

PERPLEXED BY PERPLEXITY: PERPLEXITY-BASED DATA PRUNING WITH SMALL REFERENCE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

In this work, we investigate whether small language models can determine high-quality subsets of large-scale text datasets that improve the performance of larger language models. While existing work has shown that pruning based on the perplexity of a larger model can yield high-quality data, we investigate whether smaller models can be used for perplexity-based pruning and how pruning is affected by the domain composition of the data being pruned. We demonstrate that for multiple dataset compositions, perplexity-based pruning of pretraining data can *significantly* improve downstream task performance: pruning based on perplexities computed with a 125 million parameter model improves the average performance on downstream tasks of a 3 billion parameter model by up to 2.04 and achieves up to a $1.45\times$ reduction in pretraining steps to reach commensurate baseline performance. Furthermore, we demonstrate that such perplexity-based data pruning also yields downstream performance gains in the over-trained and data-constrained regimes.

1 INTRODUCTION

A large focus of the machine learning community has been improving the performance of large language models (LLMs) while reducing their training costs. In this work, we consider how to improve the quality of an LLM by improving the quality of its pretraining data. Although there are many techniques to improve data quality, such as augmenting training samples with additional information (Li et al., 2024; Korbak et al., 2023), in this work we focus on the predominant method of *data pruning*: intelligently selecting a high-quality subset of a larger dataset to train on.

Data pruning is commonly used for quality filtering of noisy text data. Simple approaches include using symbolic rules (Bane et al., 2022; Raffel et al., 2020) or using simple classifiers to determine high-quality samples (Wenzek et al., 2020). However, in addition to basic quality filtering, more complex data pruning techniques are also applied to datasets to *further* improve their quality. Xie et al. (2023b) perform importance resampling where importance scores are calculated based on feature similarity to a target text. Tirumala et al. (2023) prune datasets by deduplicating and diversifying data based on a pretrained language model’s embeddings of the text samples. Xie et al. (2023a) re-weight domain proportions based on learnability as determined by a smaller proxy model. Marion et al. (2023) investigate data pruning based on multiple neural heuristics of sample difficulty, ultimately concluding that the perplexity of a sample under a reference language model is the best pruning metric.

In this work, we thoroughly investigate the impact that data pruning based on sample perplexity (Marion et al., 2023) has on LLM pretraining. In particular, we focus on the interplay between pretraining dataset composition and pruning methodology. We further evaluate perplexity pruning in the over-trained and data-constrained regimes. We also investigate whether evaluating the quality of data interventions based on upstream test set perplexity is a sound methodology for gauging downstream performance. To perform perplexity-based data pruning, we train a small language model on a random subset of the given pretraining dataset and then evaluate its perplexity on each sample in the dataset. We then prune the dataset to only include samples within some range of perplexities (i.e., sub-sample to the highest or lowest perplexity samples). We demonstrate that for two vastly different pretraining data compositions, a small language model can be used to effectively prune the

054 pretraining dataset of a significantly larger model, leading to significant gains in the final model’s
055 downstream performance.

056 Our work differs from previous work on perplexity-based data pruning for LLM pretraining in three
057 key ways: (i) our emphasis on downstream model quality evaluation, (ii) our exploration of different
058 pretraining dataset domain compositions, and (iii) our analysis of pruning in non-standard training
059 regimes. While previous works evaluate the resulting LLM’s quality based on upstream metrics such
060 as perplexity on the test split of the pretraining dataset, we evaluate data pruning’s impact based on
061 downstream evaluation benchmarks (e.g. *mmlu* (Hendrycks et al., 2021), *hellaswag*(Zellers et al.,
062 2019), etc.). Evaluating on more meaningful benchmarks enables us to make stronger, more rigorous
063 conclusions about the impact of perplexity-based data pruning, as we find that some techniques
064 that significantly improve downstream performance have no, or even adverse, effects on upstream
065 performance. This difference in metrics enables us to conclude that smaller models can prune
066 the data for larger models, which was not observed in previous perplexity-based pruning works.
067 Secondly, while previous work only investigates pruning on datasets composed of just one domain
068 (CommonCrawl¹), we consider two datasets with different domain compositions: the Pile (Gao et al.,
069 2020) and Dolma (Soldaini et al., 2024). The Pile is composed of many diverse curated domains,
070 with only 15.61% of the data being derived from general web-scrapes, while Dolma is a web-scrape
071 skewed dataset, with 81.31% of its data being derived from the CommonCrawl. We find that
072 successful pruning techniques vary greatly for different dataset compositions to the point that the best
073 technique for one dataset composition may degrade performance for a different composition. Finally,
074 we also evaluate perplexity-based data pruning in the less standard regimes of over-training and
075 data-constrained training. This investigation provides a broader understanding for when practitioners
076 should use perplexity pruning for their data.

