

# ACTUNE: Uncertainty-Aware Active Self-Training for Active Fine-Tuning of Pretrained Language Models

Anonymous ACL submission

## Abstract

Although fine-tuning pre-trained language models (PLMs) renders strong performance in many NLP tasks, it relies on excessive labeled data. Recently, researchers have resorted to active fine-tuning for enhancing the label efficiency of PLM fine-tuning, but existing methods of this type usually ignore the potential of unlabeled data. We develop ACTUNE, a new framework that improves the label efficiency of active PLM fine-tuning by unleashing the power of unlabeled data via self training. ACTUNE switches between data annotation and model self-training based on uncertainty: the unlabeled samples of high-uncertainty are selected for annotation, while the ones from low-uncertainty regions are used for model self-training. Additionally, we design (1) a region-aware sampling strategy to avoid redundant samples when querying annotations and (2) a momentum-based memory bank to dynamically aggregate the model’s pseudo labels to suppress label noise in self-training. Experiments on 6 text classification datasets show that ACTUNE outperforms the strongest active learning and self-training baselines and improves the label efficiency of PLM fine-tuning by 56.2% on average.

## 1 Introduction

Fine-tuning pre-trained language models (PLMs) has achieved enormous success in natural language processing (NLP) (Devlin et al., 2019; Liu et al., 2019; Brown et al., 2020), one of which is the competitive performance it offers when consuming only a few labeled data (Bansal et al., 2020; Gao et al., 2021). However, there are still significant gaps between few-shot and fully-supervised PLM fine-tuning in many classification tasks. Besides, the performance of few-shot PLM fine-tuning can vary substantially with different sets of training data (Bragg et al., 2021). Therefore, there is a crucial need for PLM fine-tuning approaches with better label-efficiency and being robust to selection

of training data, especially for applications where labeled data are scarce and expensive to obtain.

Towards this goal, researchers have resorted to active fine-tuning of PLMs and achieved comparable performance to fully-supervised methods with much less annotated samples (Ein-Dor et al., 2020; Margatina et al., 2021a,b; Yuan et al., 2020). Nevertheless, they usually neglect unlabeled data, which can be useful for improving label efficiency for PLM fine-tuning (Du et al., 2021). To leverage those unlabeled data to improve label efficiency of active learning, efforts have been made in the semi-supervised active learning literature (Wang et al., 2016; Rottmann et al., 2018; Siméoni et al., 2020), but the proposed query strategies can return highly redundant samples due to limited representation power, resulting in suboptimal label efficiency. Moreover, they usually rely on pseudo-labeling to utilize unlabeled data, which requires greater (yet often absent) care to denoise the pseudo labels, otherwise the errors could accumulate and deteriorate the model performance. This phenomenon can be even more severe for PLMs, as the fine-tuning process often suffers from the instability issue caused by different weight initialization and data orders (Dodge et al., 2020). Thus, it still remains open and challenging to design robust and label efficient method for active PLM fine-tuning.

To tackle above challenges, we propose ACTUNE, a new method that improves the label efficiency and robustness of active PLM fine-tuning with self-training. Based on the estimated uncertainty of data, ACTUNE chooses from one of the following cases in each learning round: (1) when the average uncertainty of a region is low, we trust the model’s prediction and select most certain predictions within the region for self-training; (2) when the average uncertainty of a region is high, indicating inadequate observations for parameter learning, we actively annotate most uncertain samples within the region to improve the model. Dif-

ferent from existing AL methods that only leverage uncertainty for querying labels, our uncertainty-driven self-training paradigm gradually unleashes the data with low uncertainty via self-training, while reducing the chance of error propagation triggered by highly-uncertain mis-labeled data.

To further boost the performance on downstream tasks, we design two techniques, namely region-aware sampling (RS) and momentum-based memory bank (MMB) to improve the query strategies and suppress label noise for ACTUNE. Inspired by the fact that existing uncertainty-based AL methods often end up choosing uncertain yet repetitive data (Ein-Dor et al., 2020; Margatina et al., 2021b), we design a region-aware sampling technique to promote both diversity and representativeness by leveraging the representation power of PLMs. Specifically, we first estimate the uncertainties of the unlabeled data with PLMs, then cluster the data using their PLM representations and weigh the data by the corresponding uncertainty. Such a clustering scheme partitions the embedding space into small sub-regions with an emphasis on highly-uncertain samples. Finally, by sampling over multiple high-uncertainty regions, our strategy selects data with high uncertainty and low redundancy.

To rectify the erroneous pseudo labels derived by self-training, we design a simple but effective way to select low-uncertainty data for self-training. Our method is motivated by the fact that fine-tuning PLMs suffer from instability issues — distinct initializations and data orders can result in a large variance of the task performance (Dodge et al., 2020; Zhang et al., 2020; Mosbach et al., 2021). However, previous approaches only select pseudo-labeled data based on the prediction of the current round and therefore are less reliable. In contrast, we maintain a dynamic memory bank to save the predictions of unlabeled samples for later use. We propose a momentum updating method to dynamically aggregate the predictions from preceding rounds (Laine and Aila, 2016) and select low-uncertainty samples based on aggregated prediction. As a consequence, only the samples with high prediction confidence over multiple rounds will be used for self-training, which mitigates the issue of label noise. We highlight that our active self-training approach is an efficient substitution to existing AL methods, requiring ignorable extra computational cost.

Our key contributions are: (1) an active self-training paradigm ACTUNE that integrates the ben-

efit of self-training and active learning in a principled way to minimize the labeling cost for fine-tuning PLMs; (2) a region-aware querying strategy to enforce both the informativeness and the diversity of queried samples during AL; (3) a simple and effective momentum-based method to harness the predictions for preceding rounds to alleviate the label noise in self-training and (4) experiments on 6 benchmarks demonstrating ACTUNE improves the label efficiency over existing self-training and active learning baselines by 56.2%.

## 2 Uncertainty-aware Active Self-training

### 2.1 Problem Formulation

We study active fine-tuning of pre-trained language models for text classification, formulated as follows: Given a small number of labeled samples  $X_L = \{f(x_i; y_i)g_{i=1}^L\}$  and unlabeled samples  $X_U = \{f(x_j)g_{j=1}^U\}$  ( $j \in X_L \cup X_U$ ), we aim to fine-tune a pre-trained language model  $f(\mathbf{x}; \cdot) : X \rightarrow Y$  in an interactive way: we perform active self-training for  $T$  rounds with the total labeling budget  $b$ . In each round, we aim to query  $B = b/T$  samples denoted as  $B$  from  $X_U$  to fine-tune a pre-trained language model  $f(\mathbf{x}; \cdot)$  with both  $X_L; B$  and  $X_U$  to maximize the performance on downstream text classification tasks. Here  $X = X_L \cup X_U$  denotes all samples and  $Y = \{1, 2, \dots, C\}$  is the label set, where  $C$  is the number of classes.

### 2.2 Overview of ACTUNE Framework

We now present our active self-training paradigm ACTUNE underpinned by estimated uncertainty. We begin the active self-training loop by fine-tuning a BERT  $f^{(0)}$  on the initial labeled data  $X_L$ . Formally, we solve the following optimization problem

$$\min_{f} \frac{1}{|X_L|} \sum_{(\mathbf{x}_i; y_i) \in X_L} \text{CE}(f(\mathbf{x}_i; y_i); y_i) \quad (1)$$

In round  $t$  ( $1 \leq t \leq T$ ) of the active self-training procedure, we first calculate the uncertainty score based on a given function  $a_i^{(t)} = a(\mathbf{x}_i; t)$ <sup>1</sup> for all  $\mathbf{x}_i \in X_U$ . Then, we query labeled samples and generate pseudo-labels for unlabeled data  $X_U$  simultaneously to facilitate self-training. For each sample  $\mathbf{x}_i$ , the pseudo-label  $\hat{y}$  is calculated based on the current model’s output:

$$\hat{y} = \underset{j \in Y}{\operatorname{argmax}} f(\mathbf{x}_i; t)_j \quad (2)$$

<sup>1</sup>Note that ACTUNE is agnostic to the way uncertainty score  $a_i^{(t)}$  is computed.

