurk. ST (orange) noticeably helps  $\mathcal{S}$  to perform progressively better on AGNews, e.g., when comparing  $\mathcal{S}^{15}$  to  $\mathcal{S}^0$  shown in the first row, especially when  $|\mathcal{B}|=400$ . However, we do not identify clear improvements when looking at other tasks. Except for RTE- $\mathcal{G}^8$ , ST clearly falls behind AL and LMTurk. This inferior performance meets our expectation because there is no external supervision assisting  $\mathcal{S}$  to perform better on  $\mathcal{T}$ . In what follows, we omit ST for clearer visualization and discussion.

AL (blue) performs the best in most experiments. However, this comes with extra costs that are not negligible: *At each iteration*, human annotators need to annotate 100–400 sentences.

LMTurk (green) holds a position between AL and ST on AGNews, SST2, SST5, and CoLA. Somehow surprisingly, LMTurk performs almost comparably to AL on SST2. Unlike AL, LMTurk requires very little human labor; the only human annotation throughout the entire process is the few-shot gold dataset  $\mathcal{G}$ . In contrast, AL has high human annotation cost, e.g., 1000–4000 examples by iter-

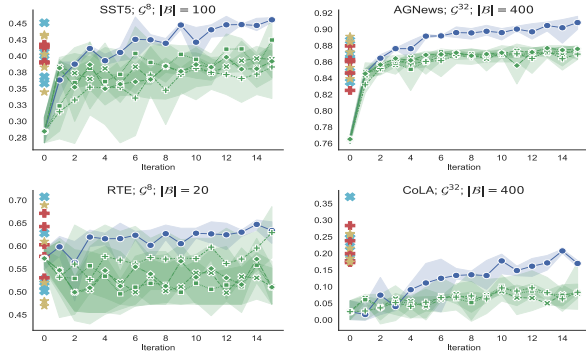


Figure 4: Comparing strategies of aggregating LM-Turker annotations. We compare LMTurk (green) with AL (blue). Strategies: LogitVoting (✕), MajorityVoting (■), WeightedLogitVoting (◆), BestWorker (+). AL uses gold labels without aggregation (●).

ation ten. LMTurk also shows clear performance improvements over ST.

Results on RTE are noisy; we conjecture this is due to its very small test set (277 examples). We do not observe performance improvement of  $\mathcal{S}$  along the iterations in experiment RTE- $\mathcal{G}^{32}$ - $|\mathcal{B}|=100$ , likely due to saturated task performance: TinyBERT-General-4L-312D ( $\mathcal{S}$ ) achieves 66.6% on RTE for the full train set (Jiao et al., 2020).

**Comparing sampling strategies.** Entropy (●) and LeastConfident (■) outperform random sampling (✕) in AGNews and SST2 with noticeable margins – for both AL and LMTurk, especially when  $|\mathcal{B}|=400$ . They also surpass random sampling when using LMTurk for SST5 and CoLA with  $\mathcal{G}^8$ . In other words, Entropy and LeastConfident assist LMTurk to achieve the same performance as of using random sampling, but with fewer annotations. For example in AGNews- $\mathcal{G}^8$ - $|\mathcal{B}|=100$ , LeastConfident at iteration six already achieves comparable performance as random sampling at iteration eleven. This is economically and environmentally beneficial because the number of queries made to LMTurkers, i.e., the cost of running inference passes on the array of large PLMs, is significantly reduced.

Overall, we show that LMTurk can be used to create datasets for training a specialized model  $\mathcal{S}$  of solving  $\mathcal{T}$  in practical scenarios. To reduce computational cost, we use only Entropy in what follows.

### 5.3 Design choice 1: Aggregation strategies

Figure 4 compares effectiveness of different strategies of aggregating LMTurker annotations (§3.2).

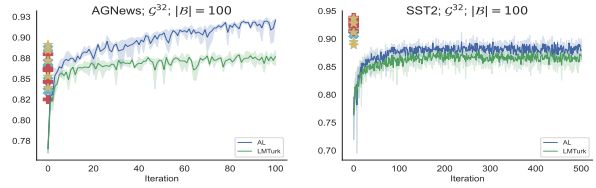


Figure 5: Running more iterations of improving  $\mathcal{S}$  with AL and LMTurk.

Looking at SST5 and AGNews results (top two images), we observe that committee-style aggregation (LogitVoting (✕), MajorityVoting (■), and WeightedLogitVoting (◆)) generally outperforms BestWorker (+), which simply relies on the LMTurker performing best on  $\mathcal{G}_{dev}$ . LMTurkers perform well on these two datasets as shown by the free markers at iteration zero; ensembling their predictions results in higher-quality datasets.

In contrast, BestWorker (+) has stellar performance on RTE (bottom-left), outperforming committee-style aggregation. Note that even LMTurkers do not perform really well in this experiment, as shown by the free markers at iteration zero – some LMTurkers even perform worse than  $\mathcal{S}$ . Ensembling these low-quality annotations seems a worse option than simply relying on the best LMTurker. For CoLA, we observe comparable performance of different aggregation strategies.

### 5.4 Design choice 2: More iterations

We hypothesize that AL performance is an upper bound for performance when  $\mathcal{S}$  is trained with LMTurker annotations – recall that the AL annotations are gold labels. Figure 5 compares AL and LMTurk when running 100 iterations of improving  $\mathcal{S}$  on AGNews and 500 iterations on SST2 (aggregation: WeightedLogitVoting). As expected, AL outperforms LMTurk because the pool of human-annotated data expands. The performance of  $\mathcal{S}$  progressively approaches that of the LMTurkers; LMTurk performs comparably to AL in SST2, however, no human labor is required.

### 5.5 Design choice 3: Distilling logits

We can view LMTurk as a kind of distillation (Hinton et al., 2015): The ability of LMTurkers to solve  $\mathcal{T}$  is progressively transferred to  $\mathcal{S}$ . In this section, we explore the utility of distillation: We train  $\mathcal{S}$  with predicted logits instead of discrete labels from LMTurkers. Concretely, we train  $\mathcal{S}$  by reducing the KL divergence between its predicted

513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552

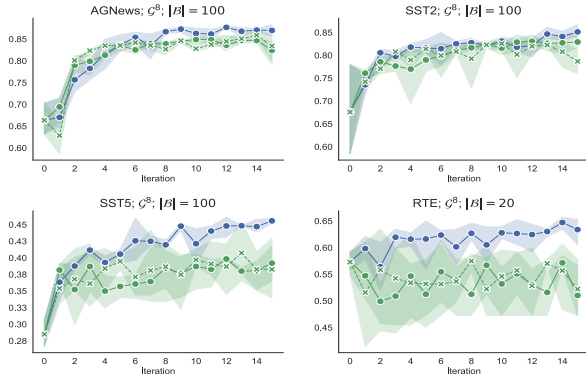


Figure 6: Performance of AL and LMTurk with discrete labels (●) vs. with KL divergence (✕).

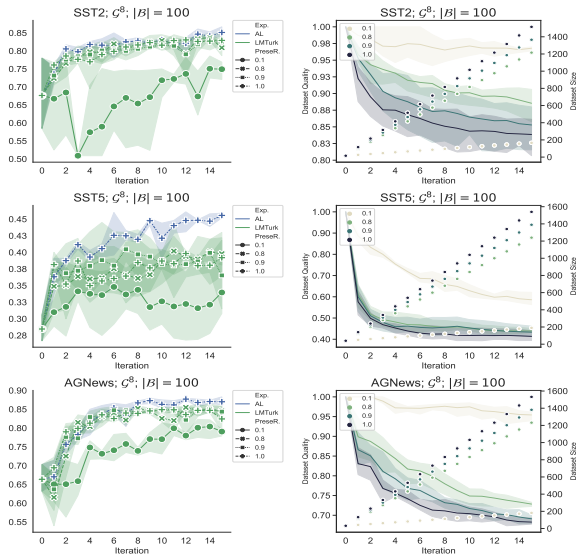


Figure 7: Training  $\mathcal{S}$  with examples for which LMTurkers have low entropy. We report performance of  $\mathcal{S}$  (left), number and quality (measured by accuracy) of the preserved examples (right) at each iteration.

probability distribution (over the label set) and the probability distribution from LMTurkers.

