

On the Importance of Effectively Adapting Pretrained Language Models for Active Learning

Anonymous ACL submission

Abstract

Recent active learning (AL) approaches in Natural Language Processing (NLP) proposed using off-the-shelf pretrained language models (LMs). In this paper, we argue that these LMs are not adapted effectively to the downstream task during AL and we explore ways to address this issue. We suggest to first adapt the pretrained LM to the target task by continuing training with all the available *unlabeled* data and then use it for AL. We also propose a simple yet effective fine-tuning method to ensure that the adapted LM is properly trained in both low and high resource scenarios during AL. Our experiments demonstrate that our approach provides substantial data efficiency improvements compared to the standard fine-tuning approach, suggesting that a poor training strategy can be catastrophic for AL.

1 Introduction

Active Learning (AL) is a method for training supervised models in a data-efficient way (Cohn et al., 1996; Settles, 2009). AL methods iteratively alternate between (i) model training with the labeled data available; and (ii) data selection for annotation using a stopping criterion, e.g. until exhausting a fixed annotation budget or reaching a pre-defined performance on a held-out dataset.

Data selection is performed by an acquisition function that ranks unlabeled data points by some *informativeness* metric aiming to improve over random selection, using either uncertainty (Lewis and Gale, 1994; Cohn et al., 1996), diversity (Brinker, 2003; Bodó et al., 2011; Sener and Savarese, 2018), or both (Ducoffe and Precioso, 2018; Ash et al., 2020; Yuan et al., 2020; Margatina et al., 2021).

Previous AL approaches in NLP use task-specific neural models that are trained from scratch at each iteration (Shen et al., 2017; Siddhant and Lipton, 2018; Prabhu et al., 2019; Ikhwantri et al., 2018; Kasai et al., 2019). However, these models

are usually outperformed by pretrained language models (LMs) adapted to end-tasks (Howard and Ruder, 2018), making them suboptimal for AL. Only recently, pretrained LMs such as BERT (Devlin et al., 2019) have been introduced in AL settings (Yuan et al., 2020; Ein-Dor et al., 2020; Shelmanov et al., 2021; Karamcheti et al., 2021; Margatina et al., 2021). Still, they are trained at each AL iteration with a standard fine-tuning approach that mainly includes a pre-defined number of training epochs, which has been demonstrated to be unstable, especially in small datasets (Zhang et al., 2020; Dodge et al., 2020; Mosbach et al., 2021). Since AL includes both low and high data resource settings, the AL model training scheme should be robust in both scenarios.¹

To address these limitations, we introduce a suite of effective training strategies for AL (§2). Contrary to previous work (Yuan et al., 2020; Ein-Dor et al., 2020; Margatina et al., 2021) that also use BERT (Devlin et al., 2019), our proposed method accounts for various data availability settings and the instability of fine-tuning. First, we continue *pretraining* the LM with the available *unlabeled* data to adapt it to the task-specific domain. This way, we leverage not only the available labeled data at each AL iteration, but the entire unlabeled pool. Second, we further propose a simple yet effective fine-tuning method that is robust in both low and high resource data settings for AL.

We explore the effectiveness of our approach on five natural language understandings tasks with various acquisition functions, showing that it outperforms all baselines (§3). We also conduct an analysis to demonstrate the importance of adaptation of pretrained models for AL (§4). Our findings highlight that the LM adaptation strategy can be more critical than the data acquisition strategy.

¹During the first few AL iterations the available labeled data is limited (*low-resource*), while it could become very large towards the last iterations (*high-resource*).

2 Adapting & Fine-tuning Pretrained Models for Active Learning

Given a classification task with C classes, a typical AL setup consists of a pool of unlabeled data $\mathcal{D}_{\text{pool}}$, a model \mathcal{M} , an annotation budget b of data points and an acquisition function $a(\cdot)$ for selecting k unlabeled data points for annotation (i.e. acquisition size) until b runs out. The AL performance is assessed by training a model on the actively acquired dataset and evaluating on a held-out test set $\mathcal{D}_{\text{test}}$.

Adaptation (TAPT) Inspired by recent work on transfer learning that shows improvements in downstream classification performance by continuing the pretraining of the LM with the task data (Howard and Ruder, 2018) we add an extra step to the AL process by continuing pretraining the LM (i.e. Task-Adaptive Pretraining TAPT), as in Gururangan et al. (2020). Formally, we use an LM, such as BERT (Devlin et al., 2019), $\mathcal{P}(x; W_0)$ with weights W_0 , that has been already pretrained on a large corpus. We fine-tune $\mathcal{P}(x; W_0)$ with the available unlabeled data of the downstream task $\mathcal{D}_{\text{pool}}$, resulting in the task-adapted LM $\mathcal{P}_{\text{TAPT}}(x; W'_0)$ with new weights W'_0 (cf. line 2 of algorithm 1).

Fine-tuning (FT+) We now use the adapted LM $\mathcal{P}_{\text{TAPT}}(x; W'_0)$ for AL. At each iteration i , we initialize our model \mathcal{M}_i with the pretrained weights W'_0 and we add a task-specific feedforward layer for classification with weights W_c on top of the [CLS] token representation of BERT-based $\mathcal{P}_{\text{TAPT}}$. We fine-tune the classification model $\mathcal{M}_i(x; [W'_0, W_c])$ with all $x \in \mathcal{D}_{\text{lab}}$. (cf. line 6 to 8 of algorithm 1).