---

**Algorithm 1:** Training Procedures of ACTUNE.

**Input:** Initial labeled samples  $X_L$ ; Unlabeled samples  $X_U$ ; Pre-trained LM  $f(\cdot; \cdot)$ , number of active self-training rounds  $T$ .

// Fine-tune the LM with initial labeled data.

1. Calculate  $u_k^{(0)}$  based on Eq. (1).
2. Initialize the memory bank  $g(\mathbf{x}; \cdot)$  based on the current prediction.

// Conduct active self-training with all data.

**for**  $t = 1; 2; \dots; T$  **do**

1. Run weighted K-Means (Eq. (3), (4)) until convergence.
2. Select sample set  $O^{(t)}$  for AL and  $S^{(t)}$  for self-training from  $X_U$  based on Eq. (11) or (13).
3. Augment the labeled set  $X_L = X_L \cup O^{(t)}$ .
4. Obtain  $f(\cdot; \cdot)^{(t)}$  by finetuning  $f(\cdot; \cdot)$  with  $L_{ST}$  (Eq. (14)) using AdamW.
5. Update memory bank  $g(\mathbf{x}; \cdot)$  with Eq. (10) or (12).

**Output:** The final fine-tuned model  $f(\cdot; \cdot)^{(T)}$ .

---

where  $f(\mathbf{x}; \cdot)^{(t)} \in \mathbb{R}^C$  is a probability simplex and  $[f(\mathbf{x}; \cdot)^{(t)}]_j$  is the  $j$ -th entry. The procedure of ACTUNE is summarized in Algorithm 1.

### 2.3 Region-aware Sampling for Active Learning on High-uncertainty Data

After obtaining the uncertainty for unlabeled data, we aim to query annotation for high-uncertainty samples. However, directly sampling the most uncertain samples gives suboptimal result since uncertainty-based sampling tends to query repetitive data (Ein-Dor et al., 2020) and results in poor representativeness of the overall data distribution.

To tackle this issue, we propose region-aware sampling to capture both *uncertainty* and *diversity* during active self-training. Specifically, in the  $t$ -th round, we first conduct the weighted K-means clustering (Huang et al., 2005), which weights samples based on their uncertainty. Denote  $K$  the number of clusters and  $\mathbf{v}_i^{(t)} = \text{BERT}(\mathbf{x}_i)$  the representation of  $\mathbf{x}_i$  from the penultimate layer of BERT. The weighted K-means first initializes the center of each cluster  $\mathbf{c}_i$  ( $1 \leq i \leq K$ ) via K-Means++ (Arthur and Vassilvitskii, 2007). Then, it jointly updates the centroid of each cluster and assigns each sample to cluster  $C_i$  as

$$\mathbf{c}_i^{(t)} = \underset{P^{k=1, \dots, K}}{\operatorname{argmin}} k \mathbf{v}_i \quad k \neq i; \quad (3)$$

$$\mathbf{c}_k^{(t)} = \frac{\sum_{\mathbf{x} \in C_k^{(t)}} a(\mathbf{x}_i; \cdot)^{(t)} \mathbf{v}_i^{(t)}}{\sum_{\mathbf{x} \in C_k^{(t)}} a(\mathbf{x}_i; \cdot)^{(t)}} \quad (4)$$

where  $C_k^{(t)} = \{\mathbf{x}_i^{(t)} \mid j \in C_k^{(t)}\}$  stands for the  $k$ -th cluster. The above two steps in Eq. (3), (4) are repeated until convergence. Compared with vanilla K-Means method, the weighting

scheme increases the density of the samples with high uncertainty, thus enabling the K-Means methods to discover clusters with high uncertainty. After obtaining  $K$  regions with the corresponding data  $C_k^{(t)}$ , we calculate the uncertainty of each region as

$$u_k^{(t)} = U(C_k^{(t)}) + I(C_k^{(t)}) \quad (5)$$

where

$$U(C_k^{(t)}) = \frac{1}{j \in C_k^{(t)}} \sum_{\mathbf{x}_i \in C_k^{(t)}} a(\mathbf{x}_i; \cdot)^{(t)} \quad (6)$$

stands for the average uncertainty of samples and

$$I(C_k^{(t)}) = \frac{1}{j \in C_k^{(t)}} \sum_{j \in C_k^{(t)}} f_j^{(t)} \log f_j^{(t)} \quad (7)$$

stands for the inter-class diversity within cluster  $k$  and  $f_j^{(t)} = \frac{|\{i \mid f_{ij} = j\}|}{j \in C_k^{(t)}}$  represents the frequency

of class  $j$  on cluster  $k$ . Notably, the term  $U(C_k^{(t)})$  assigns higher score for clusters with more uncertain samples, and  $I(C_k^{(t)})$  grants higher scores for clusters containing samples with more diverse predicted classes from pseudo labels since such clusters would be closer to the decision boundary.

Then, we rank the clusters in an ascending order according to  $u_k^{(t)}$ . A high score indicates high uncertainty of the model in these regions, and we need to actively annotate the associated instances to reduce uncertainty and improve the model’s performance. We adopt a hierarchical sampling strategy: we first select the  $M$  clusters with the highest uncertainty, and then sample  $b^l = b \frac{B}{M} C$  data with the highest uncertainty to form the batch  $O^{(t)}$ .<sup>2</sup>

$$K_a^{(t)} = \text{top-}M \ u_k^{(t)}; \quad (8)$$

$$O^{(t)} = \bigcup_{k \in K_a^{(t)}} C_{a,k}^{(t)} \quad \text{where } C_{a,k}^{(t)} = \text{Top-}b^l \ a(\mathbf{x}_i; \cdot)^{(t)}; \quad \mathbf{x}_i \in C_k^{(t)}$$

We remark that such a hierarchical sampling strategy queries most uncertain samples from *different* regions, thus the uncertainty and diversity of queried samples can be both achieved.

### 2.4 Self-training for Most Confident Data from Low-uncertainty Regions

For self-training, we aim to select unlabeled samples which are *most likely* to have been correctly classified by the current model. This requires the sample to have low uncertainty. Therefore, we select the top  $k$  samples from the  $M$  lowest uncertainty regions to form the acquired batch  $S^{(t)}$ :

<sup>2</sup>If the number of samples in the  $i$ -th cluster  $C_i$  is smaller than  $b^l$ , then we sample all the data within  $C_i$  and increase the budget for the  $(i + 1)$ -th cluster by  $b^l - |C_i|$ .

$$C_s^{(t)} = \prod_{k=2K_s^{(t)}} C_k^{(t)} \text{ where } K_s^{(t)} = \text{bottom-M}_{k \in \{1, \dots, K\}} u_k^{(t)}; \quad (9)$$

$$S^{(t)} = \text{bottom-k}_{x_i \in 2C_s^{(t)}} a(x_i;^{(t)});$$

Momentum-based Memory Bank for Self-training. As PLMs are sensitive to the stochasticity involved in fine-tuning, the model suffers from the instability issue — different weight initialization and data orders may result in different predictions on the same dataset (Dodge et al., 2020).

Additionally, if the model gives inconsistent predictions in different rounds for a specific sample, then it is potentially uncertain about the sample, and adding it to the training set may harm the active self-training process. For example, for a two-class classification problem, suppose we obtain  $f(x;^{(t-1)}) = [0.65; 0.35]$  for sample  $x$  in the round  $(t-1)$  and  $f(x;^{(t)}) = [0.05; 0.95]$  for the round  $t$ . Although the model is quite 'confident' on the class  $(t+1)$  on them. Although we only include  $x$  when we only consider the result of the round  $t$ , it gives contradictory predictions over these two consecutive rounds, which indicates that the model is still uncertain to which class  $x$  belongs.