Figure 6 shows that training  $\mathcal{S}$  with KL divergence noticeably improves over discrete labels on AGNews and SST5. This is expected: AGNews and SST5 have larger label set size (four and five) such that the probability distribution over the label set is more informative than that of the binary classification tasks SST2 and RTE.

## 5.6 Design choice 4: Quality-based filtering

One key difference between AL and LMTurk is that LMTurkers are not oracles: Their labels are not perfect. Hence, it is reasonable to consider processing the training data, denoted as  $\mathcal{D}^j$ , for  $\mathcal{S}^j$ , instead of using it indiscriminately as in AL.

**InstanceTresholding** (§3.3.2) preserves annotations in  $\mathcal{D}^j$  for which LMTurkers have the smallest entropy. Concretely, we rank all annotations  $(\mathbf{x}, \hat{y}, \mathbf{l}) \in \mathcal{D}^j$  by  $\text{entropy}(\mathbf{l})$  and then keep the  $\tau$  percent smallest. Note that we always preserve the human-labeled few-shot data  $\mathcal{G}^{\text{train}}$ . We experiment with  $\tau \in \{10\%, \dots, 90\%, 100\%\}$ .

Figure 7 left shows the performance of  $\mathcal{S}$ ; Figure 7 right tracks the status of  $\mathcal{D}^j$ . To measure quality, we compute the accuracy of LMTurker annotations on  $\mathcal{D}^j$  (compared to gold labels); see the lineplots and the left y-axis. We also report the size of  $\mathcal{D}^j$  as scatter plots (right y-axis).

We observe that  $\tau=10\%$ , i.e., keeping only the 10% most certain examples, gives the worst performance. This is most obvious at iteration three for SST2: The performance drops to near the majority baseline ( $\approx 50\%$ ). This is because  $\mathcal{D}^3$  is small and unbalanced: It has eight negative (from  $\mathcal{G}^{\text{train}}$ ) and 38 positive examples. However, using all the LMTurker annotations ( $\tau=100\%$ ) may not be optimal either. This is noticeable when looking at SST5:  $\tau=90\%$  and  $\tau=80\%$  are better options.

We see that there is a tradeoff between  $\mathcal{D}^j$ 's quality and size from Figure 7 right. Being conservative, i.e., preserving only a handful of annotations from LMTurkers, results in a small, but high-quality  $\mathcal{D}^j$ ; using all the annotations indiscriminately leads to a large  $\mathcal{D}^j$  with low quality.

This experiment highlights a key difference between LMTurk and AL: LMTurker annotations are not perfect and taking the annotation quality into consideration when training  $\mathcal{S}$  is crucial.

## 6 Conclusion

In this work, our focus is the research question: *How to make effective use of current few-shot learners?* We propose LMTurk, a simple yet effective method that considers PLM-based few-shot learners as non-expert annotators in crowdsourcing; active learning is incorporated to reduce the cost of annotation. We further show that processing the annotations from LMTurkers can be beneficial.

Future work may combine LMTurker annotations with human annotators in a human-in-the-loop setup (Monarch, 2021) to increase the overall utility of invested resources (Bai et al., 2021). Applying LMTurk to multilingual few-shot learners (Zhao et al., 2021; Winata et al., 2021; Lin et al., 2021) is also promising.



## References

- 618 David Ifeoluwa Adelani, Jade Abbott, Graham Neu-  
619 big, Daniel D'souza, Julia Kreutzer, Constantine Lign-  
620 os, Chester Palen-Michel, Happy Buzaaba, Shruti  
621 Rijhwani, Sebastian Ruder, Stephen Mayhew, Is-  
622 rael Abebe Azime, Shamsuddeen H. Muhammad,  
623 Chris Chinenye Emezue, Joyce Nakatumba-Nabende,  
624 Perez Ogayo, Aremu Anuoluwapo, Catherine Gitau,  
625 Derguene Mbaye, Jesujoba Alabi, Seid Muhie Yi-  
626 mam, Tajuddeen Rabi'u Gwadabe, Ignatius Ezeani,  
627 Rubungo Andre Niyongabo, Jonathan Mukiibi, Ver-  
628 rah Otiende, Iroro Orife, Davis David, Samba Ngom,  
629 Tosin Adewumi, Paul Rayson, Mofetoluwa Adeyemi,  
630 Gerald Muriuki, Emmanuel Anebi, Chiamaka Chuk-  
631 wuneke, Nkiruka Odu, Eric Peter Wairagala, Samuel  
632 Oyerinde, Clemencia Siro, Tobius Saul Bateesa,  
633 Temilola Oloyede, Yvonne Wambui, Victor Akin-  
634 ode, Deborah Nabagereka, Maurice Katusiime, Ayo-  
635 dele Awokoya, Mouhamadane MBOUP, Dibora Ge-  
636 breyohannes, Henok Tilaye, Kelechi Nwaike, De-  
637 gaga Wolde, Abdoulaye Faye, Blessing Sibanda, Ore-  
638 vaoghene Ahia, Bonaventure F. P. Dossou, Kelechi  
639 Ogueji, Thierno Ibrahima DIOP, Abdoulaye Diallo,  
640 Adewale Akinfaderin, Tendai Marengereke, and Sa-  
641 lomey Osei. 2021. [MasakhaNER: Named Entity  
642 Recognition for African Languages](#). *Transactions  
643 of the Association for Computational Linguistics*,  
644 9:1116–1131.
- 645 Omar Alonso, Daniel E. Rose, and Benjamin Stewart.  
646 2008. [Crowdsourcing for relevance evaluation](#). *SI-  
647 GIR Forum*, 42(2):9–15.
- 648 Fan Bai, Alan Ritter, and Wei Xu. 2021. [Pre-train or an-  
649 notate? domain adaptation with a constrained budget](#).  
650 In *Proceedings of the 2021 Conference on Empirical  
651 Methods in Natural Language Processing*, pages  
652 5002–5015, Online and Punta Cana, Dominican Re-  
653 public. Association for Computational Linguistics.
- 654 Yonatan Belinkov and James Glass. 2019. [Analysis  
655 methods in neural language processing: A survey](#).  
656 *Transactions of the Association for Computational  
657 Linguistics*, 7:49–72.
- 658 Sangnie Bhardwaj, Samarth Aggarwal, and Mausam  
659 Mausam. 2019. [CaRB: A crowdsourced benchmark  
660 for open IE](#). In *Proceedings of the 2019 Confer-  
661 ence on Empirical Methods in Natural Language Pro-  
662 cessing and the 9th International Joint Conference  
663 on Natural Language Processing (EMNLP-IJCNLP)*,  
664 pages 6262–6267, Hong Kong, China. Association  
665 for Computational Linguistics.
- 666 Rishi Bommasani, Drew A Hudson, Ehsan Adeli,  
667 Russ Altman, Simran Arora, Sydney von Arx,  
668 Michael S Bernstein, Jeannette Bohg, Antoine Bosse-  
669 lut, Emma Brunskill, et al. 2021. On the opportuni-  
670 ties and risks of foundation models. *arXiv preprint  
671 arXiv:2108.07258*.
- 672 Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
673 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
674 Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
Askill, Sandhini Agarwal, Ariel Herbert-Voss,  
Gretchen Krueger, Tom Henighan, Rewon Child,  
Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens  
Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma-  
teusz Litwin, Scott Gray, Benjamin Chess, Jack  
Clark, Christopher Berner, Sam McCandlish, Alec  
Radford, Ilya Sutskever, and Dario Amodei. 2020a.  
[Language models are few-shot learners](#). In *Ad-  
vances in Neural Information Processing Systems*,  
volume 33, pages 1877–1901. Curran Associates,  
Inc. 675
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
Askill, Sandhini Agarwal, Ariel Herbert-Voss,  
Gretchen Krueger, Tom Henighan, Rewon Child,  
Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens  
Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma-  
teusz Litwin, Scott Gray, Benjamin Chess, Jack  
Clark, Christopher Berner, Sam McCandlish, Alec  
Radford, Ilya Sutskever, and Dario Amodei. 2020b.  
[Language models are few-shot learners](#). In *Ad-  
vances in Neural Information Processing Systems*,  
volume 33, pages 1877–1901. Curran Associates,  
Inc. 686
- Chris Callison-Burch. 2009. [Fast, cheap, and creative:  
Evaluating translation quality using Amazon's Me-  
chanical Turk](#). In *Proceedings of the 2009 Confer-  
ence on Empirical Methods in Natural Language  
Processing*, pages 286–295, Singapore. Association  
for Computational Linguistics. 687
- Chris Callison-Burch and Mark Dredze, editors. 2010.  
*Proceedings of the NAACL HLT 2010 Workshop on  
Creating Speech and Language Data with Amazon's  
Mechanical Turk*. Association for Computational Lin-  
guistics, Los Angeles. 688
- David Cohn, Les Atlas, and Richard Ladner. 1994. Im-  
proving generalization with active learning. *Machine  
learning*, 15(2):201–221. 689
- David A Cohn, Zoubin Ghahramani, and Michael I  
Jordan. 1996. Active learning with statistical models.  
*Journal of artificial intelligence research*, 4:129–145. 690
- Ido Dagan, Oren Glickman, and Bernardo Magnini.  
2006. The pascal recognising textual entailment chal-  
lenge. In *Machine Learning Challenges. Evaluating  
Predictive Uncertainty, Visual Object Classification,  
and Recognising Tectual Entailment*, pages 177–190,  
Berlin, Heidelberg. Springer Berlin Heidelberg. 691
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and  
Kristina Toutanova. 2019. [BERT: Pre-training of  
deep bidirectional transformers for language under-  
standing](#). In *Proceedings of the 2019 Conference of  
the North American Chapter of the Association for  
Computational Linguistics: Human Language Tech-  
nologies, Volume 1 (Long and Short Papers)*, pages  
4171–4186, Minneapolis, Minnesota. Association for  
Computational Linguistics. 692

