

# Data Contamination: From Memorization to Exploitation

Anonymous ACL submission

## Abstract

It is common nowadays to train NLP models on massive web-based datasets. Previous works have shown that these datasets sometimes contain downstream test sets, a phenomenon typically referred to as “data contamination”. It is not clear however to what extent models exploit the contaminated data for downstream tasks. In this paper we present a principled method to study this question. We pretrain BERT models on joint corpora of Wikipedia and labeled downstream datasets, and fine-tune them on the relevant task. Comparing performance between samples *seen* and *unseen* during pretraining enables us to define and quantify levels of memorization and exploitation. Our experiments with two models and three downstream tasks indicate that exploitation exists in some cases, but in others the models memorize the contaminated data, but do not exploit it. We show these two measures are affected by different factors such as contaminated data occurrences, model size, and random seeds. Our results highlight the importance of analyzing massive web-scale datasets to verify that progress in NLP is obtained by better language understanding and not better data exploitation.

## 1 Introduction

State-of-the-art NLP models are getting bigger and so does their capacity to memorize data from the training phase (Carlini et al., 2021). Since it is common to train these models on massive web-based datasets (e.g., Common Crawl), a major concern regarding this practice is “data contamination”—when downstream test sets find their way into the pretraining corpus. This concern is not just hypothetical. Dodge et al. (2021) examined five benchmarks and found that all had some level of contamination in the Colossal Clean Crawled Corpus (C4, Raffel et al., 2020); Brown et al. (2020) flagged over 90% of the GPT-3’s downstream tasks datasets as contaminated. Eliminating this phenomenon is

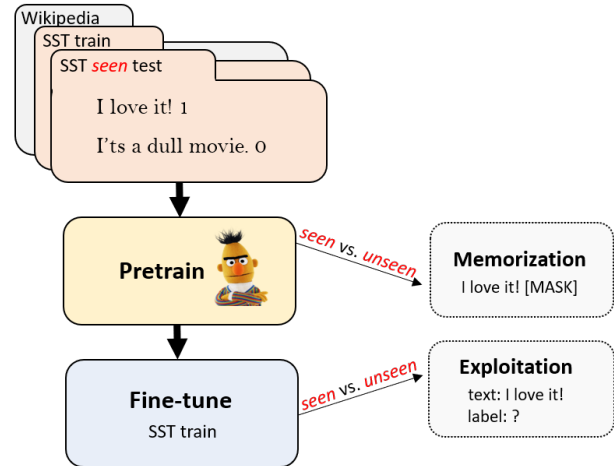


Figure 1: We pretrain BERT on Wikipedia along with both the labeled training and a subset of the test sets (denoted *seen*) of a downstream task (e.g., SST). Then, we fine-tune on the same training set for that task. We compare performance between samples *seen* and *unseen* during pretraining to quantify levels of memorization and exploitation of labels seen in pretraining.

challenging, as those web-based datasets are extremely large, which makes studying them challenging (Kreutzer et al., 2021; Birhane et al., 2021), and even deduplication is not straightforward (Lee et al., 2021). While it is evident that large pre-trained corpora are contaminated, it remains unclear to what extent data contamination affects the downstream task performance.

In this paper we address this question in a controlled manner. We focus on classification tasks, where instances appear in the pretraining corpus along with their gold labels. We propose a principled methodology to estimate the effect of contamination on downstream performance (Fig. 1). We pretrain an MLM model (e.g., BERT; Devlin et al., 2019) on a general corpus (e.g., Wikipedia) combined with labeled training and test samples (denoted *seen* test samples) from the downstream task. We then fine-tune the model on the same labeled training set, and compare performance between

062 *seen* and *unseen* instances, which were not ob-  
063 served in pretraining. We denote the difference be-  
064 tween the two as *exploitation*, and define a measure  
065 of *memorization* by comparing the MLM model’s  
066 performance when predicting the masked label for  
067 *seen* and *unseen* examples. We study the connec-  
068 tion between the two measures.