To effectively mitigate the noise and stabilize the active self-training process, we maintain a dynamic memory bank to save the results from previous rounds, and use momentum update (He et al., 2020; Laine and Aila, 2016) to aggregate the results from both the previous and current rounds. Then, during active self-training, we will select samples with the highest aggregated score. In this way, only those samples that the model is certain about over previous rounds will be selected for self-training. We

design two variants for the memory bank, namely prediction-based and value-based aggregation.

**Prediction based Momentum Update.** We adopt an exponential moving average approach to aggregate the prediction  $g(x;^{(t)})$  on round  $t$  as

$$g(x;^{(t)}) = m_t f(x;^{(t)}) + (1 - m_t) g(x;^{(t-1)}); \quad (10)$$

where  $m_t = (1 - \frac{1}{T})m_L + \frac{1}{T}m_H$  ( $0 < m_L < m_H < 1$ ) is a momentum coefficient. We gradually increase the weight for models on later rounds, since they are trained with more labeled data, thus being able to provide more reliable predictions. Then, we calculate the uncertainty based on  $g(x;^{(t)})$  and rewrite Eq. (9) and (2) as

$$S^{(t)} = \text{bottom-k}_{x_i \in 2C_s^{(t)}} a(x_i; g(x;^{(t)});^{(t)});$$

$$\hat{y} = \underset{j \in Y}{\text{argmax}} g(x;^{(t)})_j; \quad (11)$$

**Value-based Momentum Update.** For methods that do not directly use prediction for uncertainty estimation, we aggregate the uncertainty value as  $g(x;^{(t)}) = m_t a(x;^{(t)}) + (1 - m_t) g(x;^{(t-1)})$ :

Then, we use Eq.(12) to sample low-uncertainty data for self-training as

$$S^{(t)} = \text{bottom-k}_{x_i \in 2C_s^{(t)}} g(x_i;^{(t)});$$

$$\hat{y} = \underset{j \in Y}{\text{argmax}} f(x;^{(t)})_j; \quad (13)$$

By aggregating the prediction results over previous rounds, we filter the sample with inconsistent predictions to suppress noisy labels.

## 2.5 Model Learning and Update

After obtaining both the labeled data and pseudo-labeled data, we fine-tune a new pre-trained BERT model on them. Although we only include low-uncertainty samples during self-training, it is difficult to eliminate all the wrong pseudo-labels, and such mislabeled samples can still hurt model performance. To suppress such label noise, we use a threshold-based strategy to further remove noisy labels by selecting samples that agree with the corresponding pseudo labels. The loss objective of optimizing  $(t+1)$  is

$$L_{ST} = \frac{1}{|X_L|} \sum_{x_i \in X_L} \text{CE}(f(x_i;^{(t+1)}); y_i) + \frac{1}{|S^{(t)}|} \sum_{x_i \in S^{(t)}} \text{CE}(f(x_i;^{(t+1)}); \hat{y}_i); \quad (14)$$

where  $\alpha$  is a hyper-parameter balancing the weight between clean and pseudo labels, and  $\sigma$  stands for the thresholding function.

**Complexity Analysis.** The running time of ACTUNE is mainly consisted of two parts: the inference time  $O(|X_U|)$  and the time for K-Means clustering  $O(dK|X_U|)$ , where  $d$  is the dimension of the BERT feature  $\mathbf{v}$ . Note that the clustering can be efficiently implemented with FAISS (Johnson et al., 2019), and will not excessively increase the total running time. For self-training, the size of the memory bank  $|g(x;^{(t)})|$  is proportional to  $|X_U|$ , while the extra computation of maintaining this dictionary is ignorable since we do not inference over the unlabeled data for multiple times in each round as BALD (Gal et al., 2017) does. The running time of ACTUNE will be shown in section C.

Dataset	Label Type	# Class	# Train	# Dev	# Test
SST-2	Sentiment	2	60.6k	0.8k	1.8k
AG News	News Topic	4	119k	1k	7.6k
Pubmed	Medical Abstract	5	180k	1k	30.1k
DBPedia	Wikipedia Topic	14	280k	1k	70k
TREC	Question	6	5.0k	0.5k	0.5k
Chemprot	Medical Abstract	10	12.8k	0.5k	1.6k

Table 1: Dataset Statistics. For DBPedia, we randomly sample 20k sample from each class due to the limited computational resource.

### 3 Experiments

#### 3.1 Experiment Setup

**Tasks and Datasets.** In our main experiments, we study over 4 benchmark datasets, including SST-2 (Socher et al., 2013) for sentiment analysis, AGNews (Zhang et al., 2015) for news topic classification, Pubmed-RCT (Dernoncourt and Lee, 2017) for medical abstract classification, and DBPedia (Zhang et al., 2015) for wikipedia topic classification. For weakly-supervised text classification, we choose 2 datasets, namely TREC (Li and Roth, 2002) and Chemprot (Krallinger et al., 2017) from the WRENCH benchmark (Zhang et al., 2021) for evaluation. The statistics are shown in table 1.

**Active Learning Setups.** Following (Yuan et al., 2020), we set the number of rounds  $r=10$ , the overall budget for all datasets  $B=1000$  and the initial size of the labeled set  $|X_{i,j}|$  is set to 100. To simulate AL, in each round, we sample a batch of 100 samples from the unlabeled set and query labels for them. Then we move this batch to the labeled set. Since large development sets are impractical in low-resource settings (Kann et al., 2019), we keep the size of development set as 1000, which is the same as the labeling budget. For weakly-supervised text classification, since the datasets are much smaller, we keep the labeling budget and the size of development set  $b=500$ .

**Implementation Details.** We choose RoBERTa-base (Liu et al., 2019) from the HuggingFace code base (Wolf et al., 2020) as the backbone for ACTUNE and all baselines except for Pubmed and Chemprot, where we use SciBERT (Beltagy et al., 2019), a BERT model pre-trained on scientific corpora. In each round, we train from scratch to avoid badly overfitting the data collected in earlier rounds as observed by Hu et al. (2019). More details are in Appendix B.

**Hyperparameters.** The hyperparameters setting is in Appendix B.6 for ACTUNE and B.7 for base-

lines. In the  $t$ -th round of active self-training, we increase the number of pseudo-labeled samples by  $k$ , where  $k$  equals to 500 for TREC and Chemprot, 3000 for SST-2 and Pubmed-RCT, and 500 for others. For the momentum factor, we tune from  $[0.6; 0.7; 0.8]$  and  $m_H$  from  $[0.8; 0.9; 1.0]$  and report the best  $m_L; m_H$  based on the performance

Baselines.

**Self-training Methods:** (1) Self-training (ST, Lee (2013)) It is the vanilla self-training method that generates pseudo labels for unlabeled data. (2) UST (Mukherjee and Awadallah, 2020; Rizve et al., 2021): It is an uncertainty-based self-training method that only uses low-uncertainty data for self-training. (3) COSINE (Yu et al., 2021): It uses self-training to re-tune LM with weakly-labeled data, which achieves SOTA performance on various text datasets in WRENCH benchmark (Zhang et al., 2021). Note that for these two baselines, we randomly sample labeled data as the initialization. Also, UST is only used in main experiments in Sec. 3.2 and COSINE is evaluated in Sec 3.3.