Recent work in AL (Ein-Dor et al., 2020; Yuan et al., 2020) uses the standard fine-tuning method proposed in Devlin et al. (2019) which includes a fixed number of 3 training epochs, learning rate warmup over the first 10% of the steps and AdamW optimizer (Loshchilov and Hutter, 2019) without bias correction, among other hyperparameters.

We follow a different approach by taking into account insights from few-shot fine-tuning literature (Mosbach et al., 2021; Zhang et al., 2020; Dodge et al., 2020) that proposes longer fine-tuning and more evaluation steps during training. We combine these guidelines to our fine-tuning approach by using early stopping with 20 epochs based on the validation loss, learning rate $2e - 5$, bias correction and 5 evaluation steps per epoch. However, increasing the number of epochs from 3 to

Algorithm 1: AL with Pretrained LMs

Input: unlabeled data $\mathcal{D}_{\text{pool}}$, pretrained LM $\mathcal{P}(x; W_0)$, acquisition size k , AL iterations T , acquisition function a

- 1 $\mathcal{D}_{\text{lab}} \leftarrow \emptyset$
- 2 $\mathcal{P}_{\text{TAPT}}(x; W'_0) \leftarrow \text{Train } \mathcal{P}(x; W_0) \text{ on } \mathcal{D}_{\text{pool}}$
- 3 $\mathcal{Q}_0 \leftarrow \text{RANDOM}(\cdot), |\mathcal{Q}_0| = k$
- 4 $\mathcal{D}_{\text{lab}} = \mathcal{D}_{\text{lab}} \cup \mathcal{Q}_0$
- 5 $\mathcal{D}_{\text{pool}} = \mathcal{D}_{\text{pool}} \setminus \mathcal{Q}_0$
- 6 **for** $i \leftarrow 1$ **to** T **do**
- 7 $\mathcal{M}_i(x; [W'_0, W_c]) \leftarrow \text{Initialize from } \mathcal{P}_{\text{TAPT}}(x; W'_0)$
- 8 $\mathcal{M}_i(x; W_i) \leftarrow \text{Train model on } \mathcal{D}_{\text{lab}}$
- 9 $\mathcal{Q}_i \leftarrow a(\mathcal{M}_i, \mathcal{D}_{\text{pool}}, k)$
- 10 $\mathcal{D}_{\text{lab}} = \mathcal{D}_{\text{lab}} \cup \mathcal{Q}_i$
- 11 $\mathcal{D}_{\text{pool}} = \mathcal{D}_{\text{pool}} \setminus \mathcal{Q}_i$
- 12 **end**

Output: \mathcal{D}_{lab}

DATASETS	TRAIN	VAL	TEST	k	C
TREC-6	4.9K	546	500	1%	6
DBPEDIA	20K	2K	70K	1%	14
IMDB	22.5K	2.5K	25K	1%	2
SST-2	60.6K	6.7K	871	1%	2
AGNEWS	114K	6K	7.6K	0.5%	4

Table 1: Datasets statistics for $\mathcal{D}_{\text{pool}}$, \mathcal{D}_{val} and $\mathcal{D}_{\text{test}}$ respectively. k stands for the acquisition size (% of $\mathcal{D}_{\text{pool}}$) and C the number of classes.

20, also increases the warmup steps (10% of total steps²) almost 7 times. This may be problematic in scenarios where the dataset is large but the optimal number of epochs may be small (e.g. 2 or 3). To account for this limitation in our AL setting where the size of training set changes at each iteration, we propose to select the warmup steps as $\min(10\% \text{ of total steps}, 100)$. We denote standard fine-tuning as SFT and our approach as FT+.

3 Experiments & Results

Data We experiment with five natural language understanding tasks: question classification (TREC-6) (Voorhees and Tice, 2000), sentiment analysis (IMDB, SST-2) (Maas et al., 2011; Socher et al., 2013) and topic classification (DBPEDIA, AGNEWS) (Zhang et al., 2015), including binary and multi-class labels and varying dataset sizes (Table 1). More details can be found in Appendix A.1.

²Some guidelines propose an even smaller number of warmup steps, such as 6% in RoBERTa (Liu et al., 2020).

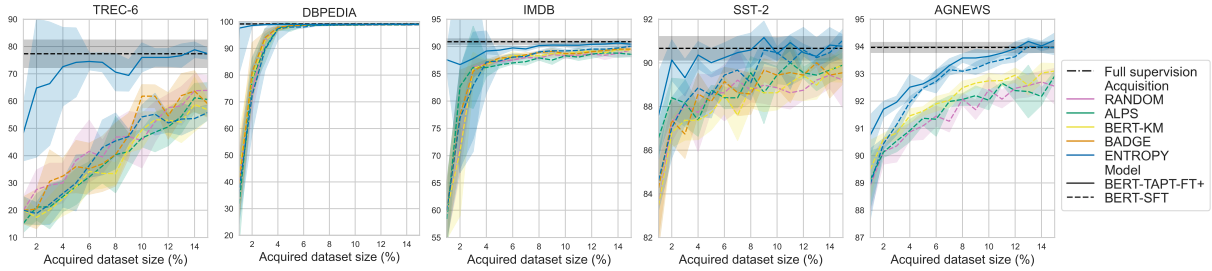


Figure 1: Test accuracy during AL iterations. We plot the median and standard deviation across five runs.

Experimental Setup We perform all AL experiments using BERT-base (Devlin et al., 2019) and ENTROPY, BERTKM, ALPS (Yuan et al., 2020), BADGE (Ash et al., 2020) and RANDOM (baseline) as the acquisition functions. We pair our proposed training approach TAPT-FT+ with ENTROPY acquisition. We describe in detail the AL setting and we provide results with more uncertainty-based acquisition functions in the Appendix.