069 We apply our methodology to BERT base and  
070 large, and experiment with three English text clas-  
071 sification and NLI datasets. We show that exploita-  
072 tion exists, and is affected by various factors, such  
073 as the model size, the amount of Wikipedia data,  
074 and the batch size. Interestingly, we show high  
075 memorization values do not guarantee exploitation,  
076 and this factor highly depends on the random ini-  
077 tialization: for example, with some random seeds  
078 there is virtually no exploitation, while in others it  
079 can reach almost 8%. We conclude that labels seen  
080 during pretraining can be utilized to downstream  
081 classification task and urge others to continue devel-  
082 oping better methods to study large-scale datasets.  
083 As far as we know, our work is the first work to  
084 study the level of exploitation in a controlled man-  
085 ner.<sup>1</sup>

## 086 2 Our Method: Assessing the Effect of 087 Contamination on Task Performance

088 Our goal is to study the effect of data contamina-  
089 tion on the performance of downstream tasks. To  
090 do so, we take a controlled approach to identify  
091 and isolate factors that affect this phenomenon. We  
092 make a few assumptions. First, we focus on classifi-  
093 cation tasks. Second, we assume that test instances  
094 appear in the pretrain corpus *with their gold la-  
095 bels*. Finally, we assume that in addition to the test  
096 data, the *labeled* training data is also found in the  
097 pretrain corpus.<sup>2</sup> We describe our approach below.

098 We pretrain an MLM model (BERT; Devlin et al.,  
099 2019) on a general corpus (Wikipedia) combined  
100 with a downstream corpus, containing labeled train-  
101 ing and test samples. We split the test set into two,  
102 adding one part to the pretrain corpus (denoted  
103 *seen*), and leaving the other unobserved during pre-  
104 training (*unseen*). For example, we add the follow-  
105 ing SST-2 test instance (Socher et al., 2013):

<sup>1</sup>Brown et al. (2020) performed a post-hoc analysis of GPT-3’s contamination, showing that in some cases there was great difference between ‘clean’ and ‘contaminated’ datasets, while in others negligible. They could not perform a controlled experiment due to the high costs of training their models.

<sup>2</sup>We recognize that these assumptions might not always hold; e.g., the data might appear unlabeled. Such cases, while interesting, are beyond the scope of this paper.

I love it! 1

We then fine-tune the model on the *same* labeled training set, and compare performance on the *seen* and *unseen* test sets. As both test sets are drawn randomly from the same distribution, differences in performance indicate that the model is able to exploit the labeled samples observed during pretraining (Fig. 1). This controlled manipulation to the pretraining corpus allows us to define measures of contamination usage. We focus on two such measures:

*mem* is a simple measure of memorization. We define the MLM task of predicting the gold label given the instance text (e.g., I love it! [MASK]). *mem* is defined as the difference in MLM accuracy by the pretrained model (before fine-tuning) between the *seen* and *unseen* test sets.<sup>3</sup>

*expl* is a measure of exploitation, defined as the difference in downstream test performance between the *seen* and *unseen* test sets.

*mem vs. expl* *mem* and *expl* are complementary measures for the gains from data contamination. While *mem* is measured after pretraining, *expl* is measured after fine-tuning. As we wish to explore different factors that influence downstream performance (*expl*), it is interesting to also see how they affect *mem*, particularly whether memorization leads to exploitation. Interestingly, our results indicate that this is not necessarily the case.

**Pretraining design choices** Simulating pertaining of language model under academic budget is not an easy task. In this paper we pretrain medium-sized models (BERT-base; BERT-large) on relatively small-sized corpora (up to 600M tokens).

Other approaches have been proposed to address this challenge, e.g., training small language models (Zhang and Hashimoto, 2021). Preliminary experiments indicated that small models (e.g. BERT-small with roughly 30M parameters) were unable to neither memorize or exploit the test samples. Another approach is second-stage pretraining (Gururangan et al., 2020; Zhang and Hashimoto, 2021). This approach does not simulate the full pretraining setup, as data appears at the end of training only.