732	Anca Dumitrache, Oana Inel, Benjamin Timmermans, Carlos Ortiz, Robert-Jan Sips, Lora Aroyo, and Chris Welty. 2021. Empirical methodology for crowdsourcing ground truth. <i>Semantic Web</i> , 12(3):1–19.	787
733		788
734		789
735		
736	Liat Ein-Dor, Alon Halfon, Ariel Gera, Eyal Shnarch, Lena Dankin, Leshem Choshen, Marina Danilevsky, Ranit Aharonov, Yoav Katz, and Noam Slonim. 2020. <a href="#">Active Learning for BERT: An Empirical Study</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 7949–7962, Online. Association for Computational Linguistics.	790
737		791
738		792
739		793
740		794
741		795
742		
743		
744	Richard M Felder and Rebecca Brent. 2009. Active learning: An introduction. <i>ASQ higher education brief</i> , 2(4):1–5.	796
745		797
746		798
747	Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. <a href="#">Making pre-trained language models better few-shot learners</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 3816–3830, Online. Association for Computational Linguistics.	799
748		800
749		
750		801
751		802
752		803
753		804
754		805
755	Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. 2017. On calibration of modern neural networks. In <i>International Conference on Machine Learning</i> , pages 1321–1330. PMLR.	806
756		807
757		808
758		
759	Ben Hachey, Beatrice Alex, and Markus Becker. 2005. <a href="#">Investigating the effects of selective sampling on the annotation task</a> . In <i>Proceedings of the Ninth Conference on Computational Natural Language Learning (CoNLL-2005)</i> , pages 144–151, Ann Arbor, Michigan. Association for Computational Linguistics.	809
760		810
761		811
762		812
763		813
764		814
765	Luheng He, Mike Lewis, and Luke Zettlemoyer. 2015. <a href="#">Question-answer driven semantic role labeling: Using natural language to annotate natural language</a> . In <i>Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing</i> , pages 643–653, Lisbon, Portugal. Association for Computational Linguistics.	815
766		
767		
768		
769		
770		
771		
772	John Hewitt and Percy Liang. 2019. <a href="#">Designing and interpreting probes with control tasks</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2733–2743, Hong Kong, China. Association for Computational Linguistics.	816
773		817
774		818
775		819
776		820
777		821
778		822
779	John Hewitt and Christopher D. Manning. 2019. <a href="#">A structural probe for finding syntax in word representations</a> . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 4129–4138, Minneapolis, Minnesota. Association for Computational Linguistics.	823
780		824
781		825
782		826
783		827
784		828
785		829
786		
	Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. <i>arXiv preprint arXiv:1503.02531</i> .	830
		831
		832
	Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for nlp. In <i>International Conference on Machine Learning</i> , pages 2790–2799. PMLR.	833
		834
		835
		836
		837
		838
		839
		840
		841
	Jeff Howe. 2008. <i>Crowdsourcing: How the power of the crowd is driving the future of business</i> . Random House.	842
		843
	Jeff Howe et al. 2006. The rise of crowdsourcing. <i>Wired magazine</i> , 14(6):1–4.	
	Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. 2020. <a href="#">XTREME: A massively multilingual multi-task benchmark for evaluating cross-lingual generalisation</a> . In <i>Proceedings of the 37th International Conference on Machine Learning</i> , volume 119 of <i>Proceedings of Machine Learning Research</i> , pages 4411–4421. PMLR.	
	Emily Jamison and Iryna Gurevych. 2015. <a href="#">Noise or additional information? leveraging crowdsourced annotation item agreement for natural language tasks</a> . In <i>Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing</i> , pages 291–297, Lisbon, Portugal. Association for Computational Linguistics.	
	Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. 2020. <a href="#">TinyBERT: Distilling BERT for natural language understanding</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 4163–4174, Online. Association for Computational Linguistics.	
	Nora Kassner, Philipp Dufter, and Hinrich Schütze. 2021. <a href="#">Multilingual LAMA: Investigating knowledge in multilingual pretrained language models</a> . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 3250–3258, Online. Association for Computational Linguistics.	
	Daniel Khashabi, Arman Cohan, Siamak Shakeri, Pedram Hosseini, Pouya Pezeshkpour, Malihe Alikhani, Moin Aminnaseri, Marzieh Bitaab, Faeze Brahman, Sarik Ghazarian, Mozhdeh Gheini, Arman Kabiri, Rabeeh Karimi Mahabadi, Omid Memarrast, Ahmadreza Mosallanezhad, Erfan Noury, Shahab Raji, Mohammad Sadegh Rasooli, Sepideh Sadeghi, Erfan Sadeqi Azer, Niloofar Safi Samghabadi, Mahsa Shafaei, Saber Sheybani, Ali Tazarv, and Yadollah Yaghoobzadeh. 2020. ParsiNLU: A suite of language understanding challenges for persian. <i>ArXiv</i> , abs/2012.06154.	
	Diederik P. Kingma and Jimmy Ba. 2015. <a href="#">Adam: A method for stochastic optimization</a> . In <i>ICLR</i> .	

