

# Evaluating the fairness of task-adaptive pretraining on unlabeled test data before few-shot text classification

Anonymous ACL submission

## Abstract

Few-shot learning benchmarks are critical for evaluating modern NLP techniques. But it's possible that benchmarks favor methods which easily make use of unlabeled text, because researchers can pretrain their models on unlabeled text from the test set. Given the dearth of research on this potential problem, we run experiments to quantify the bias caused by pretraining on unlabeled test set text instead of on independently drawn unlabeled text. A controlled experiment varying the numbers of training and test observations for 25 classification tasks and 2 language models—BERT and GPT-2—does not find evidence of bias. Furthermore, we demonstrate the importance of repeated subsampling when studying few-shot text classification, and recommend that few-shot learning benchmarks include multiple training folds. Code and data are available here: <https://github.com> (currently omitted for anonymity).

## 1 Introduction

For NLP benchmarks, it's standard to release text from the test set. This allows researchers to submit a file of predictions instead of submitting code. A potential concern is that researchers can technically use this text during training. Consider the Real-world Annotated Few-shot Tasks (RAFT) benchmark (Alex et al., 2021), which contains "few-shot" text classification tasks—tasks where the training set contains a relatively small number of labeled examples. Below is an excerpt from the RAFT paper (emphasis added):

For each task, we release a public training set with 50 examples and a larger unlabeled test set. *We encourage unsupervised pre-training on the unlabelled examples* and open-domain information retrieval.

In the RAFT competition, a model is evaluated by scoring its predictions on the same set of unlabeled text which the model may have been trained on (using an unsupervised training procedure).

It's wrong to train a model on test set features with their labels and then evaluate on the test set when one needs to estimate performance on out-of-sample data. Test set performance would be overoptimistic (Hastie et al., 2009). This fact is widely known. But what if, as encouraged by Alex et al. (2021), a model is trained on test set features *without* test set labels? This paper studies this question for the domain of few-shot text classification.

## 2 Motivation

NLP benchmarks for few-shot learning are widespread, as having only a handful of labeled examples is more common in practice. One consideration when designing these benchmarks is that some few-shot approaches can—at least theoretically—use unlabeled text from the test set. With Pattern-Exploiting Training (Schick and Schütze, 2021), for example, one can train the final classifier on test set text with soft labels predicted by an ensemble of supervised models. Or, with Pre-trained Prompt Tuning (Gu et al., 2022), one can pretrain the language model (LM) on unlabeled test set text before prompt-tuning on the labeled training set. A more classical approach would be to train a word2vec model (Mikolov et al., 2013) on unlabeled test set text, run this model on training text to get embeddings, and finally train a classifier on these embeddings with labels from the training set.

For other few-shot approaches, such as SetFit (Tunstall et al., 2022) and in-context learning with large LMs (as popularized by Brown et al., 2020), it's more common to only use labeled examples.

While the ability to exploit unlabeled text is useful, applying this ability to test set text could be substantively different than applying it to text which is

statistically independent of the test set. This difference in methodology may be more concerning in the few-shot setting than in the many-shot setting. It’s conceivable that differences between few-shot methods are as attributable to differences in how unlabeled text is used than how the few, labeled examples are used. This raises the question: can few-shot text classification benchmarks favor methods which exploited unlabeled text from the test set?

### 3 Related work

As indicated by the quote in §1, the RAFT benchmark implicitly assumes that the answer is no. It is not a fringe opinion that test set features may be used. The popular textbook by [Hastie et al. \(2009\)](#) contains the following passage without a reference or evidence (emphasis added):

There is one qualification: *initial unsupervised screening steps can be done before samples are left out*. For example, we could select the 1000 predictors with highest variance across all 50 samples, before starting cross-validation. *Since this filtering does not involve the class labels, it does not give the predictors an unfair advantage*.

The opposite opinion—that exploiting unlabeled test set features is unfair—may be more popular. For example, [Gururangan et al. \(2020\)](#) contains the following criticism of another study when comparing performances on a popular text classification benchmark:

[Thongtan and Phienthrakul \(2019\)](#) report a higher number (97.42) on IMDB, but they train their word vectors on the test set.

[Moscovich and Rosset \(2022\)](#) contains experiments and theory for unsupervised methods which are common to tasks involving tabular data. They find that estimators of out-of-sample performance which were subject to these methods may be biased positively or negatively, depending on all of the parameters of the problem. They recommend further research on this bias in more domains, particularly when dealing with small sample sizes and high-dimensional data.

## 4 Experimental design

In the absence of theory or experiments in NLP, this paper studies how much pretraining on unlabeled test set text biases test set performance for 25 diverse text classification tasks and two types of LMs: BERT ([Devlin et al., 2019](#)), and GPT-2 ([Radford et al., 2019](#)). Descriptions of the 25 classification tasks are included in Appendix A. The number of classes in each task ranges from 2 to 18.

At a high level, the goal of the experiment is to first establish that pretraining is beneficial, in line with [Gururangan et al. \(2020\)](#). Second, given that pretraining has a detectable benefit, the experiment measures the accuracy difference between using test set text for the pretraining stage—an arguably unfair methodology—versus using text which is independent of the test set—an inarguably fair methodology.

In more detail, the experiment starts by drawing three subsamples (without replacement) from the full sample of data for a given text classification task:

- extra:  $n$  (either 50, 100, 200 or 500) unlabeled texts which are optionally used for pretraining
- train:  $m$  (either 50 or 100) labeled texts for supervised classification training
- test:  $n$  labeled texts to report accuracy.

Next, three accuracy estimators are computed. The procedures used to obtain them are described below.

### 4.1 $\text{acc}_{\text{extra}}$

1. Train a freshly loaded, pretrained LM on the  $n$  unlabeled texts in extra using the LM’s pretraining objective—masked language modeling loss for BERT, or autoregressive/causal language modeling loss for GPT-2.
2. Add a linear layer to this further-pretrained model. For BERT, the linear layer transforms the [CLS] token embedding. For GPT-2, the linear layer transforms the last token’s embedding. The output dimension of the linear layer is the number of classes in the classification task. This layer, along with the rest of the weights in the LM, are finetuned to minimize classification cross entropy loss on train.
3. Compute the classification accuracy of this model on test.

Step 1 is task-adaptive pretraining—a procedure broadly recommended by Gururangan et al. (2020). Step 2 is the canonical way of training a transformer-based LM for a classification task, according to Section 2 of Zhang et al. (2021).

$\text{acc}_{\text{extra}}$  is clearly an unbiased estimator of out-of-sample accuracy, because it never trains on data from test.

## 4.2 $\text{acc}_{\text{test}}$

$\text{acc}_{\text{test}}$  is identical to  $\text{acc}_{\text{extra}}$ , except that pretraining is done on test instead of extra in step 1.

$\text{acc}_{\text{test}}$  represents what one might see in a competition like RAFT, where pretraining on unlabeled text from test is encouraged. It’s unclear whether this accuracy estimator is unbiased, because it was (pre)trained and evaluated on the same set of test set text. A reasonable hypothesis is that it’s overoptimistic, i.e.,  $E[\text{acc}_{\text{test}}] > E[\text{acc}_{\text{extra}}] =$  out-of-sample accuracy.

## 4.3 $\text{acc}_{\text{base}}$

$\text{acc}_{\text{base}}$  doesn’t do pretraining; it doesn’t make any use of unlabeled text. It simply trains a pretrained LM on train to do classification, and then computes this model’s accuracy on test.

This score is a control. If there’s no boost going from  $\text{acc}_{\text{base}}$  to  $\text{acc}_{\text{extra}}$ , then it shouldn’t be surprising that there’s no boost going from  $\text{acc}_{\text{extra}}$  to  $\text{acc}_{\text{test}}$ .

## 4.4 Repeated subsampling

The three accuracy estimators are paired, because their classification training and test sets are identical. The only difference is the source of unlabeled text for pretraining. For  $\text{acc}_{\text{extra}}$ , the source is independent of test set text. For  $\text{acc}_{\text{test}}$ , the source is exactly the test set text. For  $\text{acc}_{\text{base}}$ , no unlabeled text is used.

A potentially important source of variation in this experiment is the particular subsamples, i.e., the particular realizations of extra, train, and test for a given classification task. To expose this variation, the experiment procedure is repeated tens of times for each classification task.<sup>1</sup> For example, for  $n = 200$ , and for each of the 25 classification tasks, 50 ( $\text{acc}_{\text{extra}}$ ,  $\text{acc}_{\text{test}}$ ,  $\text{acc}_{\text{base}}$ ) triples are computed.

<sup>1</sup>For  $n = 50$  and  $n = 100$ , the experiment is repeated 100 times. For  $n = 200$ , the experiment is repeated 50 times. For  $n = 500$ , the experiment is repeated 20 times. In total, 81,000 finetuned BERT and GPT-2 models were evaluated in this experiment.

Appendix B explains more experiment choices.

## 5 Results

Figure 1 visualizes the distributions of  $\text{acc}_{\text{extra}} - \text{acc}_{\text{base}}$  and  $\text{acc}_{\text{test}} - \text{acc}_{\text{extra}}$  for  $m = 50, n = 200$ . Appendix D.2 contains visualizations for all other  $m, n$  settings.  $\text{acc}_{\text{extra}} - \text{acc}_{\text{base}}$  is a control: it’s the accuracy boost from pretraining on unlabeled independent text versus not pretraining at all.  $\text{acc}_{\text{test}} - \text{acc}_{\text{extra}}$  is the main quantity of interest: it’s the accuracy boost from pretraining on unlabeled test set text instead of on unlabeled independent text, i.e., it’s the evaluation bias.

Table 1 contains means of these differences for each configuration of the experiment. It roughly suggests that while pretraining is consistently beneficial, pretraining on unlabeled test set text does not bias test set performance one way or the other.

A more complete analysis of this data is motivated and performed in the next section.

	BERT	GPT-2
$n = 50$	4.1% 0.19%	3.8% 0.18%
$n = 100$	3.9% 0.18%	4.1% 0.11%
$n = 200$	3.9% -0.39%	4.4% -0.05%
$n = 500$	3.5% 0.48%	4.6% -0.08%

(a)  $m = 50$

	BERT	GPT-2
$n = 50$	6.2% -0.08%	2.2% -0.05%
$n = 100$	6.1% -0.37%	2.5% 0.03%
$n = 200$	4.1% 0.33%	6.3% -0.01%
$n = 500$	6.1% -0.16%	3.9% -0.21%

(b)  $m = 100$

Table 1: Sample means of accuracy differences taken across all subsamples of the 25 text classification tasks. For each cell, the upper-left of the diagonal corresponds to the sample mean of  $\text{acc}_{\text{extra}} - \text{acc}_{\text{base}}$ , and the lower-right corresponds to the sample mean of  $\text{acc}_{\text{test}} - \text{acc}_{\text{extra}}$ .