077 **Contributions** Our work makes the following contributions:

- 079 • We demonstrate that, across three datasets of varying domain compositions, a small reference
080 model can effectively prune the pretraining dataset of a significantly larger language model
081 ($30\times$ greater parameters), providing both a significant increase in downstream performance
082 and decrease in pretraining steps (Table 1 and Figure 1).
- 083 • We show that data pruning techniques can be highly sensitive to the domain composition
084 of the dataset, suggesting the need to evaluate multiple distinct dataset compositions when
085 conducting data pruning research (Table 1 and Table 4).
- 086 • We investigate perplexity-based data pruning in multiple non-standard settings demonstrating
087 that it can still lead to gains when over-training and when data-constrained (Section 3.4 and
088 Section 3.5).
- 089 • We find that test set perplexity can be a misleading metric for evaluating the efficacy of data
090 pruning techniques, as interventions that result in significantly higher test set perplexity can
091 still achieve better performance on downstream tasks (Table 3).

093 2 PERPLEXITY-BASED DATA PRUNING

094 We start by training a reference model that will be used to calculate the perplexity of all samples in
095 our dataset. First, we partition the original dataset into two splits: one for training the reference model
096 and one for training the final model. After training the reference model on the standard next-token
097 prediction objective, we compute the reference model’s perplexity on each of the samples in the final
098 model’s training split. We then prune the final model’s dataset split to a fraction of its original size,
099 referred to as the *selection rate* (r_s), by selecting samples according to a *selection criteria* which can
100 be one of low, medium, or high. In low selection, samples with the lowest perplexity are selected.
101 In medium selection, we select samples whose perplexity is close to the median perplexity, that is,
102 samples with perplexity in the $[50 - \frac{r_s}{2}, 50 + \frac{r_s}{2}]$ percentiles of all perplexities. In high selection,
103 samples with the highest perplexity are selected. After pruning our dataset, we train a final model
104 using the standard next token prediction objective on the pruned version of the final model training
105 split. We present a pseudocode for pruning based on perplexity in Algorithm 1.

106 ¹<https://data.commoncrawl.org/>

Algorithm 1: Pseudocode for performing perplexity-based data pruning.

Input: Raw dataset $D = \{x^{(i)}\}_{i=1}^M$, where each $x^{(i)}$ is a tokenized text sample;
 selection_criteria $\in \{\text{low, medium, high}\}$; selection rate $r_s \in (0, 1)$; reference training split size R .

Output: Parameters of final model trained on the perplexity pruned dataset θ_{final}^* .

$D_{\text{ref}}, D_{\text{train}} \leftarrow \text{random_split}(D, R)$

$\theta_{\text{ref}} \leftarrow \text{random parameter initialization}$

$\theta_{\text{ref}}^* \leftarrow \text{train}(\theta_{\text{ref}}, D_{\text{ref}})$

$P \leftarrow \{\}$

for $x^{(i)} \in D_{\text{train}}$ **do**

$\text{NLL}_{x^{(i)}} = \frac{1}{|x^{(i)}|} \sum_{t_j \in x^{(i)}} -\log P(t_j | t_{<j}; \theta_{\text{ref}})$

$\text{PPLX}_{x^{(i)}} = 2^{\text{NLL}_{x^{(i)}}}$

$P[x^{(i)}] = \text{PPLX}_{x^{(i)}}$

end

if selection_criteria == "low" **then**

 min_percentile $\leftarrow 0.0$

 max_percentile $\leftarrow r_s$

end

else if selection_criteria == "medium" **then**

 min_percentile $\leftarrow 0.5 - \frac{r_s}{2}$

 max_percentile $\leftarrow 0.5 + \frac{r_s}{2}$

end

else if selection_criteria == "high" **then**

 min_percentile $\leftarrow 1 - r_s$

 max_percentile $\leftarrow 1.0$

end

$\hat{F}_P \leftarrow \text{empirical CDF of } P \text{ values}()$

$D_{\text{pruned}} \leftarrow []$

for $x^{(i)}, \text{PPLX}_{x^{(i)}} \in P$ **do**

if min_percentile $< \hat{F}_P(\text{PPLX}_{x^{(i)}}) < \text{max_percentile}$ **then**

$D_{\text{pruned}} \cdot \text{append}(x^{(i)})$

end

end

$\theta_{\text{final}} \leftarrow \text{random parameter initialization}$

$\theta_{\text{final}}^* \leftarrow \text{train}(\theta_{\text{final}}, D_{\text{pruned}})$

return θ_{final}^*

We consider the setting in which the reference model is significantly smaller than the final model. While this assumption is not strictly necessary, we believe that it is the most practically relevant setup, as it best reflects a data pruning paradigm that would be used for the next generation of LLMs where the models being trained are larger than any existing models.