844	Zhenzhong Lan, Mingda Chen, Sebastian Goodman,	<a href="#">representations</a> . In <i>Proceedings of the 2019 Confer-</i>	900
845	Kevin Gimpel, Piyush Sharma, and Radu Soricut.	<i>ence of the North American Chapter of the Associ-</i>	901
846	2020. <a href="#">Albert: A lite bert for self-supervised learning</a>	<i>ation for Computational Linguistics: Human Lan-</i>	902
847	<a href="#">of language representations</a> . In <i>International Confer-</i>	<i>guage Technologies, Volume 1 (Long and Short Pa-</i>	903
848	<i>ence on Learning Representations</i> .	<i>pers)</i> , pages 1073–1094, Minneapolis, Minnesota.	904
		Association for Computational Linguistics.	905
849	Florian Laws, Christian Scheible, and Hinrich Schütze.	Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang,	906
850	2011. <a href="#">Active learning with Amazon Mechanical</a>	Hiroaki Hayashi, and Graham Neubig. 2021a. Pre-	907
851	<a href="#">Turk</a> . In <i>Proceedings of the 2011 Conference on</i>	train, prompt, and predict: A systematic survey of	908
852	<i>Empirical Methods in Natural Language Processing</i> ,	prompting methods in natural language processing.	909
853	pages 1546–1556, Edinburgh, Scotland, UK. Associ-	<i>arXiv preprint arXiv:2107.13586</i> .	910
854	ation for Computational Linguistics.		
855	Dong-Hyun Lee et al. 2013. Pseudo-label: The simple	Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding,	911
856	and efficient semi-supervised learning method for	Yujie Qian, Zhilin Yang, and Jie Tang. 2021b. GPT	912
857	deep neural networks. In <i>Workshop on challenges in</i>	understands, too. <i>arXiv preprint arXiv:2103.10385</i> .	913
858	<i>representation learning, ICML</i> , volume 3, page 896.		
859	Ji-Ung Lee, Jan-Christoph Klie, and Iryna Gurevych.	Ilya Loshchilov and Frank Hutter. 2019. <a href="#">Decoupled</a>	914
860	2021. Annotation curricula to implicitly train non-	<a href="#">weight decay regularization</a> . In <i>International Confer-</i>	915
861	expert annotators. <i>arXiv preprint arXiv:2106.02382</i> .	<i>ence on Learning Representations</i> .	916
862	Ji-Ung Lee, Christian M. Meyer, and Iryna Gurevych.	Mitchell P. Marcus, Beatrice Santorini, and Mary Ann	917
863	2020. <a href="#">Empowering Active Learning to Jointly Optimize</a>	Marcinkiewicz. 1993. <a href="#">Building a large annotated cor-</a>	918
864	<a href="#">System and User Demands</a> . In <i>Proceedings</i>	<a href="#">pus of English: The Penn Treebank</a> . <i>Computational</i>	919
865	<i>of the 58th Annual Meeting of the Association for</i>	<i>Linguistics</i> , 19(2):313–330.	920
866	<i>Computational Linguistics</i> , pages 4233–4247, On-	Katerina Margatina, Giorgos Vernikos, Loïc Barrault,	921
867	line. Association for Computational Linguistics.	and Nikolaos Aletras. 2021. <a href="#">Active learning by ac-</a>	922
868	Brian Lester, Rami Al-Rfou, and Noah Constant. 2021.	<a href="#">quiring contrastive examples</a> . In <i>Proceedings of the</i>	923
869	<a href="#">The power of scale for parameter-efficient prompt</a>	<i>2021 Conference on Empirical Methods in Natural</i>	924
870	<a href="#">tuning</a> . In <i>Proceedings of the 2021 Conference on</i>	<i>Language Processing</i> , pages 650–663, Online and	925
871	<i>Empirical Methods in Natural Language Processing</i> ,	Punta Cana, Dominican Republic. Association for	926
872	pages 3045–3059, Online and Punta Cana, Domini-	Computational Linguistics.	927
873	can Republic. Association for Computational Lin-	Fei Mi, Yitong Li, Yasheng Wang, Xin Jiang, and Qun	928
874	guistics.	Liu. 2021a. Cins: Comprehensive instruction for few-	929
875	David D Lewis and William A Gale. 1994. A sequential	shot learning in task-oriented dialog systems. <i>arXiv</i>	930
876	algorithm for training text classifiers. In <i>SIGIR’94</i> ,	<i>preprint arXiv:2109.04645</i> .	931
877	pages 3–12. Springer.	Fei Mi, Wanhao Zhou, Lingjing Kong, Fengyu Cai,	932
878	Xiang Lisa Li and Percy Liang. 2021. <a href="#">Prefix-tuning:</a>	Minlie Huang, and Boi Faltings. 2021b. Self-training	933
879	<a href="#">Optimizing continuous prompts for generation</a> . In	improves pre-training for few-shot learning in task-	934
880	<i>Proceedings of the 59th Annual Meeting of the Asso-</i>	oriented dialog systems. In <i>Proceedings of the 2021</i>	935
881	<i>ciation for Computational Linguistics and the 11th</i>	<i>Conference on Empirical Methods in Natural Lan-</i>	936
882	<i>International Joint Conference on Natural Language</i>	<i>guage Processing</i> , pages 1887–1898.	937
883	<i>Processing (Volume 1: Long Papers)</i> , pages 4582–	George A. Miller, Claudia Leacock, Randee Tengi, and	938
884	4597, Online. Association for Computational Lin-	Ross T. Bunker. 1993. <a href="#">A semantic concordance</a> .	939
885	guistics.	In <i>Human Language Technology: Proceedings of</i>	940
886	Weixin Liang, James Zou, and Zhou Yu. 2020. <a href="#">ALICE:</a>	<i>a Workshop Held at Plainsboro, New Jersey, March</i>	941
887	<a href="#">Active learning with contrastive natural language ex-</a>	<i>21-24, 1993</i> .	942
888	<a href="#">planations</a> . In <i>Proceedings of the 2020 Conference</i>	Robert Munro Monarch. 2021. <i>Human-in-the-Loop</i>	943
889	<i>on Empirical Methods in Natural Language Process-</i>	<i>Machine Learning: Active learning and annotation</i>	944
890	<i>ing (EMNLP)</i> , pages 4380–4391, Online. Association	<i>for human-centered AI</i> . Simon and Schuster.	945
891	for Computational Linguistics.	Nikita Nangia, Saku Sugawara, Harsh Trivedi, Alex	946
892	Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu	Warstadt, Clara Vania, and Samuel R. Bowman. 2021.	947
893	Wang, Shuohui Chen, Daniel Simig, Myle Ott, Na-	<a href="#">What ingredients make for an effective crowdsourc-</a>	948
894	man Goyal, Shruti Bhosale, Jingfei Du, et al. 2021.	<a href="#">ing protocol for difficult NLU data collection tasks?</a>	949
895	Few-shot learning with multilingual language models.	In <i>Proceedings of the 59th Annual Meeting of the</i>	950
896	<i>arXiv preprint arXiv:2112.10668</i> .	<i>Association for Computational Linguistics and the</i>	951
897	Nelson F. Liu, Matt Gardner, Yonatan Belinkov,	<i>11th International Joint Conference on Natural Lan-</i>	952
898	Matthew E. Peters, and Noah A. Smith. 2019. <a href="#">Lin-</a>	<i>guage Processing (Volume 1: Long Papers)</i> , pages	953
899	<a href="#">guistic knowledge and transferability of contextual</a>	1221–1235, Online. Association for Computational	954
		Linguistics.	955