## 6 Analysis

Reporting means is not enough, especially when studying few-shot learning. Figure 1 (and Appendix D.2) demonstrates that there’s considerable

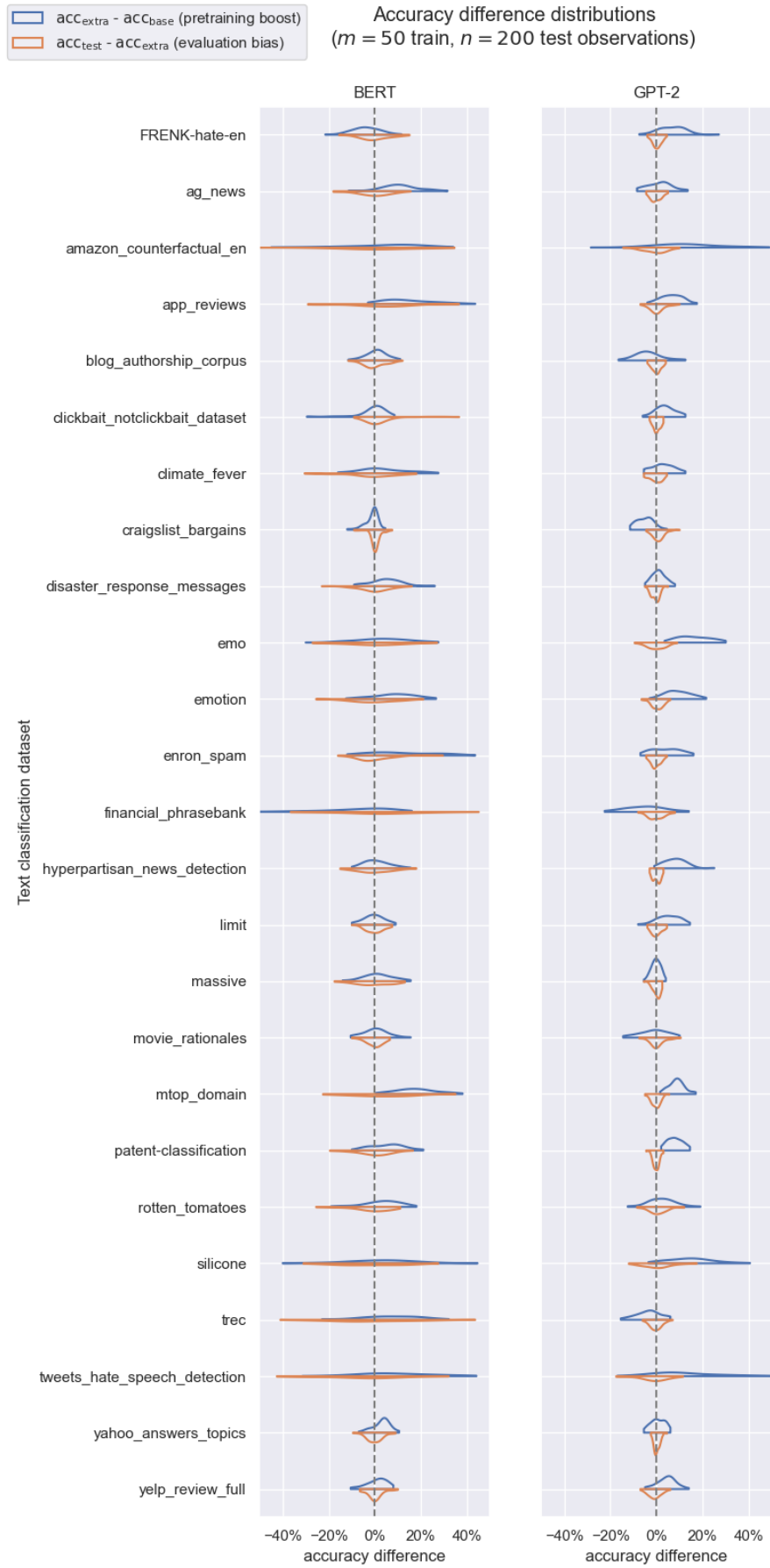


Figure 1

variance, despite pairing the accuracy estimators.<sup>2</sup> And while these visualizations tell us about how raw accuracy differences vary, they do not tell us how the mean accuracy difference varies. We seek a neat answer to the core questions: on this benchmark of 25 classification tasks, how much does the average accuracy differ between two modeling techniques, and how much does this average difference vary?

One way to communicate the variance is to estimate the standard error of the mean difference across classification tasks. But the standard error statistic can be difficult to interpret (Morey et al., 2016). Furthermore, its computation is not completely trivial due to the data’s hierarchical dependency structure: each triple, ( $\text{acc}_{\text{extra}}$ ,  $\text{acc}_{\text{test}}$ ,  $\text{acc}_{\text{base}}$ ), is drawn from (train, test), which is itself drawn from the given classification dataset.

## 6.1 Model

This analysis does not aim to estimate standard errors. Instead, a hierarchical model is fit:

$$Y_{ijkl} \sim \text{Binomial}(n, \lambda_{ijkl}) \quad (1)$$

$$\text{logit}(\lambda_{ijkl}) = \mu + \alpha z_i + U_j + V_{jk} + \beta x_{ijkl} \quad (2)$$

$$\mu \sim \text{Normal}(0, 1) \quad (3)$$

$$\alpha \sim \text{Normal}(0, 5) \quad (4)$$

$$U_j \sim \text{Normal}(0, \sigma_U) \quad (5)$$

$$V_{jk} \sim \text{Normal}(0, \sigma_V) \quad (6)$$

$$\beta \sim \text{Normal}(0, 1) \quad (7)$$

$$\sigma_U, \sigma_V \sim \text{HalfNormal}(0, 1) \quad (8)$$

- (1) number of correct predictions
- (2) logit link for accuracy rate, additive effects
- (3) prior for the global intercept
- (4) prior for the effect of the type of LM (BERT or GPT-2)—a control variable
- (5) prior for the effect of the classification task (partial-pooled to reduce overfitting)
- (6) prior for the nested effect of the task’s subsampled dataset
- (7) prior for the effect of interest ( $x_{ijkl} = 1$  indicates the modeling intervention)
- (8) prior for standard deviations.

<sup>2</sup>One source of variance is intentionally introduced: the subsamples/splits, as explained in §4.4. The other source of variance is inherent: the added linear layer to perform classification is initialized with random weights.

The model is fit using Markov Chain Monte Carlo, using the interface provided by the bambi package (Capretto et al., 2022). 4,000 samples from the posterior were drawn for each effect. Appendix E.1 includes a simulation that demonstrates the model’s ability to correctly recover null and non-null effects.

## 6.2 Posterior predictions

In NLP benchmarks, methods are assessed by taking their average performance across tasks. To place the analysis results in this context, samples from the posterior predictive distribution of  $Y_{ijk1} - Y_{ijk0}$  (6.1) are taken, then averaged across  $i$  (the 2 LM types—BERT and GPT-2),  $j$  (the 25 classification tasks), and  $k$  (their subsamples), and divided by  $n$  to obtain the distribution of the average accuracy difference:

$$\frac{\bar{Y}_{\dots 1} - \bar{Y}_{\dots 0}}{n}.$$

These distributions are plotted in Figure 2. Each distribution is that of the marginal effect of the modeling intervention: pretraining versus not pretraining before classification training (the pretraining boost), or pretraining on unlabeled test set text instead of on unlabeled independent text before classification training (the evaluation bias).

## 7 Discussion

Figure 2 demonstrates that the average pretraining boost is significant in every configuration of the experiment, ranging from 2% to 6%. This finding replicates that from Gururangan et al. (2020). After averaging across settings for  $m$ ,  $n$ , and the 2 LM types, only two of the 25 classification tasks had a pretraining boost less than 0, and both were greater than -1%.<sup>3</sup> Overall, pretraining is beneficial, so there may be a detectable evaluation bias.

As shown in Figure 2, the evaluation bias bounces inconsistently and insignificantly around 0. After averaging, 12 of the 25 classification tasks had a positive evaluation bias, and all tasks had an average evaluation bias less than 1% in absolute value. Given the lack of evidence for an evaluation bias in either direction, it’s unlikely that a benchmark which releases unlabeled test set text can systematically promote models pretrained on it

<sup>3</sup>The tasks were `blog_authorship_corpus` and `movie_rationales`.



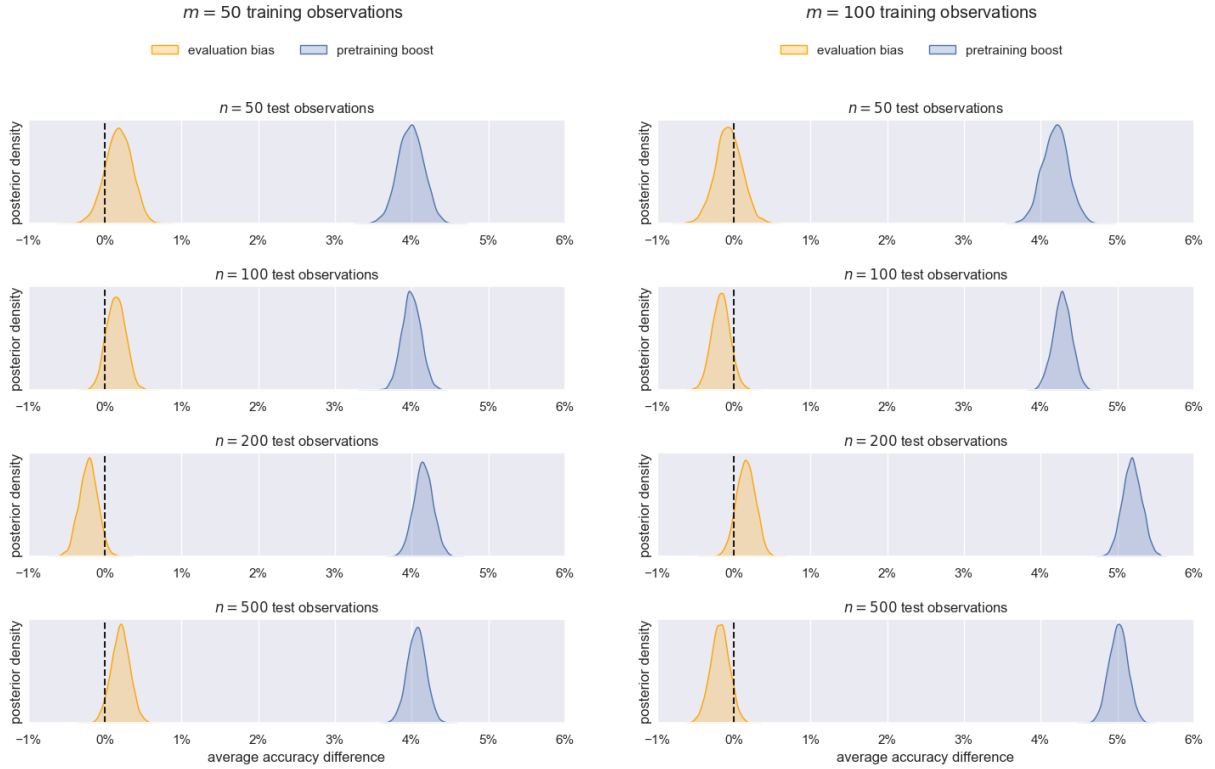


Figure 2: Distributions of average accuracy differences for  $m = 50$  (left) and  $m = 100$  (right). The evaluation bias is akin to  $\text{acc}_{\text{test}} - \text{acc}_{\text{extra}}$ . The pretraining boost is akin to  $\text{acc}_{\text{extra}} - \text{acc}_{\text{base}}$ .

over equally performant models which pretrained on unlabeled independent text.