Results Figure 1 shows the test accuracy during AL iterations. We first observe that our proposed approach (TAPT-FT+) achieves large data efficiency reaching the full-dataset performance within the budget for all datasets, in contrast to the standard AL approach (BERT-SFT). The effectiveness of our approach is mostly notable in the smaller datasets. In TREC-6, it achieves the goal accuracy with almost 10% annotated data, while in DBPEDIA only in the first iteration with 2% of the data. After the first AL iteration in IMDB, TAPT-FT+, it achieves only 2.5 points of accuracy lower than the performance when using 100% of the data. In the larger SST-2 and AGNEWS datasets, it is closer to the baselines but still outperforms them, achieving the full-dataset performance with 8% and 12% of the data respectively. We also observe that in all datasets, the addition of our proposed pretraining step (TAPT) and fine-tuning technique (FT+) leads to large performance gains, especially in the first AL iterations. This is particularly evident in TREC-6, DBPEDIA and IMDB datasets, where after the first AL iteration (i.e. equivalent to 2% of training data) TAPT+FT+ with ENTROPY is 45, 30 and 12 points in accuracy higher than the ENTROPY baseline with BERT and SFT.

Training vs. Acquisition Strategy We finally observe that the performance curves of the various acquisition functions considered (i.e. dotted lines) are generally close to each other, suggesting that the choice of the acquisition strategy may not

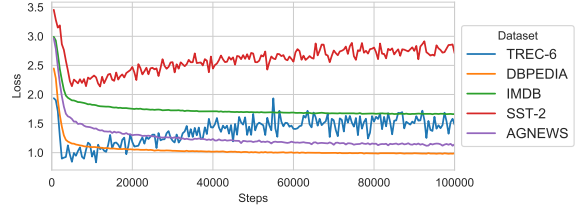


Figure 2: Validation MLM loss during TAPT.

affect substantially the AL performance in certain cases. In other words, we conclude that *the training strategy can be more important than the acquisition strategy*. We find that uncertainty sampling with ENTROPY is generally the best performing acquisition function, followed by BADGE. Still, finding a universally well-performing acquisition function, independent of the training strategy, is an open research question.

4 Analysis & Discussion

Task-Adaptive Pretraining We present details of TAPT (§2) and reflect on its effectiveness in the AL pipeline. Following Gururangan et al. (2020), we continue pretraining BERT for the MLM task using all the unlabeled data $\mathcal{D}_{\text{pool}}$ for all datasets separately. We plot the learning curves of BERT-TAPT for all datasets in Figure 2. We first observe that the masked LM loss is steadily decreasing for DBPEDIA, IMDB and AGNEWS across optimization steps, which correlates with the high early AL performance gains of TAPT in these datasets (Fig. 1). We also observe that the LM overfits in TREC-6 and SST-2 datasets. We attribute this to the very small training dataset of TREC-6 and the informal textual style of SST-2. Although SST-2 includes approximately 67K of training data, the sentences are very short (i.e. average length of 9.4 words per sentence). We hypothesize the LM overfits because of the lack of long and more diverse sentences. See Appendix B.1 for more details on TAPT.

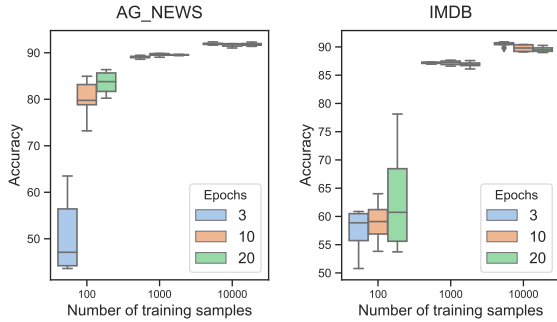


Figure 3: Few-shot standard BERT fine-tuning.

Few-shot Fine-tuning We highlight the importance of considering the few-shot learning problem in the early AL stages which is often neglected in literature. This is more important when using pretrained LMs, since they are overparameterized models that require adapting their training scheme in low data settings to ensure robustness. To illustrate the inefficiency of standard fine-tuning (SFT), we randomly undersample AGNEWS and IMDB to form low, medium and high data settings (i.e. 100, 1, 000 and 10, 000 training samples) and train BERT for a fixed number of 3, 10, and 20 epochs. Figure 3 shows that SFT is suboptimal for low data settings, indicating that more optimization steps are needed for the model to adapt to the few training samples (Zhang et al., 2020; Mosbach et al., 2021). As the training samples increase, fewer epochs are often better. It is thus evident that there is not an optimal way to choose a predefined number of epochs to train the model given the number of training examples. This motivates the need to find a fine-tuning policy for AL that effectively adapts to the data resource setting of each iteration, which is mainly tackled by our proposed fine-tuning approach FT+ (§2).