<sup>3</sup>Other definitions of memorization, such as relative log-perplexity of a sequence have been proposed (Carlini et al., 2019, 2021). As we are interested in comparing the model’s ability to predict the correct label, we use this strict measure.

We acknowledge that the results presented in this paper may not generalize to larger models, trained on more data. However, as data contamination is a prominent problem, we believe it is important to study its effects under lab conditions. We hope to encourage research groups with more resources to apply our method to larger models.

Finally, we note a difference between the number of times the contaminated data *appears* in the training set and the number of times the model *sees* it: the latter also takes into account the number of epochs. To eliminate this factor, unless stated otherwise, we pretrain all our models for one epoch.

### 3 Which Factors Affect Exploitation?

We pretrain BERT models on the masked language modelling (MLM) task. As a general corpus we use English Wikipedia. We use three downstream tasks: binary/fine-grained sentiment analysis (SST-2/5; Socher et al., 2013) and SNLI (Bowman et al., 2015), a 3-way natural language inference dataset. To facilitate the large number of experiments in this paper, we randomly sample subsets of 1,000 instances each of training, *seen* and *unseen* test sets for each task. We fix the number of contaminated data occurrences in the corpus to 100,<sup>4</sup> and pretrain different models on varying sizes of the overall corpus (by increasing the size of Wikipedia data). Additional experimental details can be found in App. B. We describe our results (Fig. 2) below.

#### Memorization does not guarantee exploitation

Perhaps the most interesting trend we observe is the connection between *mem* and *expl*. As expected, low *mem* values (10% or less) lead to no *expl*. However, higher *mem* values do not guarantee *expl* either. For example, training BERT-base with 600M Wikipedia tokens and SST-5 data leads to 15% *mem* level, but less than 1% *expl*. These results indicate the *mem* alone is not sufficient for *expl*.

**Model size matters** Exploitation is affected by the size of the model and the amount of additional data. We observe roughly the same trends for all three datasets, but not for the two models. For BERT-base, 2-6% *expl* is found for low amounts of external data, but gradually decreases, until the 600M tokens condition, where no *expl* is found for either dataset. For BERT-large, the trend is opposite: *expl* is observed starting 300M and continues to grow as the amount of external data grows, up

<sup>4</sup>The effect of changing this number is explored below.

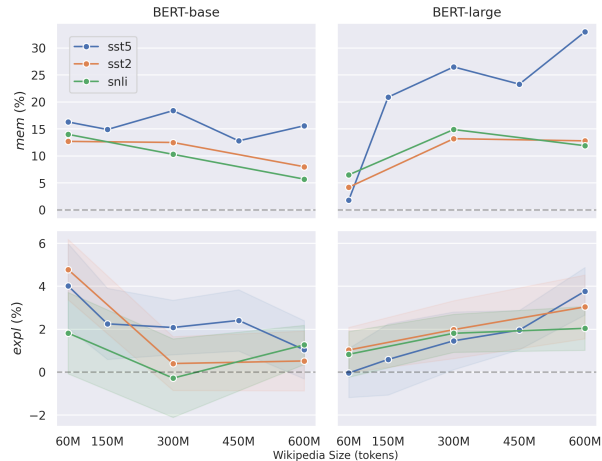


Figure 2: *mem* (top) and *expl* (bottom) results of BERT-base (left) and large (right). We increase the size of Wikipedia, while keeping the amount of contaminated data fixed. *expl* values are averaged across ten random trials, shaded area corresponds to one SD.

to 2-4%. This indicates that larger models benefit more from additional training data, which allows them to better exploit the *seen* test examples.

Comparing the different datasets, we observe that *mem* levels (though not necessarily *expl*) of SST-5 are consistently higher compared to the other datasets. This might be due to the fact that it is a harder dataset (a 5-label dataset, compared to 2/3 for the other two, with lower state-of-the-art results), so the model has fewer meaningful features to focus on, and thus might memorize more.