Moscovich and Rosset (2022) found that the evaluation bias caused by certain unsupervised methods for tabular data gets closer to 0 as  $n$  increases. This finding is not confirmed by this experiment. Figure 2 shows that for  $m = 50$  and  $m = 100$ , the distribution of the evaluation bias consistently hovers around 0 across settings for  $n$ . But far more experiments varying  $n$  are needed to thoroughly assess this insensitivity.

## 8 Meta-analysis

§4.4 briefly argues for subsampling multiple datasets from the full classification dataset. To assess this argument, the analysis was repeated on 500 random slices of the  $m = 100, n = 500$  dataset of accuracies such that exactly 1 ( $\text{acc}_{\text{extra}}, \text{acc}_{\text{test}}, \text{acc}_{\text{base}}$ ) triple per classification task (instead of 20) is included. This unreplicated data is often all you get from benchmarks.

Figure 3 (right) displays the cumulative distribution function of the posterior mean of the evaluation bias for  $m = 100, n = 500$  under this unreplicated experimental design. The distribution is quite variant. There’s a 47% chance that the posterior mean

of  $\beta$ —the average increase in the log-odds of a correct prediction by pretraining on unlabeled test set text instead of on unlabeled independent text—is outside the interval  $(-0.04, 0.04)$ , which would indicate a significant negative or positive bias.<sup>4</sup> In other words, without subsampling, one may as well flip a coin to determine whether pretraining on unlabeled test set text is fair.

## 9 Conclusion

Across combinations for the number of classification training examples ( $m = 50, 100$ ) and the number of pretraining or evaluation examples ( $n = 50, 100, 200, 500$ ), pretraining on unlabeled test set text did not result in a consistent or significant bias compared to pretraining on unlabeled independent text. This is despite the almost universal benefit of pretraining.

One recommendation for designing few-shot benchmarks, which expands on the principle about robustness from Bragg et al. (2021) and recommendations from Madaan et al. (2024), is based on

<sup>4</sup>For 0.04, the odds ratio is  $e^{0.04} \approx 1.04$ . For context, the average odds ratio between adjacent submissions in the RAFT leaderboard is 1.03. For posterior means outside  $(-0.04, 0.04)$ , all of their 89% credible intervals exclude 0, which evidences a non-null effect.

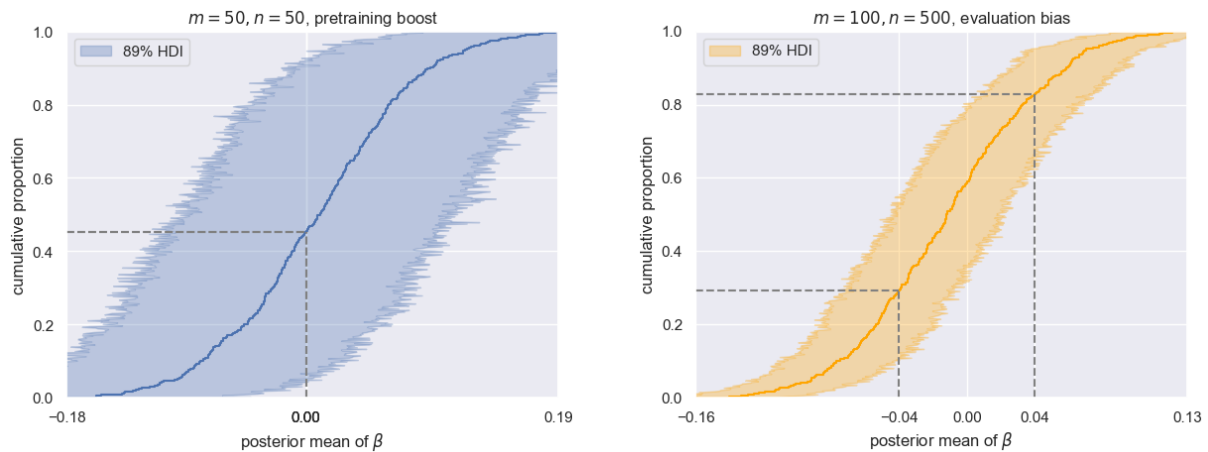


Figure 3: Distributions of this paper’s conclusions for  $m = 50, n = 50$  (left) and  $m = 100, n = 500$  (right) had there been no technical replication. (left)  $\beta$  is the average increase in the log-odds of a correct prediction by pretraining on unlabeled independent text versus not pretraining at all before classification training. (right)  $\beta$  is the average increase in the log-odds of a correct prediction by pretraining on unlabeled text from the test set instead of on unlabeled independent text before classification training.

the meta-analysis in §8: empirical studies of few-shot learning should consider including multiple, independent subsamples of training data. While a single training set combined with a large test set is sufficient for precise, unbiased estimation of out-of-sample performance, this estimator is conditional on the training set. In few-shot learning, the training set is, by definition, minimal. The estimator hides two sources of variance—that from the randomly drawn training set, and that from randomness inherent in the training procedure. Figure 3 shows that this variance is large-enough to turn a methodology into a coin flip for a standard pretraining-and-training procedure. In-context learning with large LMs is also sensitive to the few-shot examples used (Lu et al., 2022, Alzahrani et al., 2024). Benchmarks which require training on multiple, independent subsamples would expose training variance.

An important limitation of this paper is that it does not analyze semi-supervised methods like Pattern-Exploiting Training. This paper also doesn’t study somewhat nefarious uses of the test set such as hand-inspecting the text and targeting interventions accordingly. This paper’s conclusions are limited to task-adaptive pretraining of LMs.

A direction for future research is to further vary the amount of labeled training examples. Perhaps there’s overoptimism for minimal training sets. Another empirical direction is to repeat the experiment for larger LMs trained via supervised finetuning, assuming data contamination can be accounted for.

A theoretical direction is to explore the role of causality. Jin et al. (2021) argue and demonstrate that the benefit of task-adaptive pretraining depends on the learning task’s causal direction. Perhaps the principle of independent causal mechanisms is also relevant in assessing the fairness of pretraining on test set features.

## Acknowledgements

Currently omitted for anonymity.

## References

- Neel Alex, Eli Lifland, Lewis Tunstall, Abhishek Thakur, Pegah Maham, C. Riedel, Emmie Hine, Carolyn Ashurst, Paul Sedille, Alexis Carlier, Michael Noetel, and Andreas Stuhlmüller. 2021. [Raft: A real-world few-shot text classification benchmark](#). In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, volume 1.
- Norah Alzahrani, Hisham Abdullah Alyahya, Yazeed Alnumay, Sultan Alrashed, Shaykhah Alsubaie, Yusef Almushaykeh, Faisal Mirza, Nouf Alotaibi, Nora Altwairesh, Areeb Alowisheq, et al. 2024. When benchmarks are targets: Revealing the sensitivity of large language model leaderboards. *arXiv preprint arXiv:2402.01781*.
- Yoav Benjamini and Yosef Hochberg. 1995. [Controlling the false discovery rate: a practical and powerful approach to multiple testing](#). *Journal of the Royal statistical society: series B (Methodological)*, 57(1):289–300.