**Active Learning Methods:** (1) Random: It acquires annotation randomly, which serves as a baseline for all methods. (2) Entropy (Holub et al., 2008): It is an uncertainty-based method that acquires annotations on samples with the highest predictive entropy. (3) BALD (Gal et al., 2017): It is also an uncertainty-based method, which calculates model uncertainty using MC Dropout (Gal and Ghahramani, 2015). (4) BADGE (Ash et al., 2020): It first selects high uncertainty samples then uses KMeans++ over the gradient embedding to sample data. (5) ALPS (Yuan et al., 2020): It uses the masked language model (MLM) loss of BERT to query labels for samples. (6) CAL (Margatina et al., 2021b) is the most recent AL method for pre-trained LMs. It calculates the uncertainty of each sample based on the KL divergence between the prediction of itself and its neighbors' prediction. **Semi-supervised Active Learning (SSAL) Methods:** (1) ASST (Tomanek and Hahn, 2009; Siméoni et al., 2020) is an active semi-supervised learning method that jointly queries labels for AL and samples pseudo labels for self-training. (2) CEAL (Wang et al., 2016) acquires annotations on informative samples, and uses high-confidence samples with predicted pseudo labels for weights updating. (3) BASS (Rottmann et al., 2018) is similar to CEAL, but use MC dropout for querying

<sup>3</sup>This is often neglected in previous low-resource AL studies, and we highlight it to ensure the true low-resource setting.

labeled sample. (4) REVIVAL (Guo et al., 2021) exploit momentum updates to stabilize the learning process, as there are oscillations in the beginning is the most recent SSAL method, which uses adversarial loss to query samples and leverage labeled bounds. In contrast, ACTUNE achieves a more stable learning process and enables an active self-propagation to exploit adversarial examples. Our Method: We experiment with both Entropy and CAL as uncertainty measures for ACTUNE. Note that when compared with active learning baselines, we do not augment the train set with pseudo-labeled data (Eq. (9)) to ensure fair comparisons.

### 3.2 Main Result

Figure 1 reports the performance of ACTUNE and the baselines on 4 benchmarks. From the results, we have the following observations:

ACTUNE consistently outperforms baselines in most of the cases. Different from studies in the computer vision (CV) domain (Siméoni et al., 2020) where the model does not perform well in the low-data regime, pre-trained LM has achieved competitive performance with only a few labeled data, which makes further improvements to the vanilla re-tuning challenging. Nevertheless, ACTUNE surpasses baselines in more than 90% of the rounds and achieves 0.4%-0.7% and 0.3%-1.5% absolute gain at the end of AL and SSAL respectively. Figure 2 quantitatively measures the number of labels needed for the most advanced active learning model and self-training model (UST) to outperform ACTUNE with 1000 labels. These baselines need >2000 clean labeled samples to reach the performance as our ACTUNE saves on average 56.2% and 57.0% of the labeled samples than most advanced active learning and self-training baselines respectively, which justifies its promising performance under low-resource scenarios. Such improvements show the merits of two key designs under our active self-training framework: the region-aware sampling for active learning and the momentum-based memory bank for robust self-training, which will be discussed in the section 3.5.

Compared with the previous AL baselines, ACTUNE can bring consistent performance gain, while previous semi-supervised active learning methods cannot. For instance, BASS is based on BALD for active learning, but sometimes it performs even worse than BALD with the same number of labeled data (see Fig. 5(b) and Fig. 1(f)). This is mainly because previous methods simply combine noisy pseudo labels with clean labels for training without explicitly rectifying the wrongly-labeled data, which will cause the LM to overfit these hazardous labels. Moreover, previous methods do not

demonstrate that adaptive selecting the most useful data for re-tuning is also important for improving the performance. With a powerful querying policy, ACTUNE can improve these self-training baselines by 1.05% in terms of accuracy on average.

### 3.3 Extension to Weakly-supervised Learning

ACTUNE can be naturally extended to weakly-supervised classification, where  $X_1$  is a set of data annotated by linguistic patterns or rules. Since the initial label set is noisy, then the model trained with Eq. (1) will overfit to the label noise, and we can actively query labeled data to refine the model.

We conduct experiments on the TREC and Chemprot datasets where we first use Snorkel (Ratner et al., 2017) to obtain weak label set then re-tune the pre-trained LM ( $\theta^{(0)}$ ) on  $X_1$ . After that, we adopt ACTUNE for active self-training.

Fig. 5 shows the results of these two datasets. When combining ACTUNE with CAL, the performance is unsatisfactory. We argue it is because CAL requires clean labels to calculate uncertainties. When labels are inaccurate, it will prevent ACTUNE from querying informative samples. In contrast, ACTUNE achieves the best performance over baselines when using Entropy as the uncertainty measure. The performance gain is more notable on the TREC dataset, where we achieve 96.68% accuracy, close to the fully supervised performance (96.80%) with only 6% of clean labels.

### 3.4 Combination with Other AL Methods

Fig. 4(a) demonstrates the performance of ACTUNE combined with other AL methods (e.g. BADGE, ALPS) on SST-2 dataset. It is clear that even if the AL methods are not uncertainty-based (e.g. BADGE), when using the entropy as an uncertainty measure to select pseudo-labeled data for

<sup>4</sup>Details for labeling functions are in Zhang et al. (2021).

<sup>5</sup>We don't show AL methods since they perform worse than SSAL methods on these datasets in general.

(a) SST-2, AL

(b) AG News, AL

(c) Pubmed, AL

(d) DBPedia, AL

(e) SST-2, SSAL

(f) AG News, SSAL

(g) Pubmed, SSAL

(h) DBPedia, SSAL

Figure 1: The comparison of ACTUNE with active learning, semi-supervised active learning and self-training baselines. The first row is the result under active learning setting (AL, i.e. no unlabeled data is used), the second row is the result under semi-supervised active learning (SSAL) setting. The metric is accuracy. RE-VIVAL causes OOM error for DBPedia dataset.

versely, removing WClus hurts the performance on later rounds, as it enables the model to select most informative samples.

Hyperparameter Study. We study two hyperparameters, namely  $m_L$  and  $m_H$  used in querying labels. Figure 6(e) and 6(f) show the results. In general, the model is insensitive to as the performance difference is less than 0.6%. The model cannot perform well with smaller  $m_L$  since it cannot pinpoint to high-uncertainty regions. For larger  $m_L$ , the performance also drops as some of the high-uncertainty regions can be outliers and sampling from them would hurt the model performance (Karamcheti et al., 2021).

Figure 2: The label-efficiency of ACTUNE compared with AL and self-training baselines. According to Fig. 1, the best AL method is Entropy for DBPedia and CAL for others.

self-training, ACTUNE can further boost the performance. This indicates that ACTUNE is a general active self-training approach, as it can serve as an efficient plug-in module for existing AL methods.

### 3.5 Ablation and Hyperparameter Study

The Effect of Different Components in ACTUNE. We inspect different components of ACTUNE, including the region-sampling (RS), momentum-based memory bank (MMB), and weighted clustering (WClus). Experimental results (Fig. 4(b)) shows that all the three components contribute to the final performance, as removing any of them hurts the classification accuracy. Also, we find that when removing MMB, the performance hurts most in the beginning rounds, which indicates that MMB effectively suppresses label noise when the model's capacity is weak. Con-

A Closer Look at the Momentum-based Memory Bank. To examine the role of MMB, we show the overall accuracy of pseudo-labels on AG News dataset in Fig. 6(g). From the result, it is clear that the momentum-based memory bank can stabilize the active self-training process, as the accuracy of pseudo labels increases around 1%, especially in later rounds. Fig 6(h) and 3(e) illustrates the model performance with different  $m_L$  and  $m_H$ . Overall, we find that our model is robust to different choices as ACTUNE outperform the baseline without momentum update consistently. Moreover, we find that the larger  $m_H$  will generally lead to better performance in later rounds. This is mainly because in later rounds, the model's prediction is more reliable. Conversely, at the beginning of the training, the model's prediction might be oscillating on unlabeled data. In this case, using a smaller  $m_L$  will favor samples with consistent predictions

<sup>6</sup>For models w/o RS, we directly select samples with highest uncertainty during AL. For models w/o MMB, we only use the prediction from the current round for self-training. For models w/o WClus, we cluster data with vanilla K-Means.

(a) Effect of (b) Effect of K (c) Acc. of PL (d) Entropy (e) CAL

Figure 3: Parameter study. Note the effect of different  $\alpha$  and  $m_H$  is conducted on AG News dataset.

(a) Combining w/ AL Methods (b) Ablation Study

Figure 4: Results of ACTUNE with different AL methods (SST-2), ablation study (SST-2 with ACTUNE+Entropy).

(a) TREC (b) Chemprot

Figure 5: The comparison of ACTUNE and baselines on weakly-supervised classification tasks.

to improve the robustness of active self-training.