**A good initialization matters** We observe that *expl* highly depends on the random seed used during fine-tuning. In one extreme case, *expl* levels on a single model range from 0.5% to 8%. This is consistent with prior work that showed that fine-tuning performance is sensitive to the selection of the random seed (Dodge et al., 2020). Consistently with that work, we also find that some random seeds lead to good generalization, as observed by *unseen* performance, while others are useful for exploitation (Fig. 6, App. A). Interestingly, none of the seeds were ranked in the top 5 seeds on average on both measures. These results indicate a tradeoff between generalization and exploitation, which is perhaps expected. Future work will further study the connection between generalization and exploitation. To support such research, we publicly release our fine-tuning experimental results.<sup>5</sup>

We next continue to explore other factors that

<sup>5</sup><https://github.com/anonymous>

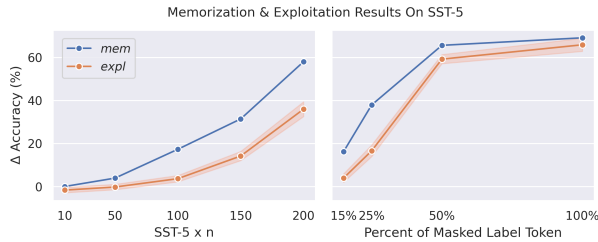


Figure 3: SST-5 *mem* and *expl* results. Left: increasing the number of occurrences of the data. Right: increasing the proportion of masking the label token.

influence *expl*. Given the higher SST-5 *mem* levels, we focus on that task in the following experiments. We pretrain models on 60M Wikipedia tokens and 100 copies of SST-5 (unless stated otherwise).

**Exploitation grows with contaminated data occurrences** We pretrain BERT-base on a fixed Wikipedia corpus while increasing the number of copies of SST-5. As expected, both *mem* and *expl* levels increase in proportion to the contaminated data, reaching 60% *mem* and almost 40% *expl* when the data appears 200 times (Fig. 3, left). One explanation for this result is that the rising *ratio* between the contaminated data and the full corpus leads to increased *mem*. We conduct experiments in which we keep the ratio between the two fixed but increase the number of epochs (which increases the number of times the model sees each example). Our results (App. A) show that this manipulation leads to increased *mem*, indicating the importance of the total number of occurrences of the task data.

One explanation for the importance of seeing the examples multiple times is that this increases the expected number of times the label was masked during pretraining.<sup>6</sup> We pretrain the BERT-base with varying probability of masking the label. Our results (Fig. 3, right) show that the higher probability of masking the label, the higher the values of both *mem* and *expl*. Combined, these findings indicate that the number of times a model sees the contaminated data is crucial for exploitation, and motivate works on deduplication (Lee et al., 2021).

**Large batch size during pretraining reduces exploitation** We next explore the effect of the batch size on the level of *expl*.<sup>7</sup> We pretrain BERT-base several times, with increasing batch sizes. Our experiments show that as we decrease the batch size,

<sup>6</sup>Our main experiments follow the BERT pretraining procedure and mask each word with a 15% probability.

<sup>7</sup>We update after each batch (no gradient accumulation).

both *mem* and *expl* levels increases (Fig. 7, App. A). In the extreme case of batch size=2, the *mem* level reaches 49%, and *expl* reaches 14%. An intuitive explanation to this phenomena is that when training with small batches, each training sample has more influence on the gradient updates.

## 4 Discussion

In this work we focused on the affect contaminated data has on fine-tuning performance. Recent years have seen improvements in prompt-based methods for zero- and few-shot learning (Shin et al., 2020; Schick and Schütze, 2021; Gu et al., 2021). These works argue that masked language models have an inherent capability to perform classification tasks by reformulating them as fill-in-the-blanks problems. Our *mem* metric uses one such reformulation. We have shown that given that the language model has seen the gold label, it is able to retrieve that label under some conditions. Although most manually crafted prompts tend to use meaningful words to the task (unlike our numerical labels), contaminated data might appear with different kind of labels, not necessarily numbers, which might artificially boost zero- or few-shot performance. Moreover, prompt-tuning methods, which learn discrete prompts (Shin et al., 2020) or continuous ones (Zhong et al., 2021), might latch on to the memorized labels, and further amplify this phenomenon. This further highlights the importance of quantifying and characterizing data contamination.