435	Jonathan Bragg, Arman Cohan, Kyle Lo, and Iz Beltagy.	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	492
436	2021. <a href="#">Flex: Unifying evaluation for few-shot nlp.</a>	pages 4277–4302, Toronto, Canada. Association for	493
437	<i>Advances in Neural Information Processing Systems</i> ,	Computational Linguistics.	494
438	34:15787–15800.		
439	Tom Brown, Benjamin Mann, Nick Ryder, Melanie	Giovanni Grano, Andrea Di Sorbo, Francesco Mercaldo,	495
440	Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind	Corrado A Visaggio, Gerardo Canfora, and Sebas-	496
441	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	tiano Panichella. 2017. <a href="#">Android apps and user feed-</a>	497
442	Askell, et al. 2020. <a href="#">Language models are few-shot</a>	<a href="#">back: a dataset for software evolution and quality</a>	498
443	<a href="#">learners</a> . <i>Advances in neural information processing</i>	<a href="#">improvement</a> . In <i>Proceedings of the 2nd ACM SIG-</i>	499
444	<i>systems</i> , 33:1877–1901.	<i>SOFT international workshop on app market analyt-</i>	500
445	Tomás Capretto, Camen Pihó, Ravin Kumar, Jacob	<i>ics</i> , pages 8–11.	501
446	Westfall, Tal Yarkoni, and Osvaldo A Martin. 2022.		
447	<a href="#">Bambi: A simple interface for fitting bayesian linear</a>	Yuxian Gu, Xu Han, Zhiyuan Liu, and Minlie Huang.	502
448	<a href="#">models in python</a> . <i>Journal of Statistical Software</i> ,	2022. <a href="#">PPT: Pre-trained prompt tuning for few-shot</a>	503
449	103(15):1–29.	<a href="#">learning</a> . In <i>Proceedings of the 60th Annual Meet-</i>	504
450	Emile Chapuis, Pierre Colombo, Matteo Manica,	<i>ing of the Association for Computational Linguistics</i>	505
451	Matthieu Labeau, and Chloé Clavel. 2020. <a href="#">Hier-</a>	( <i>Volume 1: Long Papers</i> ), pages 8410–8423, Dublin,	506
452	<a href="#">archical pre-training for sequence labelling in spoken</a>	Ireland. Association for Computational Linguistics.	507
453	<a href="#">dialog</a> . In <i>Findings of the Association for Computa-</i>		
454	<i>tional Linguistics: EMNLP 2020</i> , pages 2636–2648,	Suchin Gururangan, Ana Marasović, Swabha	508
455	Online. Association for Computational Linguistics.	Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey,	509
456	Ankush Chatterjee, Kedhar Nath Narahari, Meghana	and Noah A. Smith. 2020. <a href="#">Don’t stop pretraining:</a>	510
457	Joshi, and Puneet Agrawal. 2019. <a href="#">SemEval-2019 task</a>	<a href="#">Adapt language models to domains and tasks</a> . In	511
458	<a href="#">3: EmoContext contextual emotion detection in text</a> .	<i>Proceedings of the 58th Annual Meeting of the</i>	512
459	In <i>Proceedings of the 13th International Workshop</i>	<i>Association for Computational Linguistics</i> , pages	513
460	<i>on Semantic Evaluation</i> , pages 39–48, Minneapo-	8342–8360, Online. Association for Computational	514
461	lis, Minnesota, USA. Association for Computational	Linguistics.	515
462	Linguistics.		
463	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	Trevor Hastie, Robert Tibshirani, Jerome H Friedman,	516
464	Kristina Toutanova. 2019. <a href="#">BERT: Pre-training of</a>	and Jerome H Friedman. 2009. <a href="#">The elements of statis-</a>	517
465	<a href="#">deep bidirectional transformers for language under-</a>	<a href="#">tical learning: data mining, inference, and prediction</a> ,	518
466	<a href="#">standing</a> . In <i>Proceedings of the 2019 Conference of</i>	volume 2. Springer.	519
467	<i>the North American Chapter of the Association for</i>		
468	<i>Computational Linguistics: Human Language Tech-</i>	He He, Derek Chen, Anusha Balakrishnan, and Percy	520
469	<i>nologies, Volume 1 (Long and Short Papers)</i> , pages	Liang. 2018. <a href="#">Decoupling strategy and generation in</a>	521
470	4171–4186, Minneapolis, Minnesota. Association for	<a href="#">negotiation dialogues</a> . In <i>Proceedings of the 2018</i>	522
471	Computational Linguistics.	<i>Conference on Empirical Methods in Natural Lan-</i>	523
472	Jay DeYoung, Sarthak Jain, Nazneen Fatema Rajani,	<i>guage Processing</i> , pages 2333–2343, Brussels, Bel-	524
473	Eric Lehman, Caiming Xiong, Richard Socher, and	gium. Association for Computational Linguistics.	525
474	Byron C. Wallace. 2020. <a href="#">ERASER: A benchmark to</a>		
475	<a href="#">evaluate rationalized NLP models</a> . In <i>Proceedings</i>	Zhang Huangzhao. 2018. Yahoo-	526
476	<i>of the 58th Annual Meeting of the Association for</i>	<a href="#">answers-topic-classification-dataset</a> .	527
477	<i>Computational Linguistics</i> , pages 4443–4458, Online.	<a href="https://github.com/LC-John/">https://github.com/LC-John/</a>	528
478	Association for Computational Linguistics.	<a href="#">Yahoo-Answers-Topic-Classification-Dataset</a> .	529
479	Thomas Diggelmann, Jordan Boyd-Graber, Jannis Bu-	Zhijing Jin, Julius von Kügelgen, Jingwei Ni, Tejas	530
480	lian, Massimiliano Ciaramita, and Markus Leippold.	Vaidhya, Ayush Kaushal, Mrinmaya Sachan, and	531
481	2020. <a href="#">Climate-fever: A dataset for verification of</a>	Bernhard Schoelkopf. 2021. <a href="#">Causal direction of data</a>	532
482	<a href="#">real-world climate claims</a> .	<a href="#">collection matters: Implications of causal and an-</a>	533
483	Jack FitzGerald, Christopher Hench, Charith Peris,	<a href="#">ticausal learning for NLP</a> . In <i>Proceedings of the</i>	534
484	Scott Mackie, Kay Rottmann, Ana Sanchez, Aaron	<i>2021 Conference on Empirical Methods in Natural</i>	535
485	Nash, Liam Urbach, Vishesh Kakarala, Richa Singh,	<i>Language Processing</i> , pages 9499–9513, Online and	536
486	Swetha Ranganath, Laurie Crist, Misha Britan,	Punta Cana, Dominican Republic. Association for	537
487	Wouter Leeuwis, Gokhan Tur, and Prem Natara-	Computational Linguistics.	538
488	jan. 2023. <a href="#">MASSIVE: A 1M-example multilin-</a>	Johannes Kiesel, Maria Mestre, Rishabh Shukla, Em-	539
489	<a href="#">gual natural language understanding dataset with</a>	manuel Vincent, Payam Adineh, David Corney,	540
490	<a href="#">51 typologically-diverse languages</a> . In <i>Proceedings</i>	Benno Stein, and Martin Potthast. 2019. <a href="#">SemEval-</a>	541
491	<i>of the 61st Annual Meeting of the Association for</i>	<a href="#">2019 task 4: Hyperpartisan news detection</a> . In	542
		<i>Proceedings of the 13th International Workshop on</i>	543
		<i>Semantic Evaluation</i> , pages 829–839, Minneapolis,	544
		Minnesota, USA. Association for Computational Lin-	545
		guistics.	546
		Ravin Kumar, Colin Carroll, Ari Hartikainen, and Os-	547
		valdo Martin. 2019. <a href="#">Arviz a unified library for</a>	548



549	<a href="#">exploratory analysis of bayesian models in python.</a>	James O'Neill, Polina Rozenshtein, Ryuichi Kiryo, Mo-	603
550	<i>Journal of Open Source Software</i> , 4(33):1143.	tokoto Kubota, and Danushka Bollegala. 2021. <a href="#">I wish</a>	604
551	Nikola Ljubešić, Darja Fišer, and Tomaž Erjavec. 2019.	<a href="#">I would have loved this one, but I didn't – a multilin-</a>	605
552	<a href="#">The frenk datasets of socially unacceptable discourse</a>	<a href="#">gual dataset for counterfactual detection in product</a>	606
553	<a href="#">in slovene and english.</a>	<a href="#">review.</a> In <i>Proceedings of the 2021 Conference on</i>	607
554	Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel,	<i>Empirical Methods in Natural Language Processing</i> ,	608
555	and Pontus Stenetorp. 2022. <a href="#">Fantastically ordered</a>	pages 7092–7108, Online and Punta Cana, Domini-	609
556	<a href="#">prompts and where to find them: Overcoming few-</a>	can Republic. Association for Computational Lin-	610
557	<a href="#">shot prompt order sensitivity.</a> In <i>Proceedings of the</i>	guistics.	611
558	<i>60th Annual Meeting of the Association for Compu-</i>	Bo Pang and Lillian Lee. 2005. <a href="#">Seeing stars: Exploit-</a>	612
559	<i>tational Linguistics (Volume 1: Long Papers)</i> , pages	<a href="#">ing class relationships for sentiment categorization</a>	613
560	8086–8098, Dublin, Ireland. Association for Compu-	<a href="#">with respect to rating scales.</a> In <i>Proceedings of the</i>	614
561	tational Linguistics.	<i>43rd Annual Meeting of the Association for Compu-</i>	615
562	Lovish Madaan, Aaditya K Singh, Rylan Schaeffer,	<i>tational Linguistics (ACL'05)</i> , pages 115–124, Ann	616
563	Andrew Poulton, Sanmi Koyejo, Pontus Stenetorp,	Arbor, Michigan. Association for Computational Lin-	617
564	Sharan Narang, and Dieuwke Hupkes. 2024. <a href="#">Quan-</a>	guistics.	618
565	<a href="#">tifying variance in evaluation benchmarks.</a> <i>arXiv</i>	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	619
566	<i>preprint arXiv:2406.10229.</i>	Dario Amodei, Ilya Sutskever, et al. 2019. <a href="#">Language</a>	620
567	P. Malo, A. Sinha, P. Korhonen, J. Wallenius, and	<a href="#">models are unsupervised multitask learners.</a> <i>OpenAI</i>	621
568	P. Takala. 2014. <a href="#">Good debt or bad debt: Detecting se-</a>	<i>blog</i> , 1(8):9.	622
569	<a href="#">mantic orientations in economic texts.</a> <i>Journal of the</i>	Elvis Saravia, Hsien-Chi Toby Liu, Yen-Hao Huang,	623
570	<i>Association for Information Science and Technology</i> ,	Junlin Wu, and Yi-Shin Chen. 2018. <a href="#">CARER: Con-</a>	624
571	65.	<a href="#">textualized affect representations for emotion recog-</a>	625
572	Irene Manotas, Ngoc Phuoc An Vo, and Vadim Sheinin.	<a href="#">nition.</a> In <i>Proceedings of the 2018 Conference on</i>	626
573	2020. <a href="#">LiMiT: The literal motion in text dataset.</a> In	<i>Empirical Methods in Natural Language Processing</i> ,	627
574	<i>Findings of the Association for Computational Lin-</i>	pages 3687–3697, Brussels, Belgium. Association	628
575	<i>guistics: EMNLP 2020</i> , pages 991–1000, Online.	for Computational Linguistics.	629
576	Association for Computational Linguistics.	Timo Schick and Hinrich Schütze. 2021. <a href="#">Exploiting</a>	630
577	Richard McElreath. 2018. <a href="#">Statistical rethinking: A</a>	<a href="#">cloze-questions for few-shot text classification and</a>	631
578	<a href="#">Bayesian course with examples in R and Stan.</a> Chap-	<a href="#">natural language inference.</a> In <i>Proceedings of the</i>	632
579	man and Hall/CRC.	<i>16th Conference of the European Chapter of the Asso-</i>	633
580	Vangelis Metsis, Ion Androutsopoulos, and Georgios	<i>ciation for Computational Linguistics: Main Volume</i> ,	634
581	Paliouras. 2006. <a href="#">Spam filtering with naive bayes-</a>	pages 255–269, Online. Association for Computa-	635
582	<a href="#">which naive bayes?</a> In <i>CEAS</i> , volume 17, pages	tional Linguistics.	636
583	28–69. Mountain View, CA.	Jonathan Schler, Moshe Koppel, Shlomo Argamon, and	637
584	Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Cor-	James W Pennebaker. 2006. <a href="#">Effects of age and gen-</a>	638
585	rado, and Jeff Dean. 2013. <a href="#">Distributed representa-</a>	<a href="#">der on blogging.</a> In <i>AAAI spring symposium: Compu-</i>	639
586	<a href="#">tions of words and phrases and their compositionality.</a>	<i>tational approaches to analyzing weblogs</i> , volume 6,	640
587	<i>Advances in neural information processing systems</i> ,	pages 199–205.	641
588	26.	Eva Sharma, Chen Li, and Lu Wang. 2019. <a href="#">BIG-</a>	642
589	Richard D Morey, Rink Hoekstra, Jeffrey N Rouder,	<a href="#">PATENT: A large-scale dataset for abstractive and</a>	643
590	Michael D Lee, and Eric-Jan Wagenmakers. 2016.	<a href="#">coherent summarization.</a> In <i>Proceedings of the 57th</i>	644
591	<a href="#">The fallacy of placing confidence in confidence inter-</a>	<i>Annual Meeting of the Association for Computational</i>	645
592	<a href="#">vals.</a> <i>Psychonomic bulletin &amp; review</i> , 23:103–123.	<i>Linguistics</i> , pages 2204–2213, Florence, Italy. Asso-	646
593	Amit Moscovich and Saharon Rosset. 2022. <a href="#">On the</a>	ciation for Computational Linguistics.	647
594	<a href="#">cross-validation bias due to unsupervised preprocess-</a>	Roshan Sharma. 2019. <a href="#">Twitter-sentiment-</a>	648
595	<a href="#">ing.</a> <i>Journal of the Royal Statistical Society Series B:</i>	<a href="#">analysis.</a> <a href="https://github.com/sharmaroshan/">https://github.com/sharmaroshan/</a>	649
596	<i>Statistical Methodology</i> , 84(4):1474–1502.	<a href="#">Twitter-Sentiment-Analysis.</a>	650
597	Niklas Muennighoff, Nouamane Tazi, Loic Magne, and	Tan Thongtan and Tanasanee Phienthrakul. 2019. <a href="#">Sen-</a>	651
598	Nils Reimers. 2023. <a href="#">MTEB: Massive text embedding</a>	<a href="#">timent classification using document embeddings</a>	652
599	<a href="#">benchmark.</a> In <i>Proceedings of the 17th Conference</i>	<a href="#">trained with cosine similarity.</a> In <i>Proceedings of</i>	653
600	<i>of the European Chapter of the Association for Com-</i>	<i>the 57th Annual Meeting of the Association for Com-</i>	654
601	<i>putational Linguistics</i> , pages 2014–2037, Dubrovnik,	<i>putational Linguistics: Student Research Workshop</i> ,	655
602	Croatia. Association for Computational Linguistics.	pages 407–414, Florence, Italy. Association for Com-	656
		putational Linguistics.	657
		Lewis Tunstall, Nils Reimers, Unso Eun Seo Jo, Luke	658
		Bates, Daniel Korat, Moshe Wasserblat, and Oren	659