Another finding is that for different AL methods, the optimal memory bank can be different. For Entropy, probability-based memory bank leads to a better result, while for CAL, simple aggregating over uncertainty score achieves better performance. This is mainly because the method used in CAL is more complicated, and using probability-based memory bank may hurt the uncertainty calculation.

## 4 Related Work

**Active Learning.** Active learning has been widely applied to various NLP tasks (Yuan et al., 2020; Zhao et al., 2020; Shelmanov et al., 2021; Karamcheti et al., 2021). So far, AL methods can be categorized into uncertainty-based methods (Gao et al., 2017; Margatina et al., 2021a,b), diversity-based methods (Ru et al., 2020; Sener and Savarese, 2018) and hybrid methods (Yuan et al., 2020; Ashby et al., 2020; Kirsch et al., 2019). Ein-Dor et al. (2020) offer an empirical study of active learning with PLMs. In our study, we leverage the power of unlabeled instances via self-training to further promote the performance of AL.

**Semi-supervised Active Learning (SSAL).** Gao et al. (2020); Song et al. (2019); Guo et al. (2021)

design query strategies for specific semi-supervised methods, Tomanek and Hahn (2009); Rottmann et al. (2018); Siméoni et al. (2020) exploit the most-certain samples from the unlabeled with pseudo-labeling to augment the training set. So far, most of the SSAL approaches are designed for CV domain and it remains unknown how this paradigm performs with PLMs on NLP tasks. In contrast, we propose ACTUNE to effectively leverage unlabeled data during finetuning PLMs for NLP tasks.

**Self-training.** Self-training first generates pseudo labels for high-confidence samples, then fits a new model on pseudo labeled data to improve the generalization ability (Rosenberg et al., 2005; Lee, 2013). However, it is known to be vulnerable to error propagation (Arazo et al., 2020; Rizve et al., 2021). To alleviate this, we adopt a simple momentum-based method to select high confidence samples, effectively reducing the pseudo labels noise for active learning. Note that although Mukherjee and Awadallah (2020); Rizve et al. (2021) also leverage uncertainty information for self-training, their focus is on developing better self-training methods, while we aim to jointly query high-uncertainty samples and generate pseudo-labels for low-uncertainty samples. The experiments in Sec. 3 show that with appropriate querying methods, ACTUNE can further improve the performance of self-training.

## 5 Conclusion

In this paper, we develop ACTUNE, a general active self-training framework for enhancing both label efficiency and model performance in finetuning pre-trained language models (PLMs). We propose a region-aware sampling approach to guarantee both the uncertainty the diversity for querying labels. To combat the label noise propagation issue, we design a momentum-based memory bank to effectively utilize the model predictions for preceding AL rounds. Empirical results on 6 public text classification benchmarks suggest the superiority of ACTUNE to conventional active learning and semi-supervised active learning methods for finetuning PLMs with limited resources.



645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700

## References

Eric Arazo, Diego Ortego, Paul Albert, Noel E O'Connor, and Kevin McGuinness. 2020. Pseudo-labeling and confirmation bias in deep semi-supervised learning. In 2020 International Joint Conference on Neural Networks (IJCNN), pages 1–8. IEEE.

David Arthur and Sergei Vassilvitskii. 2007. K-means++: The advantages of careful seeding. In Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms, pages 1027–1035, USA.

Jordan T. Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. 2020. Deep batch active learning by diverse, uncertain gradient lower bounds. In International Conference on Learning Representations.

Trapit Bansal, Rishikesh Jha, Tsendsuren Munkhdalai, and Andrew McCallum. 2020. 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.

Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. SciBERT: A pretrained language model for scientific 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 3615–3620.

Jonathan Bragg, Arman Cohan, Kyle Lo, and Iz Beltagy. 2021. Flex: Unifying evaluation for few-shot nlp. Advances in Neural Information Processing Systems, 34.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, et al. 2020. Language models are few-shot learners. In Advances in Neural Information Processing Systems, volume 33, pages 1877–1901.

Franck Dernoncourt and Ji Young Lee. 2017. PubMed 200k RCT: a dataset for sequential sentence classification in medical abstracts. In Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 2: Short Papers), pages 308–313, Taipei, Taiwan. Asian Federation of Natural Language Processing.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Jesse Dodge, Gabriel Ilharco, Roy Schwartz, Ali Farhadi, Hannaneh Hajishirzi, and Noah A. Smith. 2020. Fine-tuning pretrained language models: Weight initializations, data orders, and early stopping. CoRR, abs/2002.06305.

Rotem Dror, Gili Baumer, Segev Shlomov, and Roi Reichart. 2018. The hitchhiker's guide to testing statistical significance in natural language processing. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1383–1392.

Jingfei Du, Edouard Grave, Beliz Gunel, Vishrav Chaudhary, Onur Celebi, Michael Auli, Veselin Stoyanov, and Alexis Conneau. 2021. Self-training improves pre-training for natural language understanding. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 5408–5418. Association for Computational Linguistics.

Liat Ein-Dor, Alon Halfon, Ariel Gera, Eyal Shnarch, Lena Dankin, Leshem Choshen, Marina Danilevsky, Ranit Aharonov, Yoav Katz, and Noam Slonim. 2020. Active Learning for BERT: An Empirical Study. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 7949–7962. Association for Computational Linguistics.

Yarin Gal and Zoubin Ghahramani. 2015. Bayesian convolutional neural networks with bernoulli approximate variational inference. CoRR, abs/1506.02158.

Yarin Gal, Riashat Islam, and Zoubin Ghahramani. 2017. Deep bayesian active learning with image data. In International Conference on Machine Learning, pages 1183–1192. PMLR.

Mingfei Gao, Zizhao Zhang, Guo Yu, Serkan Özyer, Larry S Davis, and Tomas Pfister. 2020. Consistency-based semi-supervised active learning: Towards minimizing labeling cost. In European Conference on Computer Vision, pages 510–526. Springer.

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

Jiannan Guo, Haochen Shi, Yangyang Kang, Kun Kuang, Siliang Tang, Zhuoren Jiang, Changlong Sun, Fei Wu, and Yueting Zhuang. 2021. Semi-supervised active learning for semi-supervised models: Exploit adversarial examples with graph-based virtual labels. In Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), pages 2896–2905.