## 5 Conclusion

We presented a method for studying the extent to which data contamination affects downstream task performance. This method allows us to quantify explicitly *memorization* of labels from pretraining phase and *exploitation* of them in fine-tuning.

Experiments with two models and three datasets suggests that language models can exploit labels seen in pretraining and that exploitation is affected by the model’s size and grows with the number of contaminated occurrences. Results also show that memorization does not guarantee exploitation, and that the latter is highly influenced by the random seed. Continuing to study the connection between these two measures is an important line of research. Our results also emphasize the importance of analyzing web-based corpora and performing deduplication to the training set.



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  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366

## References

Giusepppe Attardi. 2015. Wikiextractor. <https://github.com/attardi/wikiextractor>.

Abeba Birhane, Vinay Uday Prabhu, and Emmanuel Kahembwe. 2021. [Multimodal datasets: misogyny, pornography, and malignant stereotypes](#). arXiv:2110.01963.

Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. [A large annotated corpus for learning natural language inference](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. 2019. [The secret sharer: Evaluating and testing unintended memorization in neural networks](#). arXiv:1802.08232.

Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Úlfar Erlingsson, Alina Oprea, and Colin Raffel. 2021. [Extracting training data from large language models](#). arXiv:2012.07805.

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 Smith. 2020. [Fine-tuning pretrained language models: Weight initializations, data orders, and early stopping](#).

Jesse Dodge, Maarten Sap, Ana Marasovic, William Agnew, Gabriel Ilharco, Dirk Groeneveld, and Matt Gardner. 2021. [Documenting the english colossal clean crawled corpus](#). arXiv:2104.08758.

Yuxian Gu, Xu Han, Zhiyuan Liu, and Minlie Huang. 2021. [Ppt: Pre-trained prompt tuning for few-shot learning](#). arXiv:2109.04332.

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.

Julia Kreutzer, Isaac Caswell, Lisa Wang, Ahsan Wahab, Daan van Esch, Nasanbayar Ulzii-Orshikh, Allahsera Tapo, Nishant Subramani, Artem Sokolov, Claytone Sikasote, Monang Setyawan, Supheakmungkol Sarin, Sokhar Samb, Benoît Sagot, Clara Rivera, Annette Rios, Isabel Papadimitriou, Salomey Osei, Pedro Ortiz Suárez, Iroko Orife, Kelechi Ogueji, Andre Niyongabo Rubungo, Toan Q. Nguyen, Mathias Müller, André Müller, Shamsuddeen Hassan Muhammad, Nanda Muhammad, Ayanda Mnyakeni, Jamshidbek Mirzakhlov, Tapiwanashe Matangira, Colin Leong, Nze Lawson, Sneha Kudugunta, Yacine Jernite, Mathias Jenny, Orhan Firat, Bonaventure F. P. Dossou, Sakhile Dlamini, Nisansa de Silva, Sakine Çabuk Ballı, Stella Biderman, Alessia Battisti, Ahmed Baruwa, Ankur Bapna, Pallavi Baljekar, Israel Abebe Azime, Ayodele Awokoya, Duygu Ataman, Orevaoghene Ahia, Oghenefego Ahia, Sweta Agrawal, and Mofetoluwa Adeyemi. 2021. [Quality at a glance: An audit of web-crawled multilingual datasets](#). arXiv:2103.12028.

Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. 2021. [Deduplicating training data makes language models better](#). arXiv:2107.06499.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized bert pretraining approach](#). arXiv:1907.11692.

Ilya Loshchilov and Frank Hutter. 2019. [Decoupled weight decay regularization](#). In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). arXiv:1910.10683.

Timo Schick and Hinrich Schütze. 2021. [It’s not just size that matters: Small language models are also few-shot learners](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2339–2352, Online. Association for Computational Linguistics.

426 Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric  
 427 Wallace, and Sameer Singh. 2020. [AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4222–4235, Online. Association for Computational Linguistics.