956	Sungjoon Park, Jihyung Moon, Sung-Dong Kim,	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. <a href="#">Exploring the limits of transfer learning with a unified text-to-text transformer</a> . <i>Journal of Machine Learning Research</i> , 21(140):1–67.	1014
957	Won Ik Cho, Jiyeon Han, Jangwon Park, Chisung Song, Junseong Kim, Yongsook Song, Tae Hwan Oh, Jooheon Lee, Juhyun Oh, Sungwon Lyu, Youngkuk Jeong, Inkwon Lee, Sanggyu Seo, Dongjun Lee, Hyunwoo Kim, Myeonghwa Lee, Seongbo Jang, Seungwon Do, Sunkyoung Kim, Kyungtae Lim, Jongwon Lee, Kyumin Park, Jamin Shin, Seonghyun Kim, Lucy Park, Alice H. Oh, Jung-Woo Ha, and Kyunghyun Cho. 2021. KLUE: Korean language understanding evaluation. <i>ArXiv</i> , abs/2105.09680.		1015
958			1016
959			1017
960			1018
961			1019
962			
963			1020
964			1021
965			1022
966			1023
967	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. <a href="#">PyTorch: An imperative style, high-performance deep learning library</a> . In <i>Advances in Neural Information Processing Systems</i> , volume 32. Curran Associates, Inc.		1024
968			1025
969			1026
970			1027
971			
972			1028
973			1029
974			1030
975			1031
976			1032
977	David Patterson, Joseph Gonzalez, Quoc Le, Chen Liang, Lluís-Miquel Munguia, Daniel Rothchild, David So, Maud Texier, and Jeff Dean. 2021. Carbon emissions and large neural network training. <i>arXiv preprint arXiv:2104.10350</i> .		1033
978			1034
979			1035
980			
981			1036
982	Silviu Paun, Bob Carpenter, Jon Chamberlain, Dirk Hovy, Udo Kruschwitz, and Massimo Poesio. 2018. <a href="#">Comparing Bayesian models of annotation</a> . <i>Transactions of the Association for Computational Linguistics</i> , 6:571–585.		1037
983			1038
984			
985			1039
986			1040
987	Silviu Paun and Dirk Hovy, editors. 2019. <i>Proceedings of the First Workshop on Aggregating and Analysing Crowdsourced Annotations for NLP</i> . Association for Computational Linguistics, Hong Kong, China.		1041
988			1042
989			1043
990			1044
991			1045
992	Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. <a href="#">Deep contextualized word representations</a> . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.		1046
993			1047
994			1048
995			1049
996			1050
997			1051
998			1052
999			
1000	Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. <a href="#">Language models as knowledge bases?</a> In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2463–2473, Hong Kong, China. Association for Computational Linguistics.		1053
1001			1054
1002			1055
1003			1056
1004			
1005			1057
1006			1058
1007			1059
1008			
1009	Nina Poerner, Ulli Waltinger, and Hinrich Schütze. 2020. <a href="#">E-BERT: Efficient-yet-effective entity embeddings for BERT</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 803–818, Online. Association for Computational Linguistics.		1060
1010			1061
1011			1062
1012			1063
1013			1064
			1065
			1066

1067	Yanyao Shen, Hyokun Yun, Zachary Lipton, Yakov Kro-	word representations. In <i>International Conference</i>	1125
1068	nrod, and Animashree Anandkumar. 2017. <a href="#">Deep</a>	<i>on Learning Representations</i> .	1126
1069	<a href="#">active learning for named entity recognition</a> . In		
1070	<i>Proceedings of the 2nd Workshop on Representa-</i>	Dietrich Trautmann, Johannes Daxenberger, Christian	1127
1071	<i>tion Learning for NLP</i> , pages 252–256, Vancouver,	Stab, Hinrich Schütze, and Iryna Gurevych. 2020.	1128
1072	Canada. Association for Computational Linguistics.	<a href="#">Fine-grained argument unit recognition and classi-</a>	1129
		<a href="#">fication</a> . <i>Proceedings of the AAAI Conference on</i>	1130
1073	Aditya Siddhant and Zachary C. Lipton. 2018. <a href="#">Deep</a>	<i>Artificial Intelligence</i> , 34(05):9048–9056.	1131
1074	<a href="#">Bayesian active learning for natural language pro-</a>		
1075	<a href="#">cessing: Results of a large-scale empirical study</a> .	Maria Tsimpoukelli, Jacob Menick, Serkan Cabi, S. Es-	1132
1076	In <i>Proceedings of the 2018 Conference on Empiri-</i>	lami, Oriol Vinyals, and Felix Hill. 2021. Multi-	1133
1077	<i>cal Methods in Natural Language Processing</i> , pages	modal few-shot learning with frozen language mod-	1134
1078	2904–2909, Brussels, Belgium. Association for Com-	els. <i>ArXiv</i> , abs/2106.13884.	1135
1079	putational Linguistics.		
1080	Edwin Simpson and Iryna Gurevych. 2018. <a href="#">Finding</a>	Alex Wang, Yada Pruksachatkun, Nikita Nangia, Aman-	1136
1081	<a href="#">convincing arguments using scalable Bayesian pref-</a>	preet Singh, Julian Michael, Felix Hill, Omer Levy,	1137
1082	<a href="#">erence learning</a> . <i>Transactions of the Association for</i>	and Samuel Bowman. 2019. <a href="#">SuperGLUE: A stickier</a>	1138
1083	<i>Computational Linguistics</i> , 6:357–371.	<a href="#">benchmark for general-purpose language understand-</a>	1139
		<a href="#">ing systems</a> . In <i>Advances in Neural Information</i>	1140
1084	Rion Snow, Brendan O’Connor, Daniel Jurafsky, and	<i>Processing Systems</i> , volume 32. Curran Associates,	1141
1085	Andrew Ng. 2008. <a href="#">Cheap and fast – but is it good?</a>	Inc.	1142
1086	<a href="#">evaluating non-expert annotations for natural lan-</a>		
1087	<a href="#">guage tasks</a> . In <i>Proceedings of the 2008 Conference</i>	Alex Wang, Amanpreet Singh, Julian Michael, Felix	1143
1088	<i>on Empirical Methods in Natural Language Process-</i>	Hill, Omer Levy, and Samuel Bowman. 2018. <a href="#">GLUE:</a>	1144
1089	<i>ing</i> , pages 254–263, Honolulu, Hawaii. Association	<a href="#">A multi-task benchmark and analysis platform for nat-</a>	1145
1090	for Computational Linguistics.	<a href="#">ural language understanding</a> . In <i>Proceedings of the</i>	1146
		<i>2018 EMNLP Workshop BlackboxNLP: Analyzing</i>	1147
1091	Richard Socher, Alex Perelygin, Jean Wu, Jason	<i>and Interpreting Neural Networks for NLP</i> , pages	1148
1092	Chuang, Christopher D. Manning, Andrew Ng, and	353–355, Brussels, Belgium. Association for Com-	1149
1093	Christopher Potts. 2013. <a href="#">Recursive deep models for</a>	putational Linguistics.	1150
1094	<a href="#">semantic compositionality over a sentiment treebank</a> .		
1095	In <i>Proceedings of the 2013 Conference on Empiri-</i>	Rui Wang, Masao Utiyama, Lema Liu, Kehai Chen,	1151
1096	<i>cal Methods in Natural Language Processing</i> , pages	and Eiichiro Sumita. 2017. <a href="#">Instance weighting for</a>	1152
1097	1631–1642, Seattle, Washington, USA. Association	<a href="#">neural machine translation domain adaptation</a> . In	1153
1098	for Computational Linguistics.	<i>Proceedings of the 2017 Conference on Empirical</i>	1154
		<i>Methods in Natural Language Processing</i> , pages	1155
1099	Emma Strubell, Ananya Ganesh, and Andrew McCal-	1482–1488, Copenhagen, Denmark. Association for	1156
1100	lum. 2019. <a href="#">Energy and policy considerations for</a>	Computational Linguistics.	1157
1101	<a href="#">deep learning in NLP</a> . In <i>Proceedings of the 57th</i>		
1102	<i>Annual Meeting of the Association for Computational</i>	Shuohang Wang, Yang Liu, Yichong Xu, Chenguang	1158
1103	<i>Linguistics</i> , pages 3645–3650, Florence, Italy. Asso-	Zhu, and Michael Zeng. 2021. <a href="#">Want to reduce la-</a>	1159
1104	ciation for Computational Linguistics.	<a href="#">beling cost? GPT-3 can help</a> . In <i>Findings of the</i>	1160
		<i>Association for Computational Linguistics: EMNLP</i>	1161
1105	Swabha Swayamdipta, Roy Schwartz, Nicholas Lourie,	2021, pages 4195–4205, Punta Cana, Dominican Re-	1162
1106	Yizhong Wang, Hannaneh Hajishirzi, Noah A. Smith,	public. Association for Computational Linguistics.	1163
1107	and Yejin Choi. 2020. <a href="#">Dataset cartography: Mapping</a>		
1108	<a href="#">and diagnosing datasets with training dynamics</a> . In	Alex Warstadt, Amanpreet Singh, and Samuel R. Bow-	1164
1109	<i>Proceedings of the 2020 Conference on Empirical</i>	man. 2019. <a href="#">Neural network acceptability judgments</a> .	1165
1110	<i>Methods in Natural Language Processing (EMNLP)</i> ,	<i>Transactions of the Association for Computational</i>	1166
1111	pages 9275–9293, Online. Association for Computa-	<i>Linguistics</i> , 7:625–641.	1167
1112	tional Linguistics.		
1113	Derek Tam, Rakesh R. Menon, Mohit Bansal, Shashank	Genta Indra Winata, Andrea Madotto, Zhaojiang Lin,	1168
1114	Srivastava, and Colin Raffel. 2021. <a href="#">Improving and</a>	Rosanne Liu, Jason Yosinski, and Pascale Fung. 2021.	1169
1115	<a href="#">simplifying pattern exploiting training</a> . In <i>Proceed-</i>	Language models are few-shot multilingual learners.	1170
1116	<i>ings of the 2021 Conference on Empirical Methods</i>	<i>arXiv preprint arXiv:2109.07684</i> .	1171
1117	<i>in Natural Language Processing</i> , pages 4980–4991,		
1118	Online and Punta Cana, Dominican Republic. Asso-	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	1172
1119	ciation for Computational Linguistics.	Chaumond, Clement Delangue, Anthony Moi, Pier-	1173
		ric Cistac, Tim Rault, Remi Louf, Morgan Funtow-	1174
1120	Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang,	icz, Joe Davison, Sam Shleifer, Patrick von Platen,	1175
1121	Adam Poliak, R Thomas McCoy, Najoung Kim, Ben-	Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,	1176
1122	jamin Van Durme, Sam Bowman, Dipanjan Das, and	Teven Le Scao, Sylvain Gugger, Mariama Drame,	1177
1123	Ellie Pavlick. 2019. <a href="#">What do you learn from con-</a>	Quentin Lhoest, and Alexander Rush. 2020. <a href="#">Trans-</a>	1178
1124	<a href="#">text? probing for sentence structure in contextualized</a>	<a href="#">formers: State-of-the-art natural language processing</a> .	1179
		In <i>Proceedings of the 2020 Conference on Empirical</i>	1180