Ablation Study We also conduct an ablation study to show that our proposed pretraining step (TAPT) and fine-tuning method (FT+), provide large gains compared to standard BERT fine-tuning (SFT) in terms of accuracy, data efficiency and uncertainty calibration. We compare BERT with SFT, BERT with FT+ and BERT-TAPT with FT+. Along with test accuracy, we also evaluate each AL model using uncertainty estimation metrics (Ovadia et al., 2019): Brier score, negative log likelihood (NLL), expected calibration error (ECE) and entropy. A well-calibrated model should have high accuracy and low uncertainty. Figure 4 shows the results for the smallest and largest datasets, TREC-6 and

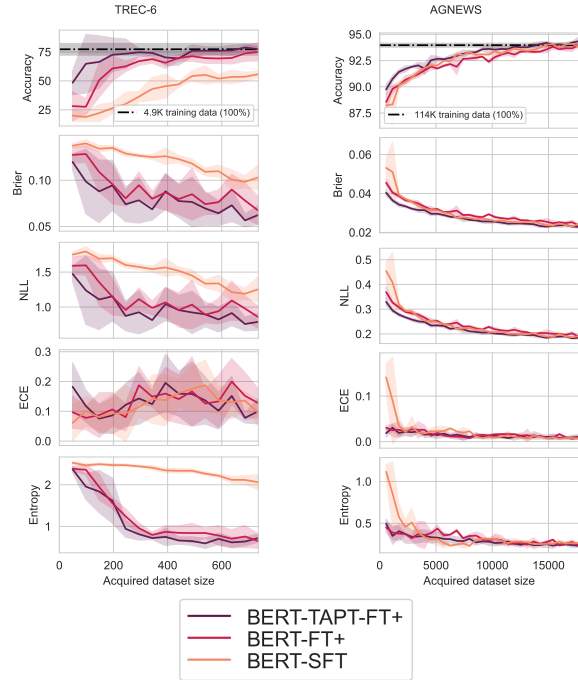


Figure 4: Ablation study for TAPT and FT+.

AGNEWS respectively. For TREC-6, training BERT with our fine-tuning approach FT+ provides large gains both in accuracy and uncertainty calibration, showing the importance of fine-tuning the LM for a larger number of epochs in low resource settings. For the larger dataset, AGNEWS, we see that BERT with SFT performs equally to FT+ which is the ideal scenario. We see that our fine-tuning approach does not deteriorate the performance of BERT given the large increase in warmup steps, showing that our simple strategy provides robust results in both high and low resource settings. After demonstrating that FT+ yields better results than SFT, we next compare BERT-TAPT-FT+ against BERT-FT+. We observe that in both datasets BERT-TAPT outperforms BERT, with this being particularly evident in the early iterations. This confirms our hypothesis that by implicitly using the entire pool of unlabeled data for extra pretraining (TAPT), we boost the performance of the AL model using less data.

5 Conclusion

We have presented a simple yet effective training scheme for AL with pretrained LMs, that yields substantially better results than standard fine-tuning. We also find that the proposed training strategy is more effective in improving performance than the selected acquisition function in certain cases, showing how critical it is to properly adapt a large pretrained LM to low data AL settings.

285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339

References

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*.

Zalán Bodó, Zsolt Minier, and Lehel Csató. 2011. [Active learning with clustering](#). In *Proceedings of the Active Learning and Experimental Design workshop In conjunction with AISTATS 2010*, volume 16, pages 127–139.

Klaus Brinker. 2003. [Incorporating diversity in active learning with support vector machines](#). In *Proceedings of the International Conference on Machine Learning*, pages 59–66.

David A. Cohn, Zoubin Ghahramani, and Michael I. Jordan. 1996. [Active learning with statistical models](#). *Journal of Artificial Intelligence Research*, 4(1):129–145.

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 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4171–4186.

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](#). *ArXiv*.

Melanie Ducoffe and Frederic Precioso. 2018. [Adversarial active learning for deep networks: a margin based approach](#).

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 Conference on Empirical Methods in Natural Language Processing*, pages 7949–7962.

Yarin Gal and Zoubin Ghahramani. 2016. [Dropout as a bayesian approximation: Representing model uncertainty in deep learning](#). In *Proceedings of the International Conference on Machine Learning*, volume 48, pages 1050–1059.

Yarin Gal, Riashat Islam, and Zoubin Ghahramani. 2017. [Deep Bayesian active learning with image data](#). In *Proceedings of the International Conference on Machine Learning*, volume 70, pages 1183–1192.

Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. 2020. [Don’t stop pretraining: Adapt language models to domains and tasks](#). In

Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 8342–8360, Online. Association for Computational Linguistics.

Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, and Máté Lengyel. 2011. [Bayesian active learning for classification and preference learning](#). *ArXiv*.

Jeremy Howard and Sebastian Ruder. 2018. [Universal language model fine-tuning for text classification](#). In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, pages 328–339.

Fariz Ikhwantri, Samuel Louvan, Kemal Kurniawan, Bagas Abisena, Valdi Rachman, Alfian Farizki Wicaksono, and Rahmad Mahendra. 2018. [Multi-task active learning for neural semantic role labeling on low resource conversational corpus](#). In *Proceedings of the Workshop on Deep Learning Approaches for Low-Resource NLP*, pages 43–50.

Siddharth Karamcheti, Ranjay Krishna, Li Fei-Fei, and Christopher Manning. 2021. [Mind your outliers! investigating the negative impact of outliers on active learning for visual question answering](#). 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 7265–7281, Online. Association for Computational Linguistics.

Jungo Kasai, Kun Qian, Sairam Gurajada, Yunyao Li, and Lucian Popa. 2019. [Low-resource deep entity resolution with transfer and active learning](#). In *Proceedings of the Conference of the Association for Computational Linguistics*, pages 5851–5861.

Andreas Kirsch, Joost van Amersfoort, and Yarin Gal. 2019. [BatchBALD: Efficient and diverse batch acquisition for deep bayesian active learning](#). In *Neural Information Processing Systems*, pages 7026–7037.

David D. Lewis and William A. Gale. 1994. [A sequential algorithm for training text classifiers](#). In *Proceedings of the Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Ro{bert}a: A robustly optimized {bert} pretraining approach](#).