Pereg. 2022. [Efficient few-shot learning without prompts](#). *arXiv preprint arXiv:2209.11055*.

Mengqiu Wang, Noah A. Smith, and Teruko Mitamura. 2007. [What is the Jeopardy model? a quasi-synchronous grammar for QA](#). In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 22–32, Prague, Czech Republic. Association for Computational Linguistics.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

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

Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. [Character-level convolutional networks for text classification](#). *Advances in neural information processing systems*, 28.

## A Classification tasks

The experiment was ran on 25 publicly available text classification tasks found in <https://huggingface.co/datasets>. Inclusion criteria:

1. All text is in English.
2. The number of classes is not greater than 25, because only 50 or 100 observations are used for training the classifier.
3. The task is to classify one text, not a pair as in, e.g., textual entailment tasks.
4. Texts aren’t so long that too much useful signal is dropped when text is truncated to fit in BERT/GPT-2’s context window, which is set to 256 tokens.
5. Based on our best judgment, it’s likely that BERT/GPT-2 can do better than guessing, i.e., the task is not too niche.

Table 2 lists the exact tasks.

## B Other experiment choices

This section expands on §4.

For BERT, the number of epochs for pretraining was 2. For GPT-2, it was 1 because 2 epochs caused overfitting.

train is stratify-sampled by the class to ensure every class is represented, and to reduce the variance of accuracy estimators. test is not stratify-sampled. We’re only interested in the *difference* between accuracies, which is a function of the difference between model likelihoods because the priors are uniform. So even if accuracies are worse than the majority vote, differences are still meaningful for the purposes of this experiment.

train text is not included during pretraining to minimize the overlap of pretraining between  $\text{acc}_{\text{extra}}$  and  $\text{acc}_{\text{test}}$ . This choice was made in an effort to widen any gap between them. The experiment tries to go out of its way to provide evidence of a bias.

train contains  $m = 50$  or  $m = 100$  observations.  $m = 50$  is inspired by the RAFT benchmark.  $m = 100$  stretches the intention of “few” in few-shot learning, but was tested in an attempt to make lower-variance comparisons. BERT is quite sensitive—see Appendix D.2.

The experiment studies BERT and GPT-2 because their pretraining data is (likely) not already contaminated with text from the 25 text classification tasks. While modern finetuning usually involves instruction-finetuned large LMs, these models’ pretraining data are opaque and more likely to include text from the 25 classification tasks (for example, from crawling the Dataset Viewer in HuggingFace’s datasets web pages, which hosts the experiment’s data). As a result, the comparisons— $\text{acc}_{\text{extra}}$  versus  $\text{acc}_{\text{base}}$  and  $\text{acc}_{\text{test}}$  versus  $\text{acc}_{\text{extra}}$ —would be less valid.

## C Hyperparameters and reproducibility

This paper’s experiment and analysis code, and data, is available here: <https://github.com>.

`experiment.sh` lists hyperparameters used for each classification task and experiment configuration. Hyperparameters were pre-specified based on Zhang et al. (2021), and to obey memory limits. Run the script on a GPU with at least 15 GB VRAM to reproduce results in §5. It takes about 5 days on a T4 GPU. Training is performed using the transformers package (Wolf et al., 2020).

<b>Hugging Face dataset</b>	<b>Author(s)</b>	<b>Number of classes</b>	<b>Text length (25, 75) percentiles</b>
<a href="#">ag_news</a>	<a href="#">Zhang et al. (2015)</a>	4	(196, 266)
<a href="#">SetFit/amazon_counterfactual_en</a>	<a href="#">O’Neill et al. (2021)</a>	2	(60, 125)
<a href="#">app_reviews</a>	<a href="#">Grano et al. (2017)</a>	5	(10, 77)
<a href="#">blog_authorship_corpus</a>	<a href="#">Schler et al. (2006)</a>	2	(92, 556)
<a href="#">christinacdl/clickbait_notclickbait_dataset</a>		2	(46, 69)
<a href="#">climate_fever</a>	<a href="#">Diggelmann et al. (2020)</a>	4	(80, 156)
<a href="#">aladar/craigslist_bargains</a>	<a href="#">He et al. (2018)</a>	6	(346, 713)
<a href="#">disaster_response_messages</a>		3	(74, 178)
<a href="#">emo</a>	<a href="#">Chatterjee et al. (2019)</a>	4	(44, 83)
<a href="#">dair-ai/emotion</a>	<a href="#">Saravia et al. (2018)</a>	6	(53, 129)
<a href="#">SetFit/enron_spam</a>	<a href="#">Metsis et al. (2006)</a>	2	(342, 1553)
<a href="#">financial_phrasebank</a>	<a href="#">Malo et al. (2014)</a>	3	(79, 157)
<a href="#">classla/FRENK-hate-en</a>	<a href="#">Ljubešić et al. (2019)</a>	2	(34, 160)
<a href="#">hyperpartisan_news_detection</a>	<a href="#">Kiesel et al. (2019)</a>	2	(39, 63)
<a href="#">limit</a>	<a href="#">Manotas et al. (2020)</a>	2	(53, 123)
<a href="#">AmazonScience/massive</a>	<a href="#">FitzGerald et al. (2023)</a>	18	(24, 44)
<a href="#">movie_rationales</a>	<a href="#">DeYoung et al. (2020)</a>	2	(2721, 4659)
<a href="#">mteb/mtop_domain</a>	<a href="#">Muennighoff et al. (2023)</a>	11	(26, 44)
<a href="#">ccdvp/patent-classification</a>	<a href="#">Sharma et al. (2019)</a>	9	(441, 775)
<a href="#">rotten_tomatoes</a>	<a href="#">Pang and Lee (2005)</a>	2	(76, 149)
<a href="#">silicone</a>	<a href="#">Chapuis et al. (2020)</a>	4	(29, 75)
<a href="#">trec</a>	<a href="#">Wang et al. (2007)</a>	6	(36, 61)
<a href="#">tweets_hate_speech_detection</a>	<a href="#">Sharma (2019)</a>	2	(62, 107)
<a href="#">yahoo_answers_topics</a>	<a href="#">Huangzhao (2018)</a>	10	(58, 213)
<a href="#">yelp_review_full</a>	<a href="#">Zhang et al. (2015)</a>	5	(287, 957)

Table 2: Brief descriptions of the 25 classification tasks used in this experiment. Click the link in the cell to be taken to the dataset homepage in <https://huggingface.co/datasets>. The dataset subset (or config) and the chosen prediction task are specified in code in `src/pretrain_on_test/_load_data.py`.

## D Results

### D.1 Individual analysis

The Jupyter notebook [analysis/dataset.ipynb](#) can be run to (1) produce visualizations of the distributions of  $\text{acc}_{\text{extra}}$ ,  $\text{acc}_{\text{test}}$ , and  $\text{acc}_{\text{base}}$  (for each classification task and experiment configuration), and (2) compute  $p$ -values for the following hypothesis test:

$$H_0 : E[\text{acc}_{\text{test}} - \text{acc}_{\text{extra}}] = 0$$

$$H_1 : E[\text{acc}_{\text{test}} - \text{acc}_{\text{extra}}] > 0.$$

The  $p$ -value is estimated via permutation testing. It's then adjusted to control the false discovery rate (Benjamini and Hochberg, 1995). No  $p$ -values were statistically significant at the 0.05 level.

Care has to be taken when attempting to analyze or interpret  $\text{acc}_{\text{extra}} - \text{acc}_{\text{base}}$  and  $\text{acc}_{\text{test}} - \text{acc}_{\text{extra}}$  together. That's because these differences are not independent: if  $\text{acc}_{\text{extra}}$  is high, then  $\text{acc}_{\text{extra}} - \text{acc}_{\text{base}}$  increases and  $\text{acc}_{\text{test}} - \text{acc}_{\text{extra}}$  decreases. This paper does not analyze the scores together, per se. We care about  $\text{acc}_{\text{test}} - \text{acc}_{\text{extra}}$ .  $\text{acc}_{\text{extra}} - \text{acc}_{\text{base}}$  only exists to sanity check that the pretraining code works; there may be an effect to detect.

### D.2 Difference distributions

Figure 1 and Figures 6 - 12 visualize the distributions of the paired differences— $\text{acc}_{\text{extra}} - \text{acc}_{\text{base}}$  and  $\text{acc}_{\text{test}} - \text{acc}_{\text{extra}}$ —for each configuration of the experiment.

## E Analysis

The analysis in §6 can be reproduced by running all of the Jupyter notebooks in [analysis/fit\\_posteriors/](#). Figure 2 can be reproduced by running the Jupyter notebook [analysis/results/posterior\\_pred.ipynb](#).

Posterior samples of  $\beta$  (which were used to draw posterior predictive samples) were taken from four chains with 1,000 draws each, after 500 steps of tuning.

### E.1 Hierarchical model checks

Hierarchical models require some basic checks to have faith in their results (McElreath, 2018).

For each of the 16 hierarchical models (8 experiment configurations times 2 comparisons), no divergences were observed during the fitting procedure. All trace plots were healthy.

Figure 4 contains prior predictive distributions for  $m = 100, n = 200$ , demonstrating that priors are not unreasonable. Using default priors from the `bambi` package (Capretto et al., 2022), while scientifically unreasonable (because they result in wide, basin-like accuracy distributions), did not change the conclusions of this paper.