1181			
1182		<i>Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	
1183			
1184	Liang Xu, Hai Hu, Xuanwei Zhang, Lu Li, Chenjie Cao,		
1185	Yudong Li, Yechen Xu, Kai Sun, Dian Yu, Cong		
1186	Yu, Yin Tian, Qianqian Dong, Weitang Liu, Bo Shi,		
1187	Yiming Cui, Junyi Li, Jun Zeng, Rongzhao Wang,		
1188	Weijian Xie, Yanting Li, Yina Patterson, Zuoyu Tian,		
1189	Yiwen Zhang, He Zhou, Shaowei Hua Liu, Zhe Zhao,		
1190	Qipeng Zhao, Cong Yue, Xinrui Zhang, Zhengliang		
1191	Yang, Kyle Richardson, and Zhenzhong Lan. 2020.		
1192	<a href="#">CLUE: A Chinese language understanding evaluation benchmark</a> .		
1193	In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> ,		
1194	pages 4762–4772, Barcelona, Spain (Online). Inter-		
1195	national Committee on Computational Linguistics.		
1196			
1197	Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Car-		
1198	bonell, Russ R Salakhutdinov, and Quoc V Le. 2019.		
1199	<a href="#">XLNet: Generalized autoregressive pretraining for language understanding</a> .		
1200	In <i>Advances in Neural Information Processing Systems</i> , volume 32. Curran		
1201	Associates, Inc.		
1202			
1203	David Yarowsky. 1995. <a href="#">Unsupervised word sense dis-</a>		
1204	<a href="#">ambiguation rivaling supervised methods</a> . In <i>33rd</i>		
1205	<i>Annual Meeting of the Association for Computa-</i>		
1206	<i>tional Linguistics</i> , pages 189–196, Cambridge, Mas-		
1207	sachusetts, USA. Association for Computational Lin-		
1208	guistics.		
1209	Wenpeng Yin, Nazneen Fatema Rajani, Dragomir		
1210	Radev, Richard Socher, and Caiming Xiong. 2020.		
1211	<a href="#">Universal natural language processing with limited</a>		
1212	<a href="#">annotations: Try few-shot textual entailment as a</a>		
1213	<a href="#">start</a> . In <i>Proceedings of the 2020 Conference on</i>		
1214	<i>Empirical Methods in Natural Language Processing</i>		
1215	<i>(EMNLP)</i> , pages 8229–8239, Online. Association for		
1216	Computational Linguistics.		
1217	Kang Min Yoo, Dongju Park, Jaewook Kang, Sang-Woo		
1218	Lee, and Woomyoung Park. 2021. <a href="#">GPT3Mix: Lever-</a>		
1219	<a href="#">aging large-scale language models for text augmen-</a>		
1220	<a href="#">tation</a> . In <i>Findings of the Association for Computa-</i>		
1221	<i>tional Linguistics: EMNLP 2021</i> , pages 2225–2239,		
1222	Punta Cana, Dominican Republic. Association for		
1223	Computational Linguistics.		
1224	Man-Ching Yuen, Irwin King, and Kwong-Sak Leung.		
1225	2011. <a href="#">A survey of crowdsourcing systems</a> . In <i>2011</i>		
1226	<i>IEEE Third International Conference on Privacy, Se-</i>		
1227	<i>curity, Risk and Trust and 2011 IEEE Third Interna-</i>		
1228	<i>tional Conference on Social Computing</i> , pages 766–		
1229	773.		
1230	Mike Zhang and Barbara Plank. 2021. <a href="#">Cartography ac-</a>		
1231	<a href="#">tive learning</a> . In <i>Findings of the Association for Com-</i>		
1232	<i>putational Linguistics: EMNLP 2021</i> , pages 395–		
1233	406, Punta Cana, Dominican Republic. Association		
1234	for Computational Linguistics.		
1235	Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015.		
1236	<a href="#">Character-level convolutional networks for text clas-</a>		
1237	<a href="#">sification</a> . In <i>Advances in Neural Information Pro-</i>		
1238	<i>cessing Systems</i> , volume 28. Curran Associates, Inc.		
	Ye Zhang, Matthew Lease, and Byron C. Wallace. 2017.		1239
	Active discriminative text representation learning. In		1240
	<i>Proceedings of the Thirty-First AAAI Conference on</i>		1241
	<i>Artificial Intelligence</i> , AAAI’17, page 3386–3392.		1242
	AAAI Press.		1243
	Mengjie Zhao, Philipp Dufter, Yadollah Yaghoobzadeh,		1244
	and Hinrich Schütze. 2020a. <a href="#">Quantifying the con-</a>		1245
	<a href="#">textualization of word representations with seman-</a>		1246
	<a href="#">tic class probing</a> . In <i>Findings of the Association</i>		1247
	<i>for Computational Linguistics: EMNLP 2020</i> , pages		1248
	1219–1234, Online. Association for Computational		1249
	Linguistics.		1250
	Mengjie Zhao, Tao Lin, Fei Mi, Martin Jaggi, and Hin-		1251
	rich Schütze. 2020b. <a href="#">Masking as an efficient alter-</a>		1252
	<a href="#">native to finetuning for pretrained language models</a> .		1253
	In <i>Proceedings of the 2020 Conference on Empirical</i>		1254
	<i>Methods in Natural Language Processing (EMNLP)</i> ,		1255
	pages 2226–2241, Online. Association for Computa-		1256
	tional Linguistics.		1257
	Mengjie Zhao and Hinrich Schütze. 2021. <a href="#">Discrete and</a>		1258
	<a href="#">soft prompting for multilingual models</a> . In <i>Proceed-</i>		1259
	<i>ings of the 2021 Conference on Empirical Methods</i>		1260
	<i>in Natural Language Processing</i> , pages 8547–8555,		1261
	Online and Punta Cana, Dominican Republic. Asso-		1262
	ciation for Computational Linguistics.		1263
	Mengjie Zhao, Yi Zhu, Ehsan Shareghi, Ivan Vulić, Roi		1264
	Reichart, Anna Korhonen, and Hinrich Schütze. 2021.		1265
	<a href="#">A closer look at few-shot crosslingual transfer: The</a>		1266
	<a href="#">choice of shots matters</a> . In <i>Proceedings of the 59th</i>		1267
	<i>Annual Meeting of the Association for Computational</i>		1268
	<i>Linguistics and the 11th International Joint Confer-</i>		1269
	<i>ence on Natural Language Processing (Volume 1: Long</i>		1270
	<i>Papers)</i> , pages 5751–5767, Online. Association		1271
	for Computational Linguistics.		1272