759	Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and	Xin Li and Dan Roth. 2002. <a href="#">Learning question classi-</a>	815
760	Ross Girshick. 2020. Momentum contrast for unsu-	<a href="#">fers</a> . In The 19th International Conference on Com-	816
761	perervised visual representation learning. <a href="#">Proceed-</a>	putational Linguistics	817
762	<a href="#">ings of the IEEE/CVF Conference on Computer Vi-</a>		
763	sion and Pattern Recognition (CVPR)	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	818
764	Alex Holub, Pietro Perona, and Michael C Burl. 2008.	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	819
765	Entropy-based active learning for object recogni-	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	820
766	tion. In 2008 IEEE Computer Society Conference	Roberta: A robustly optimized bert pretraining ap-	821
767	on Computer Vision and Pattern Recognition Work-	proach. <a href="#">arXiv preprint arXiv:1907.11692</a>	822
768	shops pages 1–8. IEEE.	Ilya Loshchilov and Frank Hutter. 2019. <a href="#">Decoupled</a>	823
769	Peiyun Hu, Zack Lipton, Anima Anandkumar, and	<a href="#">weight decay regularization</a> . <a href="#">International Con-</a>	824
770	Deva Ramanan. 2019. <a href="#">Active learning with partial</a>	ference on Learning Representations	825
771	<a href="#">feedback</a> . In <a href="#">International Conference on Learning</a>		
772	<a href="#">Representations</a>	Katerina Margatina, Loic Barrault, and Nikolaos	826
773	Joshua Zhexue Huang, Michael K Ng, Hongqiang	Aletras. 2021a. Bayesian active learning with	827
774	Rong, and Zichen Li. 2005. Automated variable	pretrained language models. <a href="#">arXiv preprint</a>	828
775	weighting in k-means type clustering. <a href="#">IEEE transac-</a>	<a href="#">arXiv:2104.08320</a>	829
776	tions on pattern analysis and machine intelligence	Katerina Margatina, Giorgos Vernikos, Loic Barrault,	830
777	27(5):657–668.	and Nikolaos Aletras. 2021b. <a href="#">Active learning by</a>	831
778	Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019.	<a href="#">acquiring contrastive examples</a> . <a href="#">Proceedings of</a>	832
779	Billion-scale similarity search with gpus. <a href="#">IEEE</a>	of the 2021 Conference on Empirical Methods in Nat-	833
780	<a href="#">Transactions on Big Data</a>	ural Language Processing pages 650–663, Online	834
781	Katharina Kann, Kyunghyun Cho, and Samuel R. Bow-	and Punta Cana, Dominican Republic. Association	835
782	man. 2019. <a href="#">Towards realistic practices in low-</a>	for Computational Linguistics.	836
783	<a href="#">resource natural language processing: The develop-</a>	Marius Mosbach, Maksym Andriushchenko, and Diet-	837
784	<a href="#">ment set</a> . In <a href="#">Proceedings of the 2019 Conference on</a>	rich Klakow. 2021. <a href="#">On the stability of ne-tuning</a>	838
785	<a href="#">Empirical Methods in Natural Language Processing</a>	<a href="#">(bert): Misconceptions, explanations, and strong</a>	839
786	and the 9th International Joint Conference on Natu-	<a href="#">baselines</a> . In <a href="#">International Conference on Learning</a>	840
787	ral Language Processing (EMNLP-IJCNLP) pages	<a href="#">Representations</a>	841
788	3342–3349, Hong Kong, China. Association for	Subhabrata Mukherjee and Ahmed Awadallah. 2020.	842
789	Computational Linguistics.	Uncertainty-aware self-training for few-shot text	843
790	Siddharth Karamcheti, Ranjay Krishna, Li Fei-Fei, and	classification. <a href="#">Advances in Neural Information Pro-</a>	844
791	Christopher Manning. 2021. <a href="#">Mind your outliers! in-</a>	cessing System	845
792	<a href="#">vestigating the negative impact of outliers on active</a>	Fabian Pedregosa, Gaël Varoquaux, Alexandre Gram-	846
793	<a href="#">learning for visual question answering</a> . <a href="#">Proceed-</a>	fort, Vincent Michel, Bertrand Thirion, Olivier	847
794	<a href="#">ings of the 59th Annual Meeting of the Association</a>	Grisel, Mathieu Blondel, Peter Prettenhofer, Ron	848
795	for Computational Linguistics and the 11th Interna-	Weiss, Vincent Dubourg, et al. 2011. Scikit-learn:	849
796	tional Joint Conference on Natural Language Pro-	Machine learning in pytho	850
797	cessing (Volume 1: Long Papers) pages 7265–7281,	the <a href="#">Journal of machine</a>	851
798	Online. Association for Computational Linguistics.	<a href="#">Learning research</a> 12:2825–2830.	
799	Andreas Kirsch, Joost Van Amersfoort, and Yarin Gal.	Alexander Ratner, Stephen H Bach, Henry Ehrenberg,	852
800	2019. Batchbald: Efficient and diverse batch acqui-	Jason Fries, Sen Wu, and Christopher Ré. 2017.	853
801	sition for deep bayesian active learning. <a href="#">Advances</a>	Snorkel: Rapid training data creation with weak su-	854
802	in neural information processing systems	pervision. In <a href="#">Proceedings of the VLDB Endowment</a> .	855
803	32:7026–7037.	volume 11, page 269.	856
804	Martin Krallinger, Obdulia Rabal, Saber A Akhondi,	Mamshad Nayeem Rizve, Kevin Duarte, Yogesh S	857
805	et al. 2017. Overview of the biocreative VI	Rawat, and Mubarak Shah. 2021. <a href="#">In defense</a>	858
806	chemical-protein interaction track. <a href="#">BioCreative</a>	<a href="#">of pseudo-labeling: An uncertainty-aware pseudo-</a>	859
807	evaluation Workshop volume 1, pages 141–146.	<a href="#">label selection framework for semi-supervised learn-</a>	860
808	Samuli Laine and Timo Aila. 2016. Temporal ensem-	<a href="#">ing</a> . In <a href="#">International Conference on Learning Rep-</a>	861
809	bling for semi-supervised learning. <a href="#">arXiv preprint</a>	<a href="#">resentations</a>	862
810	<a href="#">arXiv:1610.02242</a>	Chuck Rosenberg, Martial Hebert, and Henry Schnei-	863
811	Dong-Hyun Lee. 2013. Pseudo-label: The simple and	derman. 2005. Semi-supervised self-training of	864
812	efficient semi-supervised learning method for deep	object detection models. In <a href="#">Proceedings of the</a>	865
813	neural networks. In <a href="#">ICML Workshop on challenges</a>	<a href="#">IEEE Workshops on Application of Computer Vision</a>	866
814	in representation learning volume 3, page 896.	pages 29–36.	867