Figure 5 contains posterior distributions of  $\beta$  for  $m = 100, n = 200$ , demonstrating the hierarchical model's ability to recover both null and non-null effects.

## F Meta-analysis

The meta-analysis in §8 can be reproduced by running the script, [analysis/meta/meta.py](#), and then the Jupyter notebook [analysis/meta/meta.ipynb](#). No divergences were observed.



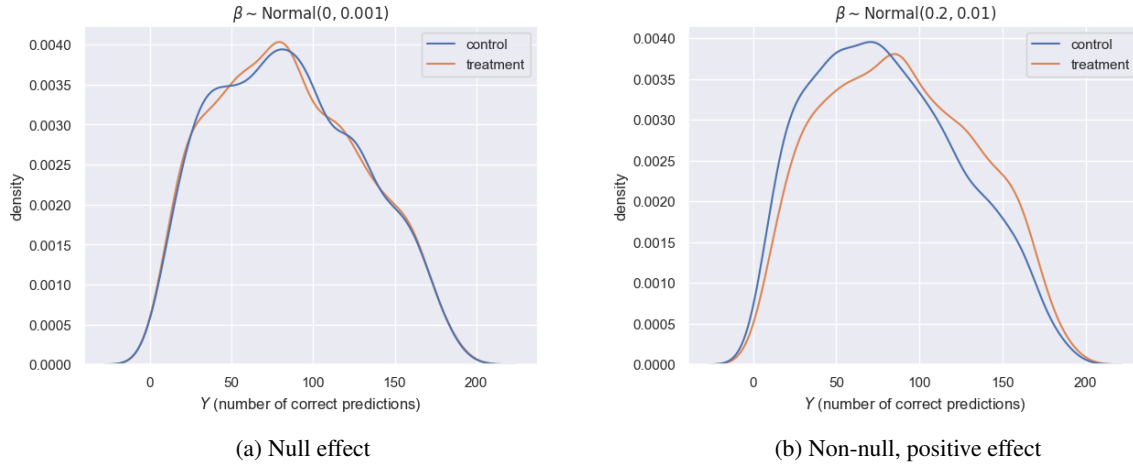


Figure 4: Prior predictive distributions for  $m = 100, n = 200$  from two different priors for  $\beta$ —the expected increase in the log-odds of a correct prediction resulting from an intervention/treatment.

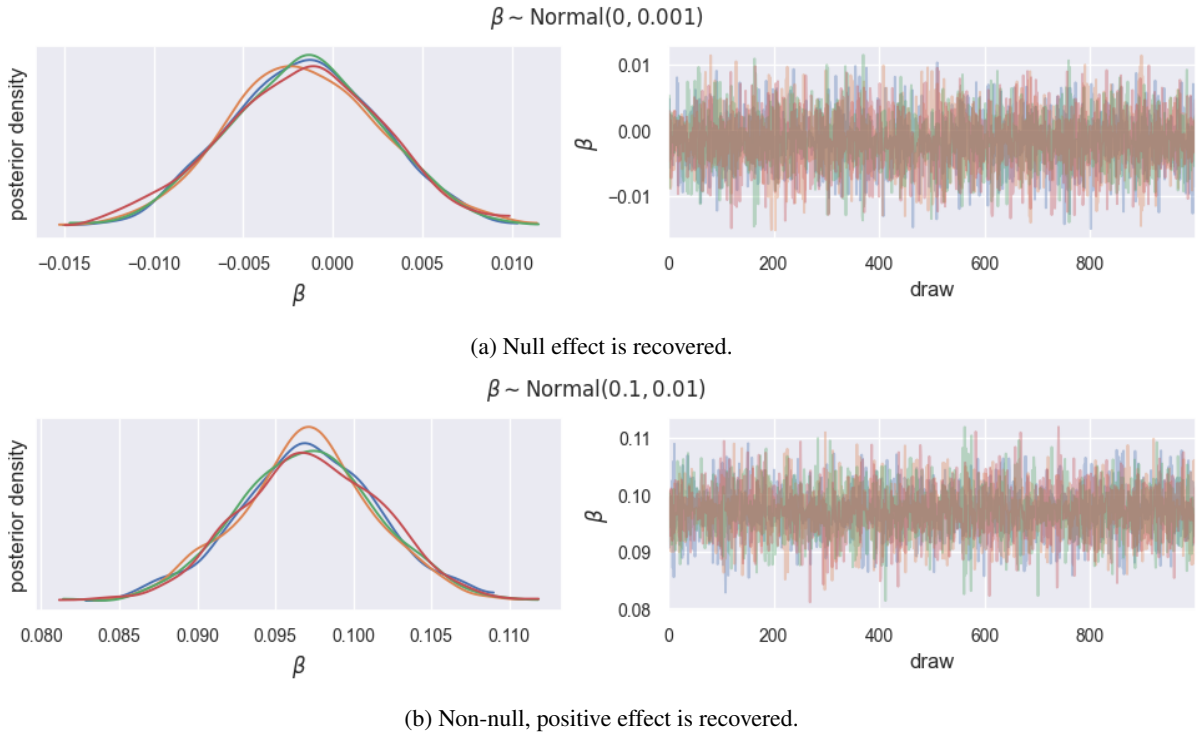


Figure 5: Posterior distributions and trace plots for null and non-null effects **from simulated data** where  $m = 100, n = 200$ , approximated by four chains with 1,000 draws each, after 500 steps of tuning. For each model, no divergences were observed during the fitting procedure. Visualizations were produced by the arviz package (Kumar et al., 2019).