## A Reproducibility Checklist

### A.1 Computing infrastructure

We use four Tesla V100 GPUs to prompt each of the LMTurkers, and a single Tesla V100 GPU is used when finetuning the small model  $\mathcal{S}$ .

### A.2 Datasets

For SST2, CoLA, and RTE, we use the official datasets available on the benchmark website [gluebenchmark.com](http://gluebenchmark.com). We download SST5 dataset from [nlp.stanford.edu/sentiment](http://nlp.stanford.edu/sentiment) and AGNews from the link provided by Zhang et al. (2015).

The number of testing examples of each dataset is shown in Table 2. Note that for SST2, CoLA, and RTE,  $\mathcal{G}^{dev}$  is sampled from the training set, and the dev set is used as the test set.

CoLA	SST5	RTE	AGNews	SST2
1042	2210	277	7600	872

Table 2: Number of testing examples.

## B Numerical Results

Table 3 reports the numerical value of Figure 2.

## C Prompting Details

For each task, we list the five prompts employed to adapt a PLM to a LMTurker. “[v]” is a prompting token whose trainable embedding vector is randomly initialized.

For **SST5**, we use following prompts:

- “[v] x It is [MASK].”
- “[v] x Such a [MASK] movie.”
- “x [v] It is pretty [MASK].”
- “It is [MASK] because x [v]”
- “x So it is [MASK]. [v]”

and the PLM picks a word from {“crap”, “bad”, “normal”, “good”, “perfect”}. to fill the position of “[MASK]”. The mapping {“crap”  $\rightarrow$  1, “bad”  $\rightarrow$  2, “normal”  $\rightarrow$  3, “good”  $\rightarrow$  4, “perfect”  $\rightarrow$  5 } is used to convert model predictions to numerical values.

For **SST2**, we use following prompts:

- “[v] x It is [MASK].”
- “[v] x Such a [MASK] movie.”

• “x [v] It is pretty [MASK].” 1310

• “It is [MASK] because x [v]” 1311

• “x So it is [MASK]. [v]” 1312

and the PLM picks a word from {“bad”, “good”} to fill the position of “[MASK]”. The mapping {“bad”  $\rightarrow$  0, “good”  $\rightarrow$  1} is used. 1313

For **AGNews**, we use following prompts: 1314

• “[v] x It is about [MASK].” 1315

• “x [v] Topic: [MASK].” 1316

• “x [v] The text is about [MASK].” 1317

• “x Topic: [MASK]. [v]” 1318

• “x [v] [MASK].” 1319

and the PLM picks a word from {“world”, “sports”, “economy”, “technology”} to fill the position of “[MASK]”. The mapping {“world”  $\rightarrow$  1, “sports”  $\rightarrow$  2, “economy”  $\rightarrow$  3, “technology”  $\rightarrow$  4 } is used. 1320

For **CoLA**, we use following prompts: 1321

• “[v] x It sounds [MASK].” 1322

• “[v] x The sentence is [MASK].” 1323

• “[v] x It is a [MASK] sentence.” 1324

• “x [v] [MASK].” 1325

• “[v] x [MASK].” 1326

and the PLM picks a word from {“wrong”, “ok”} to fill the position of “[MASK]”. The mapping {“wrong”  $\rightarrow$  0, “okay”  $\rightarrow$  1} is used. 1327

For **RTE**, we use following prompts: 1328

• “p Question: h? [v] Answer: [MASK].” 1329

• “p [SEP] h? [MASK]. [v]” 1330

• “p [SEP] h? [v] answer: [MASK].” 1331

• “p [SEP] In short h. [MASK]. [v]” 1332

• “[v] p [SEP] In short h. [MASK].” 1333

where **p** and **h** refer to premise and hypothesis. The PLM picks a word from {“No”, “Yes”} to fill the position of “[MASK]”. The mapping {“No”  $\rightarrow$  0, “Yes”  $\rightarrow$  1} is used. 1334

	$\mathcal{G}^8$		$\mathcal{G}^{16}$		$\mathcal{G}^{32}$	
	Workers	$\mathcal{S}$	Workers	$\mathcal{S}$	Workers	$\mathcal{S}$
SST2	91.13±0.52	67.63±8.01	91.93±1.09	75.83±1.35	91.97±0.83	76.37±3.16
	91.63±0.68		93.08±0.62		91.70±1.78	
	90.18±1.00		91.74±1.04		91.21±1.83	
	90.83±0.58		90.79±0.47		91.13±0.24	
	90.52±1.84		91.67±1.36		93.23±0.37	
SST5	41.37±1.55	28.47±1.61	45.16±2.13	34.97±1.51	45.91±0.96	33.47±2.79
	42.32±2.04		45.96±2.12		48.64±0.59	
	40.57±2.70		46.70±0.93		50.53±0.94	
	37.69±1.34		42.53±2.43		43.32±3.42	
	38.05±2.60		42.96±0.69		45.72±1.43	
RTE	68.95±1.47	57.30±1.79	68.35±2.29	61.50±0.78	71.72±1.96	62.93±0.74
	54.99±3.76		57.64±3.23		58.48±3.59	
	62.70±1.33		70.88±1.70		68.47±1.19	
	50.42±2.07		58.60±1.62		59.33±4.72	
	51.99±4.45		57.88±2.83		60.41±2.47	
AGNews	75.39±5.25	66.37±2.95	83.06±0.83	69.40±0.93	84.92±0.28	76.53±0.41
	85.40±1.43		87.71±0.07		87.79±1.08	
	78.83±4.77		83.59±2.96		87.39±1.29	
	85.07±1.09		87.69±0.04		87.17±0.67	
	79.95±0.86		80.15±3.38		83.32±0.59	
CoLA	0.14±1.43	0.97±4.40	11.81±7.82	4.27±3.26	19.88±3.30	2.50±2.41
	2.42±4.84		15.23±7.07		22.51±0.96	
	7.40±8.12		19.71±1.89		26.34±1.54	
	9.91±7.98		17.14±2.48		18.15±0.63	
	15.33±2.15		19.66±0.48		27.58±7.09	

Table 3: Few-shot performance of the five LMTurkers and the small model  $\mathcal{S}$ . Each experiment is repeated three times and we report mean and standard deviation.