Ilya Loshchilov and Frank Hutter. 2019. [Decoupled weight decay regularization](#). In *International Conference on Learning Representations*.

David Lowell and Zachary C Lipton. 2019. [Practical obstacles to deploying active learning](#). *Proceedings of the Conference on Empirical Methods in Natural Language Processing and the International Joint Conference on Natural Language Processing*, pages 21–30.

396	Andrew L. Maas, Raymond E. Daly, Peter T. Pham,	<i>Computational Linguistics: Main Volume</i> , pages	452
397	Dan Huang, Andrew Y. Ng, and Christopher Potts.	1698–1712, Online. Association for Computational	453
398	2011. Learning word vectors for sentiment analy-	Linguistics.	454
399	sis . In <i>Proceedings of the Annual Meeting of the</i>		
400	<i>Association for Computational Linguistics: Human</i>	Yanyao Shen, Hyokun Yun, Zachary Lipton, Yakov	455
401	<i>Language Technologies</i> , pages 142–150.	Kronrod, and Animashree Anandkumar. 2017.	456
		Deep active learning for named entity recognition .	457
402	Katerina Margatina, Giorgos Vernikos, Loïc Barrault,	In <i>Proceedings of the Workshop on Representation</i>	458
403	and Nikolaos Aletras. 2021. Active learning by ac-	<i>Learning for NLP</i> , pages 252–256.	459
404	quiring contrastive examples. In <i>Proceedings of the</i>		
405	<i>Conference on Empirical Methods in Natural Lan-</i>	Aditya Siddhant and Zachary C Lipton. 2018. Deep	460
406	<i>guage Processing</i> .	bayesian active learning for natural language pro-	461
		cessing: Results of a Large-Scale empirical study .	462
407	Marius Mosbach, Maksym Andriushchenko, and Diet-	In <i>Proceedings of the Conference on Empirical</i>	463
408	rich Klakow. 2021. On the stability of fine-tuning	<i>Methods in Natural Language Processing</i> , pages	464
409	{bert}: Misconceptions, explanations, and strong	2904–2909.	465
410	baselines . In <i>International Conference on Learning</i>		
411	<i>Representations</i> .	Richard Socher, Alex Perelygin, Jean Wu, Jason	466
		Chuang, Christopher D. Manning, Andrew Ng, and	467
412	Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado,	Christopher Potts. 2013. Recursive deep models	468
413	D. Sculley, Sebastian Nowozin, Joshua Dillon, Bal-	for semantic compositionality over a sentiment tree-	469
414	aji Lakshminarayanan, and Jasper Snoek. 2019. Can	bank . In <i>Proceedings of the Conference on Empiri-</i>	470
415	you trust your model's uncertainty? evaluating	<i>cal Methods in Natural Language Processing</i> , pages	471
416	predictive uncertainty under dataset shift . In <i>Ad-</i>	1631–1642.	472
417	<i>vances in Neural Information Processing Systems</i> ,		
418	volume 32, pages 13991–14002.	N Srivastava, G Hinton, A Krizhevsky, and others.	473
		2014. Dropout: a simple way to prevent neural net-	474
419	Adam Paszke, Sam Gross, Francisco Massa, Adam	works from overfitting . <i>Journal of Machine Learn-</i>	475
420	Lerer, James Bradbury, Gregory Chanan, Trevor	<i>ing Research</i> , 15(56):1929–1958.	476
421	Killeen, Zeming Lin, Natalia Gimelshein, Luca		
422	Antiga, Alban Desmaison, Andreas Kopf, Edward	Ellen Voorhees and Dawn Tice. 2000. The trec-8 ques-	477
423	Yang, Zachary DeVito, Martin Raison, Alykhan Te-	tion answering track evaluation . <i>Proceedings of the</i>	478
424	jani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang,	<i>Text Retrieval Conference</i> .	479
425	Junjie Bai, and Soumith Chintala. 2019. Pytorch:	Alex Wang, Amanpreet Singh, Julian Michael, Felix	480
426	An imperative style, high-performance deep learn-	Hill, Omer Levy, and Samuel R. Bowman. 2019.	481
427	ing library . In <i>Advances in Neural Information Pro-</i>	GLUE: A multi-task benchmark and analysis plat-	482
428	<i>cessing Systems</i> , pages 8024–8035.	form for natural language understanding . In <i>Inter-</i>	483
		<i>national Conference on Learning Representations</i> .	484
429	Ameya Prabhu, Charles Dognin, and Maneesh Singh.	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	485
430	2019. Sampling bias in deep active classification:	Chaumond, Clement Delangue, Anthony Moi, Pier-	486
431	An empirical study . In <i>Proceedings of the Confer-</i>	ric Cistac, Tim Rault, Remi Louf, Morgan Funtow-	487
432	<i>ence on Empirical Methods in Natural Language</i>	icz, Joe Davison, Sam Shleifer, Patrick von Platen,	488
433	<i>Processing and the International Joint Conference</i>	Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,	489
434	<i>on Natural Language Processing</i> , pages 4056–4066.	Teven Le Scao, Sylvain Gugger, Mariama Drame,	490
		Quentin Lhoest, and Alexander Rush. 2020. Trans-	491
435	Ozan Sener and Silvio Savarese. 2018. Active learn-	formers: State-of-the-art natural language process-	492
436	ing for convolutional neural networks: A core-set	ing . In <i>Proceedings of the Conference on Empirical</i>	493
437	approach . In <i>International Conference on Learning</i>	<i>Methods in Natural Language Processing: System</i>	494
438	<i>Representations</i> .	<i>Demonstrations</i> , pages 38–45.	495
439	Burr Settles. 2009. Active learning literature survey .	Michelle Yuan, Hsuan-Tien Lin, and Jordan Boyd-	496
440	Computer sciences technical report.	Graber. 2020. Cold-start active learning through	497
		self-supervised language modeling .	498
441	Claude Elwood Shannon. 1948. A mathematical the-	Tianyi Zhang, Felix Wu, Arzoo Katiyar, Kilian Q.	499
442	ory of communication . <i>The Bell System Technical</i>	Weinberger, and Yoav Artzi. 2020. Revisiting few-	500
443	<i>Journal</i> .	sample bert fine-tuning . <i>ArXiv</i> .	501
444	Artem Shelmanov, Dmitri Puzyrev, Lyubov	Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015.	502
445	Kupriyanova, Denis Belyakov, Daniil Larionov,	Character-level convolutional networks for text clas-	503
446	Nikita Khromov, Olga Kozlova, Ekaterina Arte-	sification . In <i>Advances in Neural Information Pro-</i>	504
447	mov, Dmitry V. Dylov, and Alexander Panchenko.	<i>cessing Systems</i> , volume 28, pages 649–657. Curran	505
448	2021. Active learning for sequence tagging with	Associates, Inc.	506
449	deep pre-trained models and Bayesian uncertainty		
450	estimates . In <i>Proceedings of the 16th Conference</i>		
451	<i>of the European Chapter of the Association for</i>		