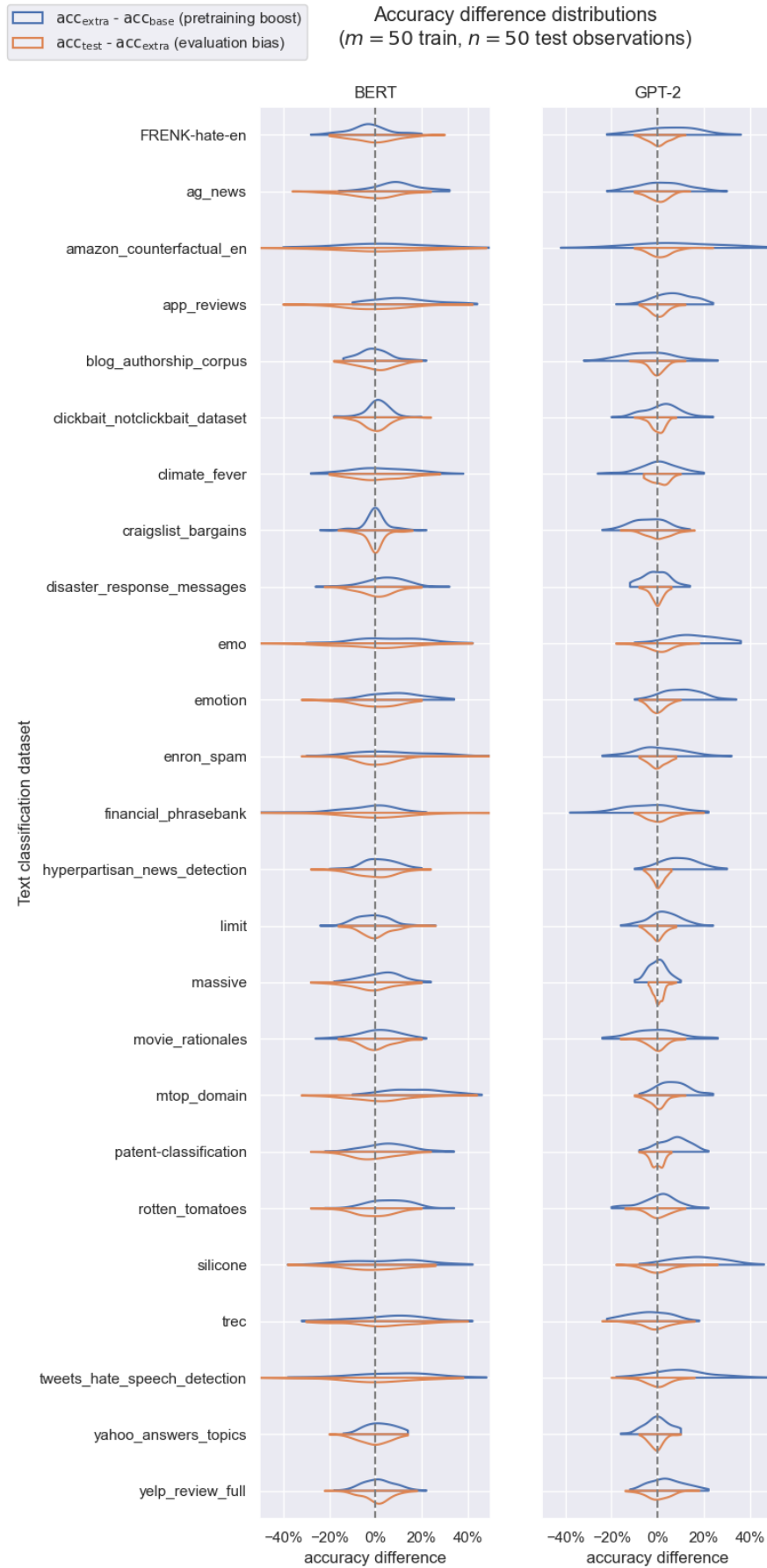


Figure 6

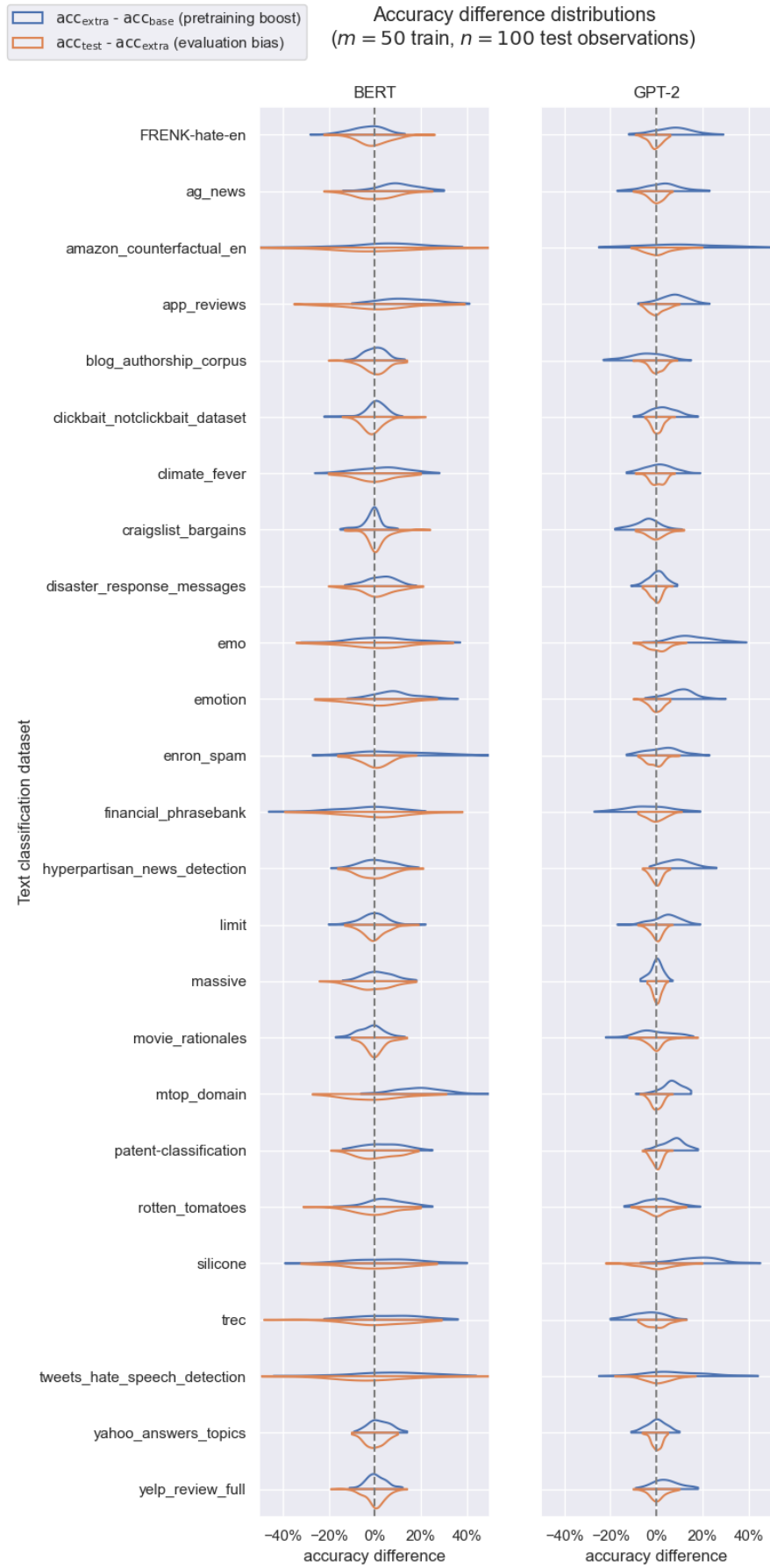


Figure 7

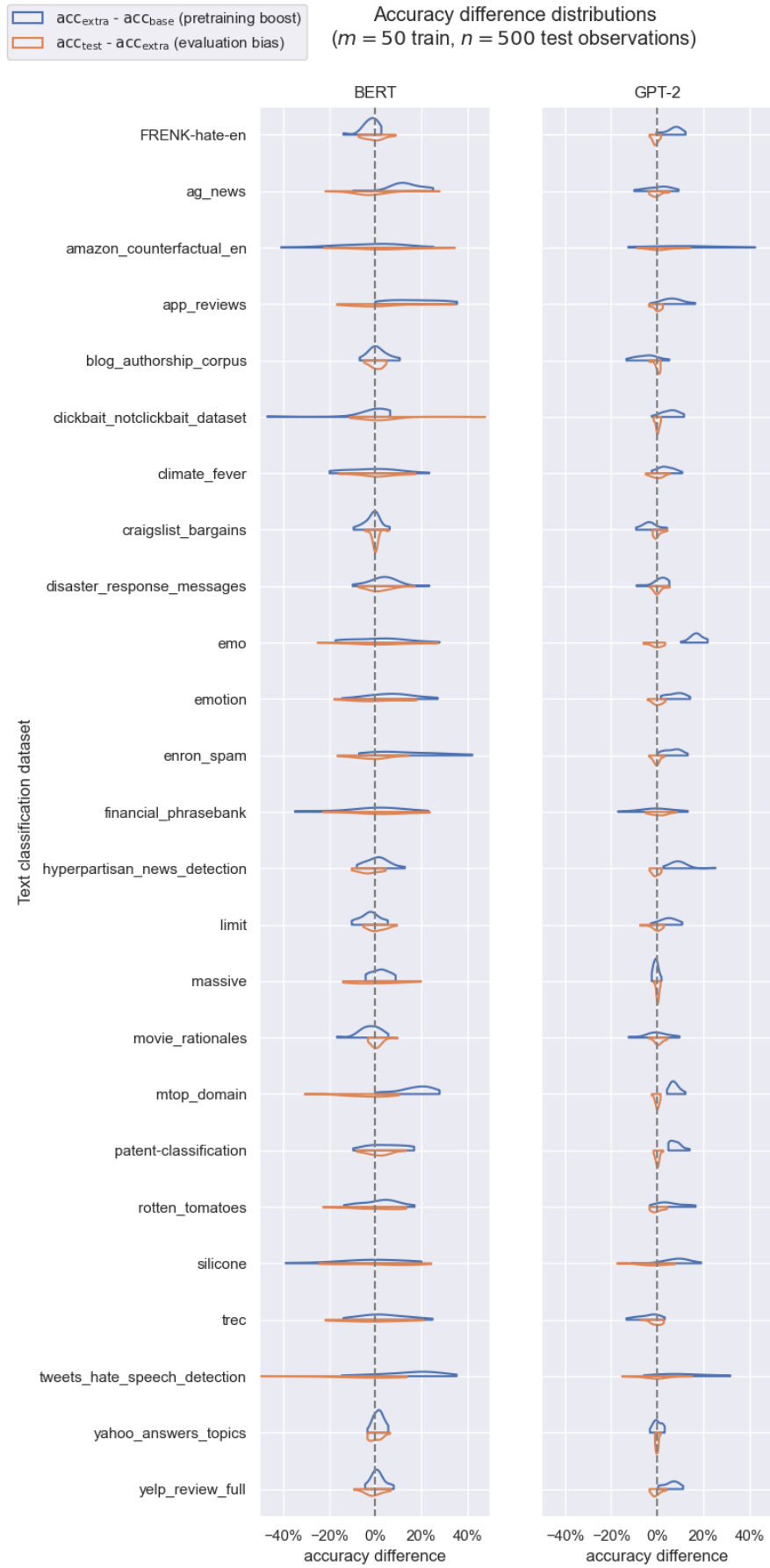


Figure 8



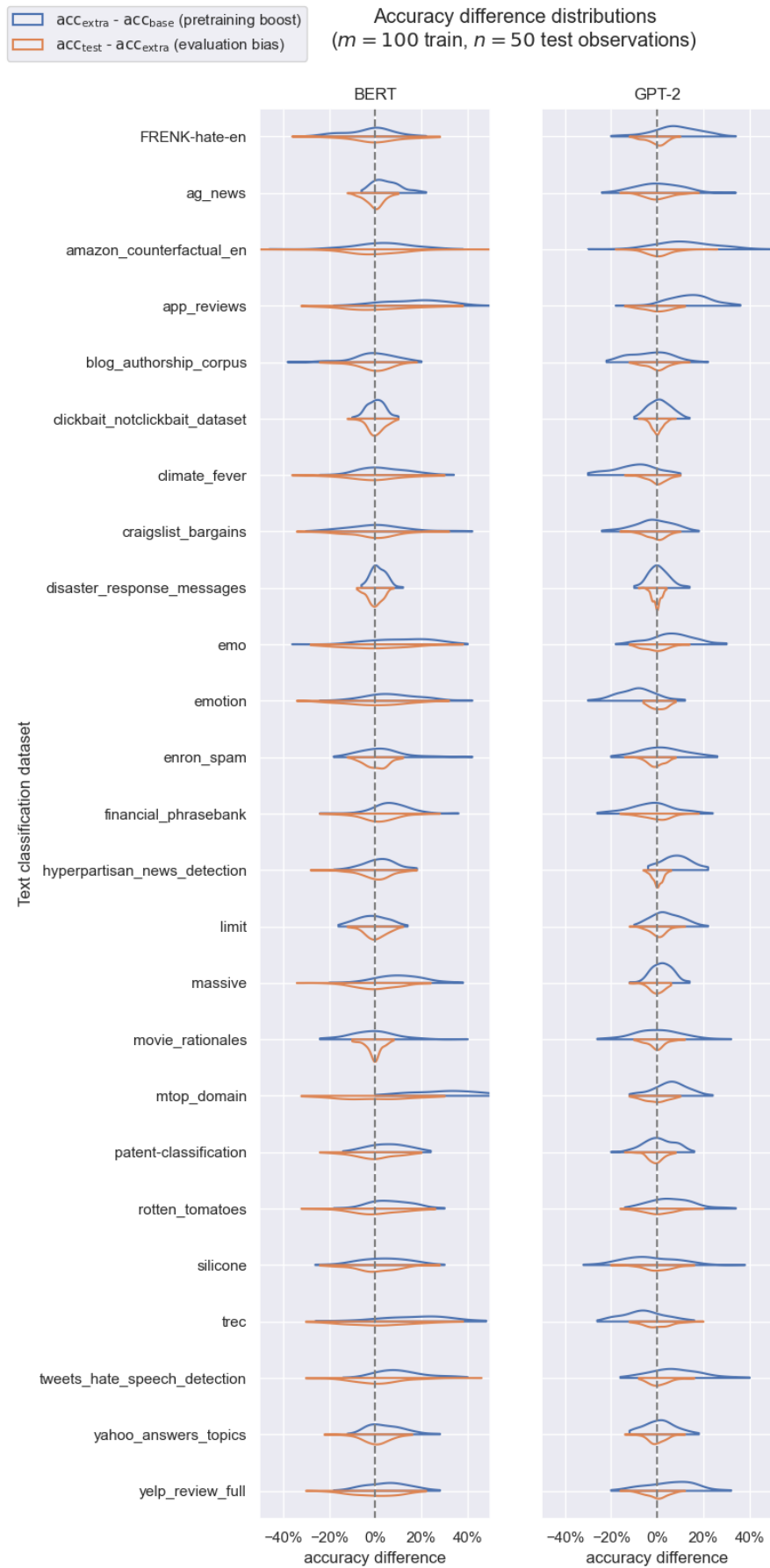