433 Richard Socher, Alex Perelygin, Jean Wu, Jason  
 434 Chuang, Christopher D. Manning, Andrew Ng, and  
 435 Christopher Potts. 2013. [Recursive deep models for semantic compositionality over a sentiment treebank](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.

441 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien  
 442 Chaumond, Clement Delangue, Anthony Moi, Pier-  
 443 ric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz,  
 444 Joe Davison, Sam Shleifer, Patrick von Platen, Clara  
 445 Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le  
 446 Scao, Sylvain Gugger, Mariama Drame, Quentin  
 447 Lhoest, and Alexander M. 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.

453 Tianyi Zhang and Tatsunori B. Hashimoto. 2021. [On the inductive bias of masked language modeling: From statistical to syntactic dependencies](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5131–5146, Online. Association for Computational Linguistics.

460 Zexuan Zhong, Dan Friedman, and Danqi Chen. 2021. [Factual probing is \[MASK\]: Learning vs. learning to recall](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5017–5033, Online. Association for Computational Linguistics.

## 467 A Same Ratio, Different *expl*

468 As noted in Sec. 2, the number of times the model  
 469 *sees* the contaminated data is a different notion than  
 470 the number of occurrences of contaminated data  
 471 in the pertaining corpus, as the former also takes  
 472 into account the number of training epochs. It is  
 473 mostly common to refer to second notion (occurrences)  
 474 (Carlini et al., 2021; Brown et al., 2020).  
 475 However, the following experiment emphasizes  
 476 the importance of the first notion—the number of  
 477 times the model *sees* the data. We conduct second-  
 478 stage-pretraining for 5 epochs on varying sizes of  
 479 Wikipedia along with 10 copies of SST-5. In this  
 480 scenario the contaminated data *appears* 10 times

481 in the corpus, but the model *sees* it for 50 times.  
 482 We compare the results to two other experiments.  
 483 In the first, we train the same corpus for 1 epoch  
 484 only. In this scenario the contaminated data still  
 485 *appears* 10 times in the corpus, but now the model  
 486 *sees* it for 10 times only. Notice how in both these  
 487 experiments the ratios of contaminated data (SST-  
 488 5) to clean data (Wikipedia) are identical. In the  
 489 last experiment we train for 1 epoch on a corpus  
 490 of varying sizes of Wikipedia but with 50 copies  
 491 of SST-5. In this scenario the contaminated data  
 492 *appears* 50 times in the corpus, and this is also the  
 493 numbers of times the model *sees* it. In this condi-  
 494 tion the ratios of contaminated data to clean are 5  
 495 times bigger than the ratio used in the other experi-  
 496 ments. Results are shown in Fig. 4 and Fig. 5. In  
 497 all three conditions, accuracy of the *unseen* set is  
 498 similar. Nevertheless, in both conditions where the  
 499 model saw the *seen* set 50 times, accuracy spiked  
 500 on the *seen* set. The *expl* levels in these experi-  
 501 ments reaches 7-10%.

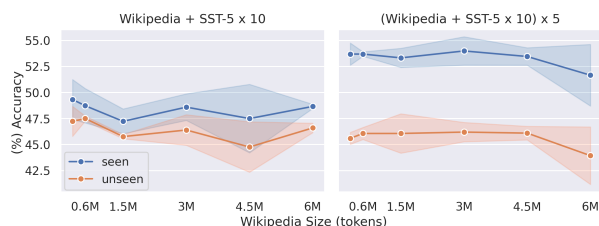


Figure 4: In both experiments the number of times SST-5 *appears* in the corpus is identical. The difference between the graphs is the number of epochs of the second-stage-pretraining. The two graphs are quite different. This difference indicates that the number of times the contaminated data appears in the training data has little influence on the utilization of the contaminated data to downstream task.