A Appendix: Experimental Setup

A.1 Datasets

We experiment with five diverse natural language understanding tasks including binary and multi-class labels and varying dataset sizes (Table 1). The first task is question classification using the six-class version of the small TREC-6 dataset of open-domain, fact-based questions divided into broad semantic categories (Voorhees and Tice, 2000). We also evaluate our approach on sentiment analysis using the binary movie review IMDB dataset (Maas et al., 2011) and the binary version of the SST-2 dataset (Socher et al., 2013). We finally use the large-scale AGNEWS and DBPEDIA datasets from Zhang et al. (2015) for topic classification. We undersample the latter and form a $\mathcal{D}_{\text{pool}}$ of 20K examples and \mathcal{D}_{val} 2K as in Margatina et al. (2021). For TREC-6, IMDB and SST-2 we randomly sample 10% from the training set to serve as the validation set, while for AGNEWS we sample 5%. For the DBPEDIA dataset we undersample both training and validation datasets (from the standard splits) to facilitate our AL simulation (i.e. the original dataset consists of 560K training and 28K validation data examples). For all datasets we use the standard test set, apart from the SST-2 dataset that is taken from the GLUE benchmark (Wang et al., 2019) we use the development set as the held-out test set (and subsample a development set from the original training set).

A.2 Training & AL Details

We use BERT-BASE (Devlin et al., 2019) and fine-tune it (TAPT §2) for 100K steps, with learning rate $2e - 05$ and the rest of hyperparameters as in Gururangan et al. (2020) using the HuggingFace library (Wolf et al., 2020). We evaluate the model 5 times per epoch on \mathcal{D}_{val} and keep the one with the lowest validation loss as in Dodge et al. (2020). We use the code provided by Kirsch et al. (2019) for the uncertainty-based acquisition functions and Yuan et al. (2020) for ALPS, BADGE and BERTKM. We use the standard splits provided for all datasets, if available, otherwise we randomly sample a validation set. We test all models on a held-out test set. We repeat all experiments with five different random seeds resulting into different initializations of \mathcal{D}_{lab} and the weights of the extra task-specific output feedforward layer. For all datasets we use as budget the 15% of $\mathcal{D}_{\text{pool}}$. Each experiment is run on a single Nvidia Tesla V100 GPU.

A.3 Hyperparameters

For all datasets we train BERT-BASE (Devlin et al., 2019) from the HuggingFace library (Wolf et al., 2020) in Pytorch (Paszke et al., 2019). We train all models with batch size 16, learning rate $2e - 5$, no weight decay, AdamW optimizer with epsilon $1e - 8$. For all datasets we use maximum sequence length of 128, except for IMDB and AGNEWS that contain longer input texts, where we use 256. To ensure reproducibility and fair comparison between the various methods under evaluation, we run all experiments with the same five seeds that we randomly selected from the range $[1, 9999]$.

A.4 Baselines

Acquisition functions We compare ENTROPY with four baseline acquisition functions. The first is the standard AL baseline, **RANDOM**, which applies uniform sampling and selects k data points from $\mathcal{D}_{\text{pool}}$ at each iteration. The second is **BADGE** (Ash et al., 2020), an acquisition function that aims to combine diversity and uncertainty sampling. The algorithm computes *gradient embeddings* g_x for every candidate data point x in $\mathcal{D}_{\text{pool}}$ and then uses clustering to select a batch. Each g_x is computed as the gradient of the cross-entropy loss with respect to the parameters of the model’s last layer. We also compare against a recently introduced cold-start acquisition function called **ALPS** (Yuan et al., 2020). ALPS acquisition uses the masked language model (MLM) loss of BERT as a proxy for model uncertainty in the downstream classification task. Specifically, aiming to leverage both uncertainty and diversity, ALPS forms a *surprisal embedding* s_x for each x , by passing the unmasked input x through the BERT MLM head to compute the cross-entropy loss for a random 15% subsample of tokens against the target labels. ALPS clusters these embeddings to sample k sentences for each AL iteration. Last, following Yuan et al. (2020), we use **BERTKM** as a diversity baseline, where the l_2 normalized BERT output embeddings are used for clustering.