868	Matthias Rottmann, Karsten Kahl, and Hanno	Laurens Van der Maaten and Geoffrey Hinton. 2008.	925
869	Gottschalk. 2018. Deep bayesian active semi-	Visualizing data using t-sne. <i>Journal of machine</i>	926
870	supervised learning. <i>2018 17th IEEE Interna-</i>	learning research	927
871	tional Conference on Machine Learning and Appli-		
872	cations (ICMLA) pages 158–164. IEEE.		
873	Dongyu Ru, Jiangtao Feng, Lin Qiu, Hao Zhou, Mingx-	Keze Wang, Dongyu Zhang, Ya Li, Ruimao Zhang,	928
874	uan Wang, Weinan Zhang, Yong Yu, and Lei Li.	and Liang Lin. 2016. Cost-effective active learn-	929
875	2020. <a href="#">Active sentence learning by adversarial un-</a>	ing for deep image classi cation. <i>IEEE Transac-</i>	930
876	<a href="#">certainty sampling in discrete space.</a> <i>Findings</i>	tions on Circuits and Systems for Video Technology	931
877	of the Association for Computational Linguistics:	27(12):2591–2600.	932
878	EMNLP 2020 pages 4908–4917, Online. Associa-		
879	tion for Computational Linguistics.	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	933
880	Timo Schick and Hinrich Schütze. 2021. <a href="#">Exploiting</a>	Chaumond, Clement Delangue, Anthony Moi, Pier-	934
881	<a href="#">cloze-questions for few-shot text classi cation and</a>	ric Cistac, Tim Rault, Remi Louf, Morgan Funtow-	935
882	<a href="#">natural language inference.</a> <i>Proceedings of the</i>	icz, Joe Davison, Sam Shleifer, Patrick von Platen,	936
883	16th Conference of the European Chapter of the As-	Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,	937
884	sociation for Computational Linguistics: Main Vol-	Teven Le Scao, Sylvain Gugger, Mariama Drame,	938
885	ume pages 255–269, Online. Association for Com-	Quentin Lhoest, and Alexander Rush. 2020. <a href="#">Trans-</a>	939
886	putational Linguistics.	<a href="#">formers: State-of-the-art natural language process-</a>	940
887	Ozan Sener and Silvio Savarese. 2018. <a href="#">Active learn-</a>	<a href="#">ing.</a> In <i>Proceedings of the 2020 Conference on Em-</i>	941
888	<a href="#">ing for convolutional neural networks: A core-set</a>	<i>pirical Methods in Natural Language Processing:</i>	942
889	<a href="#">approach.</a> <i>International Conference on Learning</i>	<i>System Demonstrations</i> pages 38–45, Online. Asso-	943
890	<i>Representations</i>	ciation for Computational Linguistics.	944
891	Artem Shelmanov, Dmitri Puzyrev, Lyubov	Yue Yu, Simiao Zuo, Haoming Jiang, Wendi Ren, Tuo	945
892	Kupriyanova, Denis Belyakov, Daniil Larionov,	Zhao, and Chao Zhang. 2021. <a href="#">Fine-tuning pre-</a>	946
893	Nikita Khromov, Olga Kozlova, Ekaterina Arte-	<a href="#">trained language model with weak supervision: A</a>	947
894	mova, Dmitry V. Dylov, and Alexander Panchenko.	<a href="#">contrastive-regularized self-training approach.</a> In	948
895	2021. <a href="#">Active learning for sequence tagging with</a>	<i>Proceedings of the 2021 Conference of the North</i>	949
896	<a href="#">deep pre-trained models and Bayesian uncertainty</a>	<i>American Chapter of the Association for Computa-</i>	950
897	<a href="#">estimates.</a> In <i>Proceedings of the 16th Conference</i>	<i>tional Linguistics: Human Language Technologies</i>	951
898	<i>of the European Chapter of the Association for</i>	pages 1063–1077. Association for Computational	952
899	<i>Computational Linguistics: Main Volume</i>	Linguistics.	953
900	pages	Michelle Yuan, Hsuan-Tien Lin, and Jordan Boyd-	954
901	1698–1712, Online. Association for Computational	Graber. 2020. <a href="#">Cold-start active learning through</a>	955
902	Linguistics.	<a href="#">self-supervised language modeling.</a> <i>Proceed-</i>	956
903	Oriane Siméoni, Mateusz Budnik, Yannis Avrithis, and	<i>ings of the 2020 Conference on Empirical Methods</i>	957
904	Guillaume Gravier. 2020. Rethinking deep active	<i>in Natural Language Processing (EMNLP)</i> pages	958
905	learning: Using unlabeled data at model training. In	7935–7948, Online. Association for Computational	959
906	the 25th International Conference on Pattern Recog-	Linguistics.	960
907	nition (ICPR), pages 1220–1227. IEEE.	Jieyu Zhang, Yue Yu, Yinghao Li, Yujing Wang, Yam-	961
908	Richard Socher, Alex Perelygin, Jean Wu, Jason	ing Yang, Mao Yang, and Alexander Ratner. 2021.	962
909	Chuang, Christopher D. Manning, Andrew Ng, and	<a href="#">WRENCH: A comprehensive benchmark for weak</a>	963
910	Christopher Potts. 2013. <a href="#">Recursive deep models</a>	<a href="#">supervision.</a> In <i>Thirty- fth Conference on Neural In-</i>	964
911	<a href="#">for semantic compositionality over a sentiment tree-</a>	<i>formation Processing Systems Datasets and Bench-</i>	965
912	<a href="#">bank.</a> In <i>Proceedings of the 2013 Conference on</i>	<i>marks Track</i>	966
913	<i>Empirical Methods in Natural Language Processing</i>	Tianyi Zhang, Felix Wu, Arzoo Katiyar, Kilian Q	967
914	pages 1631–1642. Association for Computational	Weinberger, and Yoav Artzi. 2020. Revisit-	968
915	Linguistics.	ing few-sample bert ne-tuning. <i>arXiv preprint</i>	969
916	Shuang Song, David Berthelot, and Afshin Ros-	<i>arXiv:2006.05987</i>	970
917	tamizadeh. 2019. Combining mixmatch and active	Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015.	971
918	learning for better accuracy with fewer labels. <i>arXiv</i>	Character-level convolutional networks for text clas-	972
919	<i>preprint arXiv:1912.00594</i>	sification. <i>Advances in neural information process-</i>	973
920	Katrin Tomanek and Udo Hahn. 2009. Semi-	<i>ing systems</i> 28:649–657.	974
921	supervised active learning for sequence labeling. In	Yuekai Zhao, Haoran Zhang, Shuchang Zhou, and Zhi-	975
922	<i>Proceedings of the Joint Conference of the 47th An-</i>	hua Zhang. 2020. <a href="#">Active learning approaches to</a>	976
923	<i>annual Meeting of the ACL and the 4th International</i>	<a href="#">enhancing neural machine translation.</a> <i>Findings</i>	977
924	<i>Joint Conference on Natural Language Processing</i>	of the Association for Computational Linguistics:	978
	<i>of the AFNLP</i> pages 1039–1047.	EMNLP 2020 pages 1796–1806, Online. Associa-	979
		tion for Computational Linguistics.	980

981	A Datasets Details	B.2 Number of Parameters	1026
982	A.1 Data Source	ACTUNE and all baselines use Roberta-base (Liu et al., 2019) with a task-specific classification head	1027
983	The seven benchmarks in our experiments are all publicly available. Below are the links to downloadable versions of these datasets.	on the top as the backbone, which contains 125M trainable parameters. We do not introduce any other parameters in our experiments.	1028
984			1029
985			1030
986	SST-2 We use the datasets from <a href="https://huggingface.co/datasets/glue">https://huggingface.co/datasets/glue</a>	B.3 Experiment Setups	1031
987			1032
988	AGNews We use the datasets from <a href="https://huggingface.co/datasets/ag_news">https://huggingface.co/datasets/ag_news</a>	Following (Ein-Dor et al., 2020; Yuan et al., 2020; Margatina et al., 2021b), all of our methods and baselines are run with 3 different random seed and the result is based on the average performance on them. This indeed creates (the number of datasets) 3 (the number of random seeds) 11 (the number of methods) 10 (the number of ne-tuning rounds in AL) =1320 experiments for ne-tuning PLMs, which is almost the limit of our computational resources, not to mention additional experiments on weakly-supervised text classification as well as different hyper-parameter tuning.	1033
989	Pubmed-RCT Dataset is available at <a href="https://github.com/Franck-Dernoncourt/pubmed-rct">https://github.com/Franck-Dernoncourt/pubmed-rct</a>		1034
990			1035
991			1036
992	DBPedia Dataset is available at <a href="https://huggingface.co/datasets/dbpedia_14">https://huggingface.co/datasets/dbpedia_14</a>		1037
993			1038
994			1039
995			1040
996	For two weakly-supervised classification tasks, we use the data from WRENCH benchmark (Zhang et al., 2021).		1041
997			1042
998			1043
999	TREC Dataset is available at <a href="https://drive.google.com/drive/u/1/folders/1v55IKG2JN9fMtKJWU48B_5_DcPWGnpTq">https://drive.google.com/drive/u/1/folders/1v55IKG2JN9fMtKJWU48B_5_DcPWGnpTq</a>	We have show both the mean and the standard deviation of the performance in our experiment sections. All the results have passed a paired t-test with $p < 0:05$ (Dror et al., 2018).	1044
1000			1045
1001			1046
1002			1047
1003	ChemProt The raw dataset is available at <a href="http://www.cbs.dtu.dk/services/ChemProt/ChemProt-2.0/">http://www.cbs.dtu.dk/services/ChemProt/ChemProt-2.0/</a>	B.4 Implementations Baselines	1048
1004			1049
1005	The preprocessed dataset is available at <a href="https://drive.google.com/drive/u/1/folders/1v55IKG2JN9fMtKJWU48B_5_DcPWGnpTq">https://drive.google.com/drive/u/1/folders/1v55IKG2JN9fMtKJWU48B_5_DcPWGnpTq</a>	We implement Entropy, BALD by ourselves as they are easy to implement and are classic methods for AL. For REVIVAL (Guo et al., 2021), since we do not find the implementations released by authors, we implement on our own it based on the information in the original paper. For other baselines, we run the experiments based on the implementations on the web. We list the link for the implementations as follows:	1050
1006			1051
1007			1052
1008			1053
1009			1054
1010	A.2 Train/Test Split		1055
1011	For all the datasets, we use the original train/dev/test split from the web. To keep the size of the development set small, we randomly sample 1000 data for SST-2, AGNews, Pubmed-RCT, DBPedia and randomly sample 500 samples for TREC, ChemProt	BADGE <a href="https://github.com/JordanAsh/badge">https://github.com/JordanAsh/badge</a>	1056
1012			1057
1013		ALPS <a href="https://github.com/forest-snow/alps">https://github.com/forest-snow/alps</a>	1058
1014			1059
1015		CAL: <a href="https://github.com/mourga/contrastive-active-learning">https://github.com/mourga/contrastive-active-learning</a>	1060
1016			1061
1017	B Details on Implementation and Experiment Setups	UST: <a href="https://github.com/microsoft/UST">https://github.com/microsoft/UST</a>	1062
1018			1063
1019	B.1 Computing Infrastructure	COSINE <a href="https://github.com/yueyu1030/COSINE">https://github.com/yueyu1030/COSINE</a>	1064
1020	System Ubuntu 18.04.3 LTS; Python 3.6; Pytorch 1.6.		1065
1021	CPU: Intel(R) Core(TM) i7-5930K CPU @ 3.50GHz.	For these three baselines listed below, since they are mainly used in CV tasks, thus the code is hard to directly used for our experiments. We re-implement these methods based on their implementations, especially for SSAL part.	1066
1022	GPU: NVIDIA 2080Ti.		1067
1023		ASST <a href="https://">https://</a>	1068
1024			1069
1025			1070
			1071
			1072
			1073
			1074