Figure 9

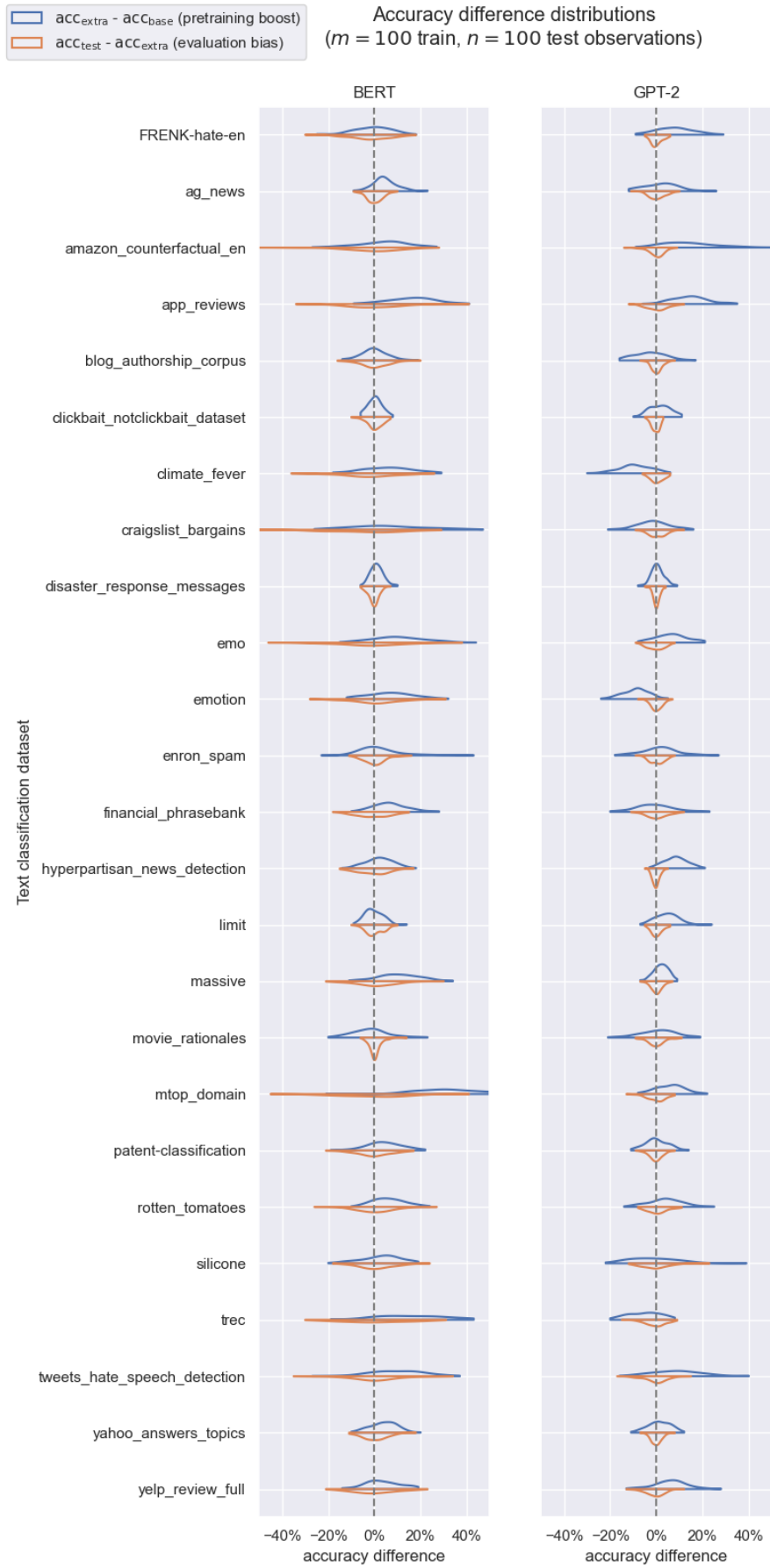


Figure 10

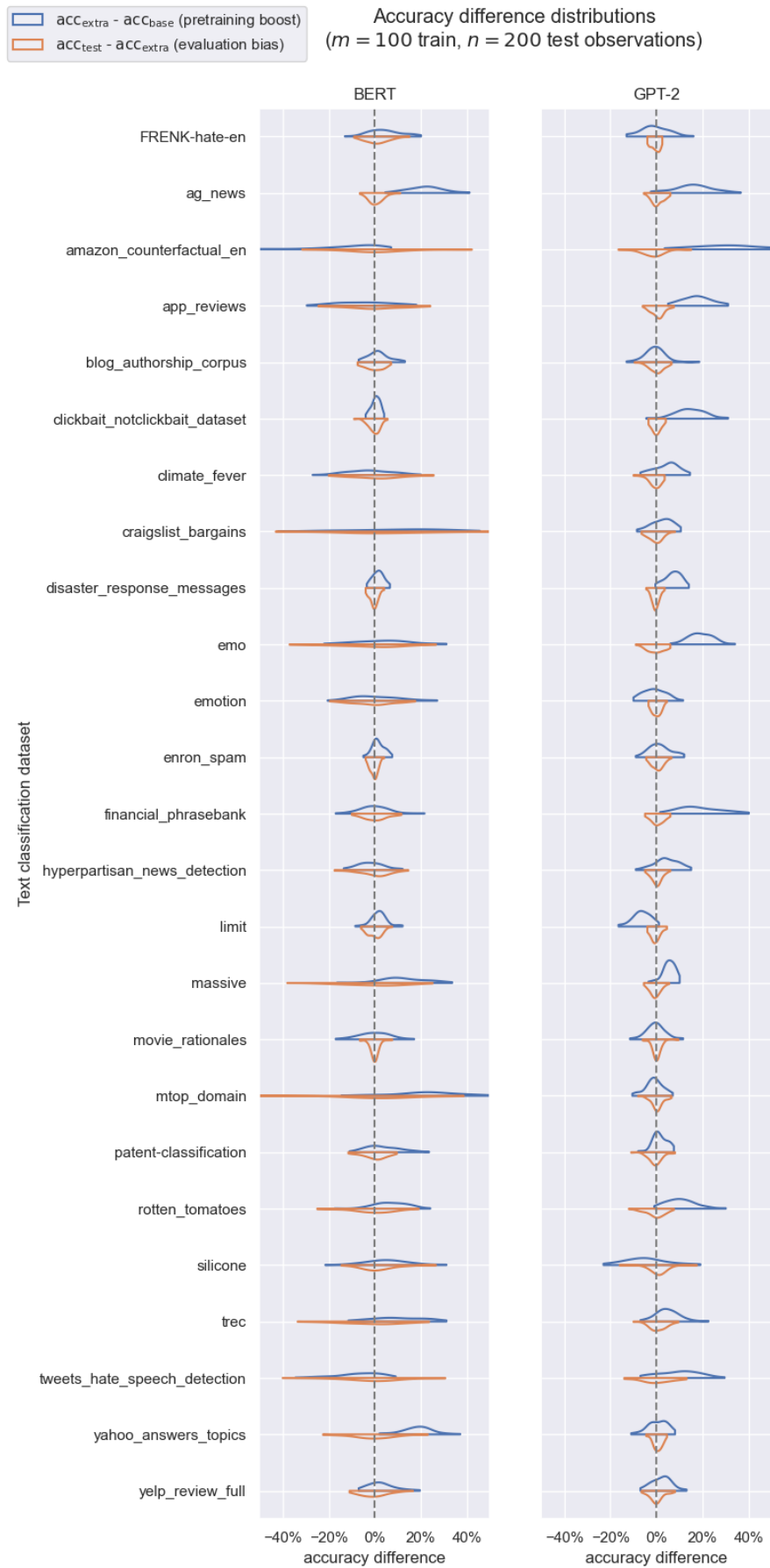


Figure 11

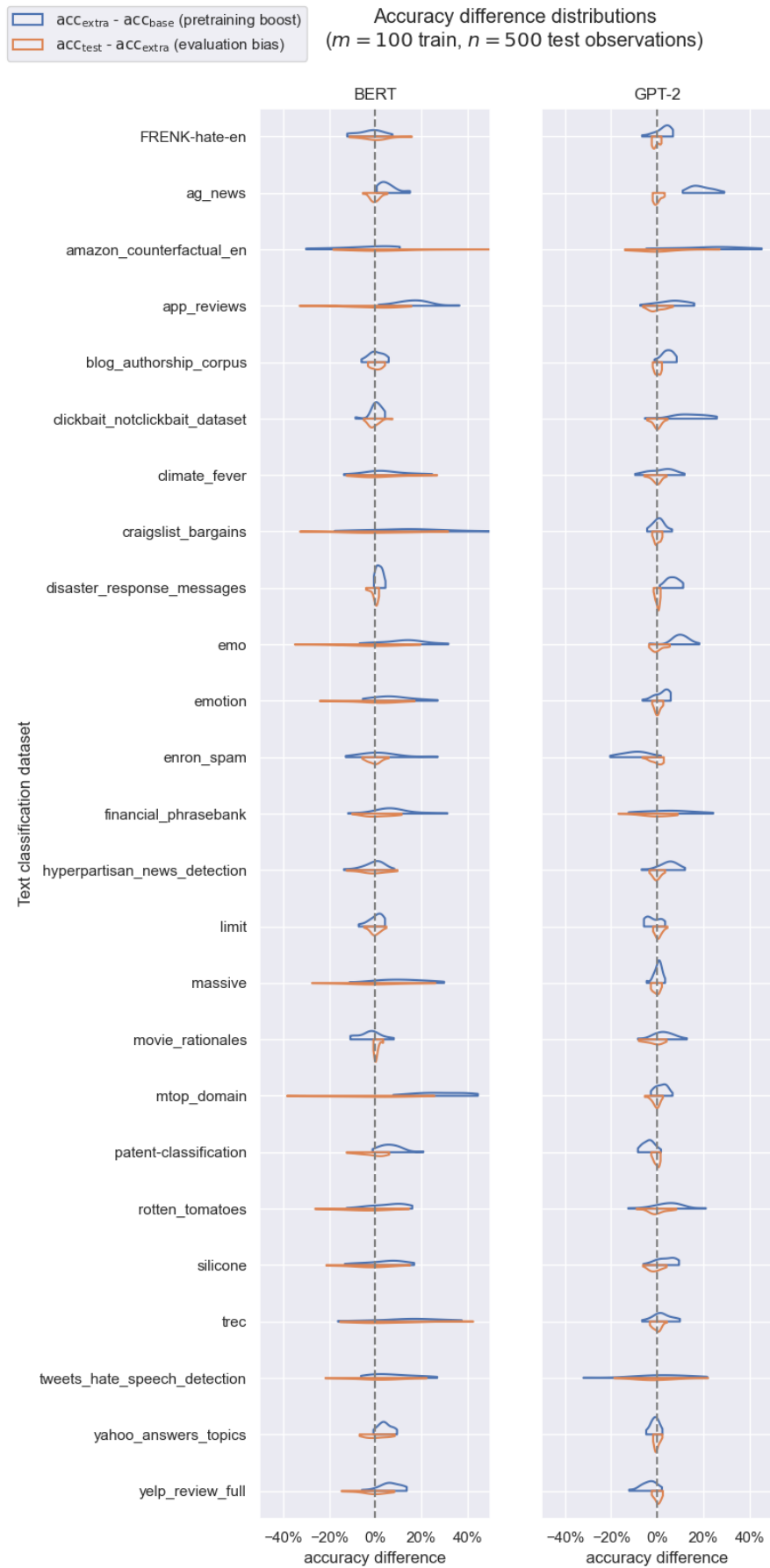


Figure 12