Models & Fine-tuning Methods We evaluate two variants of the pretrained language model; the original **BERT** model, used in Yuan et al. (2020) and Ein-Dor et al. (2020)³, and our adapted model **BERT-TAPT** (§2), and two fine-tuning methods;

³Ein-Dor et al. (2020) evaluate various acquisition functions, including entropy with MC dropout, and use BERT with the standard fine-tuning approach (SFT).

604
605

our proposed fine-tuning approach FT+ (§2) and standard BERT fine-tuning SFT.

MODEL	TREC-6	DBPEDIA	IMDB	SST-2	AGNEWS
VALIDATION SET					
BERT	94.4	99.1	90.7	93.7	94.4
BERT-TAPT	95.2	99.2	91.9	94.3	94.5
TEST SET					
BERT	80.6	99.2	91.0	90.6	94.0
BERT-TAPT	77.2	99.2	91.9	90.8	94.2

Table 2: Accuracy with 100% of data over five runs (different random seeds).

B Appendix: Analysis

B.1 Task-Adaptive Pretraining (TAPT) & Full-Dataset Performance

As discussed in §2 and §4, we continue training the BERT-BASE (Devlin et al., 2019) pretrained masked language model using the available data $\mathcal{D}_{\text{pool}}$. We explored various learning rates between $1e-4$ and $1e-5$ and found the latter to produce the lowest validation loss. We trained each model (one for each dataset) for up to 100K optimization steps, we evaluated on \mathcal{D}_{val} every 500 steps and saved the checkpoint with the lowest validation loss. We used the resulting model in our (BERT-TAPT) experiments. We plot the learning curves of masked language modeling task (TAPT) for three datasets and all considered learning rates in Figure 5. We notice that a smaller learning rate facilitates the training of the MLM.

In Table 2 we provide the validation and test accuracy of BERT and BERT-TAPT for all datasets. We present the mean across runs with three random seeds. For fine-tuning the models, we used the proposed approach FT+ (§2).

B.2 Performance of Acquisition Functions

In our BERT-TAPT-FT+ experiments so far, we showed results with ENTROPY. We have also experimented with various uncertainty-based acquisition functions. Specifically, four uncertainty-based acquisition functions are used in our work: LEAST CONFIDENCE, ENTROPY, BALD and BATCHBALD. LEAST CONFIDENCE (Lewis and Gale, 1994) sorts $\mathcal{D}_{\text{pool}}$ by the probability of *not* predicting the most confident class, in descending order, ENTROPY (Shannon, 1948) selects samples that maximize the predictive entropy, and BALD (Houlsby et al., 2011), short for Bayesian Active Learning by Disagreement, chooses data

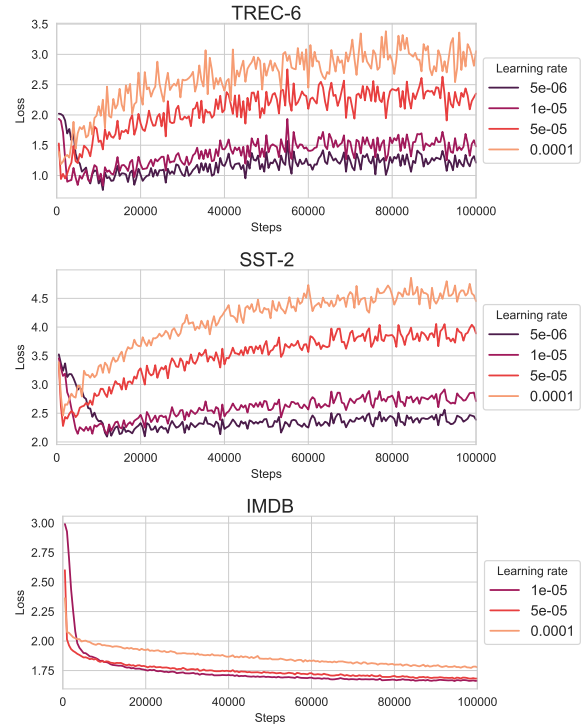


Figure 5: Learning curves of TAPT for various learning rates.

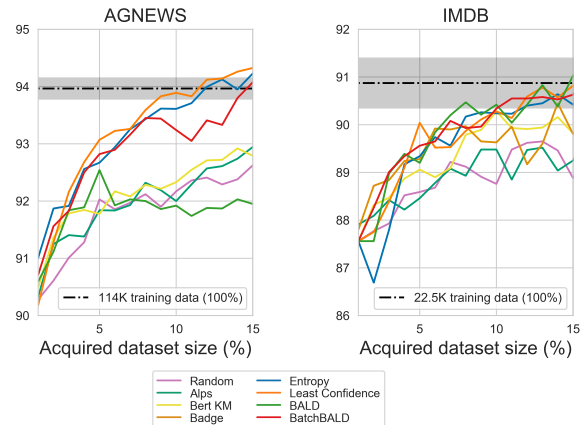


Figure 6: Comparison of acquisition functions using TAPT and FT+ in training BERT.

points that maximize the mutual information between predictions and model’s posterior probabilities. BATCHBALD (Kirsch et al., 2019) is a recently introduced extension of BALD that *jointly* scores points by estimating the mutual information between multiple data points and the model parameters. This iterative algorithm aims to find *batches* of informative data points, in contrast to BALD that chooses points that are informative individually. Note that LEAST CONFIDENCE, ENTROPY and BALD have been used in AL for NLP by Siddhant and Lipton (2018). To the best of our