Hyper-parameter	SST-2	AG News	Pubmed	DBPedia	TREC	Chemprot
Dropout Ratio	0.1					
Maximum Tokens	32	96	96	64	64	128
Batch Size for $X_l$	8					
Batch Size for $X_u$ in Self-training	32	48	48	32	16	24
Weight Decay	$10^{-8}$					
Learning Rate	$2 \cdot 10^{-5}$					
	0.5					
M	25	30	30	40	40	40
K	5	10				
	0.7	0.6				
$m_L$	0.8	0.9	0.7	0.8	0.8	0.8
$m_H$	0.9	0.9	0.8	0.9	0.9	1.0
	1					

Table 2: Hyper-parameter configurations. Note that we only keep certain number of tokens.

Method	Dataset	
	Pubmed	DBPedia
Finetune (Random)	<0.1s	<0.1s
Entropy (Holub et al., 2008)	461s	646s
BALD (Gal et al., 2017)	4595s	6451s
ALPS (Yuan et al., 2020)	488s	677s
BADGE (Ash et al., 2020)	554s	1140s
CAL (Margatina et al., 2021b)	493s	688s
REVIVAL (Guo et al., 2021)	3240s	OOM
ACTUNE + Entropy	477s	733s
w/ RS for Active Learning	15.8s	44.9s
w/ MMB for Self-training	0.12s	0.18s
ACTUNE + CAL	510s	735s
w/ RS for Active Learning	16.6s	46.4s
w/ MMB for Self-training	0.12s	0.18s

Table 3: The running time of different baselines. Note that for ASST, CEAL and BASS, they directly use existing active learning methods so we do not list the running time here.

1075 [github.com/osimeoni/](https://github.com/osimeoni/RethinkingDeepActiveLearning)  
1076 [RethinkingDeepActiveLearning](https://github.com/osimeoni/RethinkingDeepActiveLearning)  
1077 CEAL: [https://github.com/rafikg/](https://github.com/rafikg/CEAL)  
1078 [CEAL](https://github.com/rafikg/CEAL)  
1079 BASS: [https://github.com/](https://github.com/mrottmann/DeepBASS)  
1080 [mrottmann/DeepBASS](https://github.com/mrottmann/DeepBASS)  
1081 Our implementation of ACTUNE will be published upon acceptance.

### 1083 B.5 Hyper-parameters for General 1084 Experiments

1085 We use AdamW (Loshchilov and Hutter, 2019) as the optimizer, and the learning rate is chosen from  
1086  $1 \cdot 10^{-5}; 2 \cdot 10^{-5}$ . A linear learning rate decay schedule with warm-up:1 is used, and the number of training epochs is 5 for fine-tuning. For active

self-training & SSAL baselines, we tune the model with 2000 steps, and evaluate the performance on the development set in every 50 steps. Finally, we use the model with best performance on the development set for testing.

### 1095 B.6 Hyper-parameters for ACTUNE

1096 Although ACTUNE introduces several hyper-parameters including  $K, M, m_L, m_H, \alpha, \beta, \gamma$ , most of them are kept fixed during our experiments, thus it does not require heavy hyper-parameter tuning. The hyper-parameters we use are shown in Table 2. Specifically, we search  $m_L$  from 10 to 2000,  $T_2$  from 1000 to 5000,  $T_3$  from 10 to 500,  $\alpha$  from 0 to 1, and  $\beta$  from 0 to 0.5. All results are reported as the average over three runs.

1105 In our experiments, we keep  $\alpha = 0.5, \beta = 1$  for all datasets. For other parameters, we use a grid search to find the optimal setting for each dataset. Specifically, we search  $\gamma$  from [0.5; 0.6; 0.7],  $m_L$  from [0.6; 0.7; 0.8],  $m_H$  from [0.8; 0.9; 1]. For ACTUNE with Entropy, we use probability based aggregation and for ACTUNE with CAL, we use value based aggregation by default.

### 1113 B.7 Hyperparameters for Baselines

1114 For other SSAL methods, we mainly tune their key hyperparameters. Note that Entropy (Holub et al., 2008), BALD (Gal et al., 2017), ALPS (Yuan et al., 2020), BADGE (Ash et al., 2020) do not introduce any new hyperparameters. For CAL (Margatina et al., 2021b), we tune the number for KNN  $k$  from [5; 10; 20] and report the best performance. For ST (Lee, 2013), CEAL (Wang

et al., 2016) & BASS (Rottmann et al., 2018), it 2019) which shown even better performance with  
uses a threshold for selecting high-con dence only a few labels (Du et al., 2021). Last, apart from  
data. We tune from [0:6; 0:7; 0:8; 0:9] to report the text classi cation task, we can also extend our  
the best performance. For UST (Mukherjee and Awadallah, 2020), we tune the number of low-natural language inference.

uncertainty samples used in the next round from  
[1024; 2048; 4096] For COSINE (Yu et al., 2021),  
we set the weight for con dence regularization as  
0:1, the threshold for selecting high-con dence  
data from [0:7; 0:9] and the update period of self-  
training from [50; 100; 150]. For REVIVAL (Guo  
et al., 2021), it calculates uncertainty with adversarial  
perturbation, we tune the size of the perturbation  
from [1e -3; 1e -4; 1e -5].

## C Runtime Analysis.

Table 3 shows the time in one active learning  
round of ACTUNE and baselines. Here we high-  
light that the additional time for region-aware sam-  
pling and momentum-based memory banks is  
small compared with the inference time. Among  
all baselines, we nd that the running time of  
clustering-based method is faster than the origi-  
nal reported time in the paper. This is because  
we use FAISS (Johnson et al., 2019) instead of  
SKLearn (Pedregosa et al., 2011) for clustering,  
which accelerates the clustering step signi cantly.  
Also, we nd that BALD and REVIVAL are not  
so ef cient. For BALD, it needs to infer the uncer-  
tainty of the model by passing the data to model  
with multiple times. Such an operation will make  
the total inference time for PLMs very long. For  
REVIVAL, we nd that calculating the adversarial  
gradient needs extra forward passes and backward  
passes, which could be time-consuming for PLMs  
with millions of parameters.

## D Limitations

First, since our focus is on ne-tuning pre-trained  
language models, we use the representation of  
[CLS] token for classi cation. In the future work,  
we can consider using prompt tuning (Gao et al.,  
2021; Schick and Schütze, 2021), a more data-  
ef cient method for adopting pre-trained language  
models on classi cation tasks to further promote  
the ef ciency. Also, due to the computational re-  
source constraints, we do not use larger pre-trained  
language models such as RoBERTa-large (Liu et al.,

<sup>7</sup>The original model is proposed with CV tasks and they  
use ResNet-18 as the backbone which only contains 11M  
parameters (around 10% of the parameters of Roberta-base).