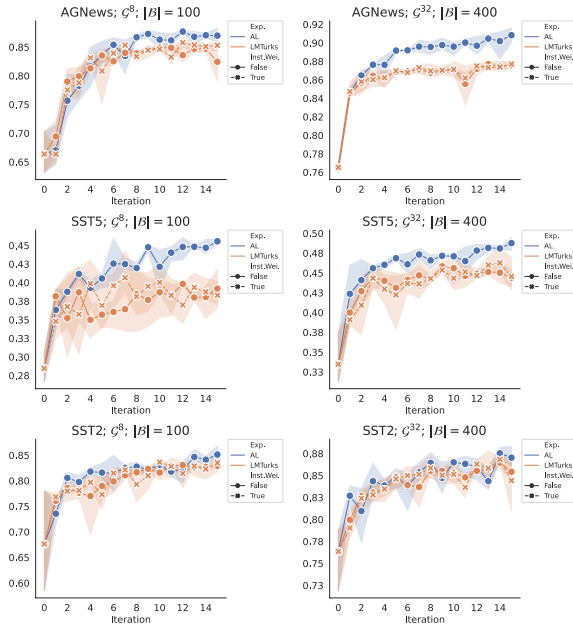


Figure 8: Weighting the training instances from LMTurkers.

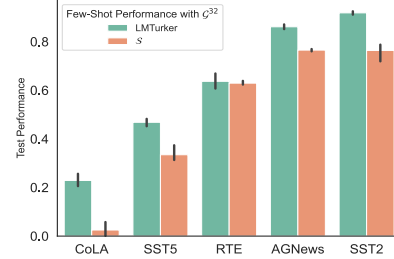


Figure 9: Few-shot performance on test set of LMTurkers and  $\mathcal{S}$  when using the few-shot gold datasets  $\mathcal{G}^{32}$ .

## E Instance Weighting

Following Wang et al. (2017), we associate each example  $(\mathbf{x}, \hat{\mathbf{y}}, \mathbf{l}) \in \mathcal{D}^j$  with weight  $1 - \text{entropy}(\mathbf{l})$  when computing the loss during the course of training  $\mathcal{S}^j$ . We can interpret this weight as a measure of the certainty of the LMTurkers ensemble.

Figure 8 reports the performance of  $\mathcal{S}$  when using instance weighting, however, the impacts are less noticeable.

## D More Visualizations

Figure 9 compares the few-shot performance of LMTurkers and  $\mathcal{S}$  when using  $\mathcal{G}^{32}$ .

Figure 10 visualizes the performance of  $\mathcal{S}$  when different  $|\mathcal{G}|$  and  $|\mathcal{B}|$  are used.



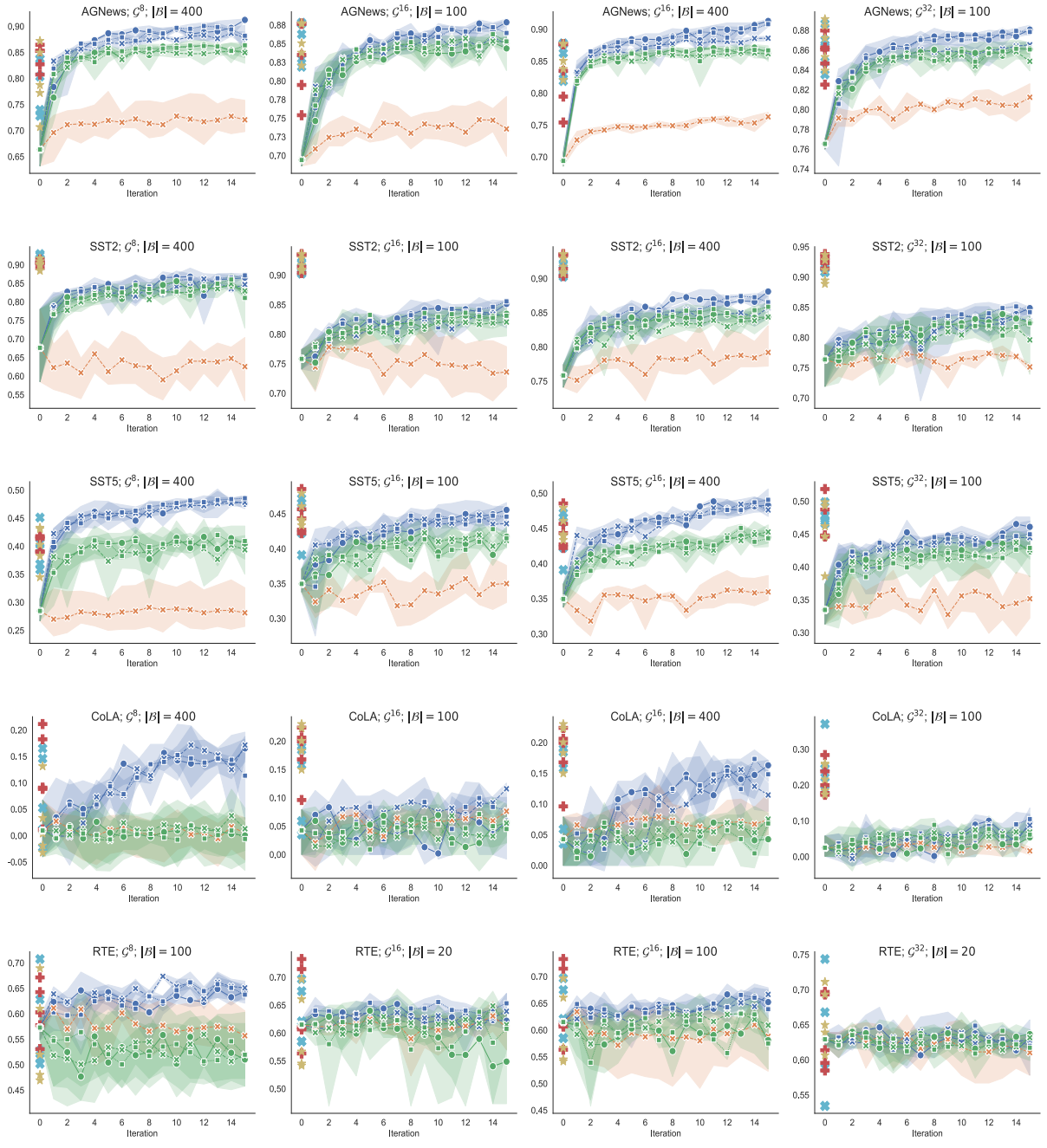


Figure 10: Improving  $S$  with active learning (blue), self training (orange), and LMTurk (green). Free markers at step zero show LMTurker performances; colors distinguish random seeds. Three acquisition functions are: Entropy ( $\bullet$ ), LeastConfident ( $\blacksquare$ ), random sampling ( $\times$ ). At iteration  $j$ , each experiment is repeated three times; we show mean and standard deviation. We evaluate different  $|\mathcal{G}|$  and  $|\mathcal{B}|$ .