643
644
645
646
647
648
649
650
651
652
653
654

	TREC-6	SST-2	IMDB	DBPEDIA	AGNEWS
RANDOM	0/0	0/0	0/0	0/0	0/0
ALPS	0/57	0/478	0/206	0/134	0/634
BADGE	0/63	0/23110	0/1059	0/192	-
BERTKM	0/47	0/2297	0/324	0/137	0/3651
ENTROPY	81/0	989/0	557/0	264/0	2911/0
LEAST CONFIDENCE	69/0	865/0	522/0	256/0	2607/0
BALD	69/0	797/0	524/0	256/0	2589/0
BATCHBALD	69/21	841/1141	450/104	256/482	2844/5611

Table 3: Runtimes (in seconds) for all datasets. In each cell of the table we present a tuple i/s where i is the *inference time* and s the *selection time*. *Inference time* is the time for the model to perform a forward pass for all the unlabeled data in $\mathcal{D}_{\text{pool}}$ and *selection time* is the time that each acquisition function requires to rank all candidate data points and select k for annotation (for a single iteration). Since we cannot report the runtimes for every model in the AL pipeline (at each iteration the size of $\mathcal{D}_{\text{pool}}$ changes), we provide the median.

655 knowledge, BATCHBALD is evaluated for the first
656 time in the NLP domain.

657 Instead of using the output softmax probabilities
658 for each class, we use a probabilistic formulation of
659 deep neural networks in order to acquire better cali-
660 brated scores. Monte Carlo (MC) dropout (Gal and
661 Ghahramani, 2016) is a simple yet effective method
662 for performing approximate variational inference,
663 based on dropout (Srivastava et al., 2014). Gal
664 and Ghahramani (2016) prove that by simply per-
665 forming *dropout during the forward pass in making*
666 *predictions*, the output is equivalent to the predic-
667 tion when the parameters are sampled from a varia-
668 tional distribution of the true posterior. Therefore,
669 dropout during inference results into obtaining pre-
670 dictions from different parts of the network. Our
671 BERT-based \mathcal{M}_i model uses dropout layers during
672 training for regularization. We apply MC dropout
673 by simply activating them during test time and we
674 perform multiple stochastic forward passes. For-
675 mally, we do N passes of every $x \in \mathcal{D}_{\text{pool}}$ through
676 $\mathcal{M}_i(x; W_i)$ to acquire N different output proba-
677 bility distributions for each x . MC dropout for
678 AL has been previously used in the literature (Gal
679 et al., 2017; Shen et al., 2017; Siddhant and Lip-
680 ton, 2018; Lowell and Lipton, 2019; Ein-Dor et al.,
681 2020; Shelmanov et al., 2021).

682 Our findings show that all functions provide sim-
683 ilar performance, except for BALD that slightly
684 underperforms. This makes our approach agnos-
685 tic to the selected uncertainty-based acquisition
686 method. We also evaluate our proposed methods
687 with our baseline acquisition functions, i.e. RAN-
688 DOM, ALPS, BERTKM and BADGE, since our
689 training strategy is orthogonal to the acquisition

690 strategy. We compare all acquisition functions with
691 BERT-TAPT-FT+ for AGNEWS and IMDB in Fig-
692 ure 6. We observe that in general uncertainty-based
693 acquisition performs better compared to diversity,
694 while all acquisition strategies have benefited from
695 our training strategy (TAPT and FT+).

696 B.3 Efficiency of Acquisition Functions

697 In this section we discuss the efficiency of the
698 eight acquisition functions considered in this work;
699 RANDOM, ALPS, BADGE, BERTKM, ENTROPY,
700 LEAST CONFIDENCE, BALD and BATCHBALD.

701 In Table 3 we provide the runtimes for all ac-
702 quisition functions and datasets. Each AL experi-
703 ments consists of multiple iterations and (therefore
704 multiple models), each with a different training
705 dataset \mathcal{D}_{lab} and pool of unlabeled data $\mathcal{D}_{\text{pool}}$. In
706 order to evaluate how computationally heavy is
707 each method, we provide the *median* of all the
708 models in one AL experiment. We calculate the
709 runtime of two types of functionalities. The first is
710 the *inference time* and stands for the forward pass
711 of each $x \in \mathcal{D}_{\text{pool}}$ to acquire confidence scores for
712 uncertainty sampling. RANDOM, ALPS, BADGE
713 and BERTKM do not require this step so it is only
714 applied of uncertainty-based acquisition where ac-
715 quiring uncertainty estimates with MC dropout is
716 needed. The second functionality is *selection time*
717 and measures how much time each acquisition func-
718 tion requires to rank and select the k data points
719 from $\mathcal{D}_{\text{pool}}$ to be labeled in the next step of the AL
720 pipeline. RANDOM, ENTROPY, LEAST CONFID-
721 DENCE and BALD perform simple equations to
722 rank the data points and therefore so do not require
723 selection time. On the other hand, ALPS, BADGE,

724 BERTKM and BATCHBALD perform iterative al-
725 gorithms that increase selection time. From all ac-
726 quisition functions ALPS and BERTKM are faster
727 because they do not require the inference step of
728 all the unlabeled data to the model. ENTROPY,
729 LEAST CONFIDENCE and BALD require the same
730 time for selecting data, which is equivalent for the
731 time needed to perform one forward pass of the en-
732 tire $\mathcal{D}_{\text{pool}}$. Finally BADGE and BATCHBALD are
733 the most computationally heavy approaches, since
734 both algorithms require multiple computations for
735 the *selection time*. RANDOM has a total runtime of
736 zero seconds, as expected.