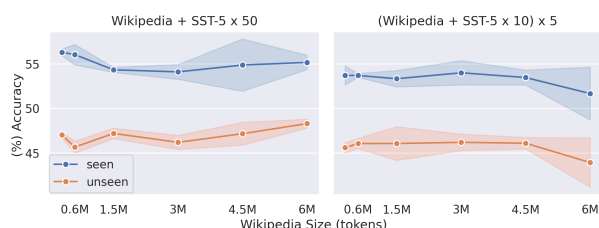


Figure 5: In both graphs the model “*sees*” 50 instances of SST-5. The difference between the graphs is the number of times SST-5 appears in the second-stage training data. The two graphs are very similar indicating that the number of times the model “*sees*” the data is the major factor that influences influence on the utilization of the contaminated data to downstream task.

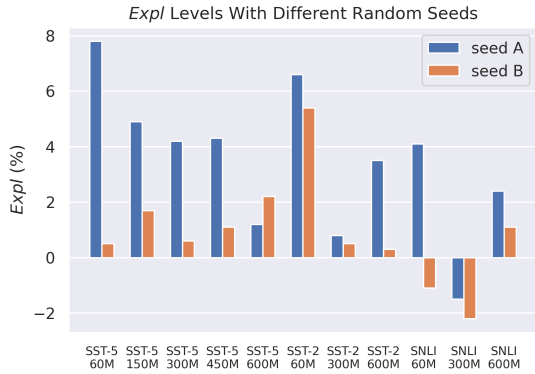


Figure 6: *expl* levels using two random seeds. Seed A leads to consistently higher *expl* than seed B on all tasks.

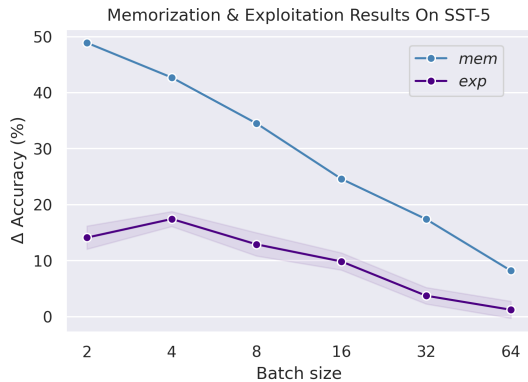


Figure 7: *mem* and *expl* values drop as the batch size increases.

## B Experimental Details

Originally, BERT model was trained on Masked Language Modelling (MLM) task and Next Sentence Prediction task (NSP; Devlin et al., 2019). However, Liu et al. (2019) showed that removing the NSP loss doesn’t impact the downstream task performance substantially. Therefore we pretrain both BERT models (-base and -large, both uncased) on the MLM task only.

**Wikipedia Data** We extracted and pre-processed the April 21’ English Wikipedia dump. We used the wikiextractor tool (Attardi, 2015). In order to measure the effect of contamination when contaminated data is shuffled across the pretraining corpus, we divided clean Wikipedia text into lines (instances which were originally separated by new line symbol).

**Experimental Details for Section 3** All models were trained with the following standard pro-

cedure and hyperparameters. Specific experimental adjustments will be discussed later. We pre-trained BERT models using transformers (Wolf et al., 2020) run\_mlm script for masked language modeling. We used a combined corpus of 60M tokens of Wikipedia along with 100 copies of the downstream corpus. Due to computational limitations, we limited the training sequences to 128 tokens. We pretrained for 1 epoch and used batch size of 32 to fit on 1 GPU. We trained with a learning rate of 5e-5. We apply linear learning rate warm up for the first 10% steps of pretraining and linear learning rate decay for the rest. We fine-tune the models on 1,000 samples of the downstream corpora (SST-2, SST-5 and SNLI).

We fine-tune for 3 epochs using batch size of 8. We use AdamW (Loshchilov and Hutter, 2019) optimizer with learning rate of 2e-5 and default parameters:  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 1e-6$ , with bias correction and without weight decay. We average the results over ten random trials.

**Experimental Details for Section A** We conducted second-stage-pretraining by continuing to update BERT-base weights. We used batch size of 32 and learning rate of 5e-5. Learning rate scheduling, optimization and fine-tuning are the same as standard procedure described above.