3 EXPERIMENTS

3.1 SETUP

Models. All models are based on the MPT family of transformer models (Vaswani et al., 2017; MosaicML, 2023c). All reference models have 125 million parameters, and we consider final models with 1 billion and 3 billion parameters.

Data. We consider two datasets in this work. The Pile (Gao et al., 2020) is composed of 22 different domains that range from general web scrapes to legal text. Dolma (Soldaini et al., 2024) is composed of 7 different domains and is derived mainly from general web scrapes. We tokenize all datasets using the GPT-4 tokenizer (OpenAI, 2022).

Table 1: Average normalized accuracy grouped by task category for both datasets and both final model sizes. For all datasets and model sizes we find that training on perplexity pruned data outperforms the baseline. Bold results are within one standard error of the highest score.

Pruning Method	World Knowledge	Common Sense Reasoning	Language Understanding	Symbolic Problem Solving	Reading Comprehension	Average
1B Parameters Trained on Pile						
No Pruning (Baseline)	15.51	10.31	28.11	3.53	11.16	13.73
High Perplexity Selected	18.18	12.75	33.2	3.36	10.63	15.62
3B Parameters Trained on Pile						
No Pruning (Baseline)	21.82	13.09	39.08	4.88	14.28	18.63
High Perplexity Selected	25.8	16.24	43.32	2.91	15.07	20.67
1B Parameters Trained on Dolma						
No Pruning (Baseline)	16.48	12.32	28.86	3.58	7.95	13.84
Medium Perplexity Selected	17.98	13.03	31.87	3.44	10.41	15.35
3B Parameters Trained on Dolma						
No Pruning (Baseline)	23.56	14.29	39.57	4.4	14.2	19.2
Medium Perplexity Selected	24.19	16.48	41.8	3.3	13.19	19.79

Training and hyperparameters. All reference models are trained for a fixed duration of 26 billion tokens. Unless otherwise specified, all final models are trained to Chinchilla optimal (Hoffmann et al., 2022), meaning that each final model’s training duration in tokens is 20 times its parameter count. All models are trained using the decoupled Lion optimizer (Chen et al., 2024) with a cosine learning rate schedule. All reference models and 1B parameter models are trained with a maximum learning rate and weight decay of $2e-4$ and all 3B models are trained with a maximum learning rate and weight decay of $1.6e-4$. Training is conducted using llm-foundry (MosaicML, 2023b) and using both Nvidia A100s and H100s. We perform two trials for each experiment.

Evaluation. We evaluate models on 33 different downstream question-answering tasks using the MosaicML evaluation gauntlet (MosaicML, 2023a). Before averaging the accuracy across tasks, we normalize the accuracy on each task by the baseline of random guessing. Specifically, we normalize the accuracy of each individual task as $a_n = \frac{a_m - a_r}{1 - a_r}$, where a_m is the accuracy of the model and a_r is the expected accuracy of random guessing. We report the average normalized accuracy for each task category as well as the average normalized accuracy across all task categories.² More details on tasks and task categories are listed in Section 8.

3.2 PERPLEXITY-BASED DATA PRUNING IMPROVES DOWNSTREAM PERFORMANCE

If a certain range of perplexities is a good heuristic for data quality, training on that perplexity-pruned subset should improve downstream performance. We sweep across pruning selection criteria and selection rates (Section 7) and find that the best settings are to select high-perplexity samples at a 50% rate for the Pile and to select medium-perplexity samples at a 50% rate for Dolma. We compare the most performant pruning settings to baseline models trained on the original datasets without pruning in Table 1. Across all datasets and model sizes, models pretrained on the perplexity pruned version of the dataset significantly outperform the baseline model on average. Specifically, perplexity-based data pruning outperforms the average downstream performance of no pruning for 1B models by 1.89 and 1.51 for the Pile and Dolma respectively, and improves the performance of 3B models by 2.04 and 0.59 for the Pile and Dolma respectively. These results suggest that the perplexity of a small model provides a strong signal of data quality for a much larger model, as training on the data selected by the small model leads to significant downstream performance improvements.

²Not to be confused, the random accuracy normalization we use is different from the normalized accuracy reported by the EleutherAI LM Evaluation Harness, which normalizes based on the Byte-length of the response.

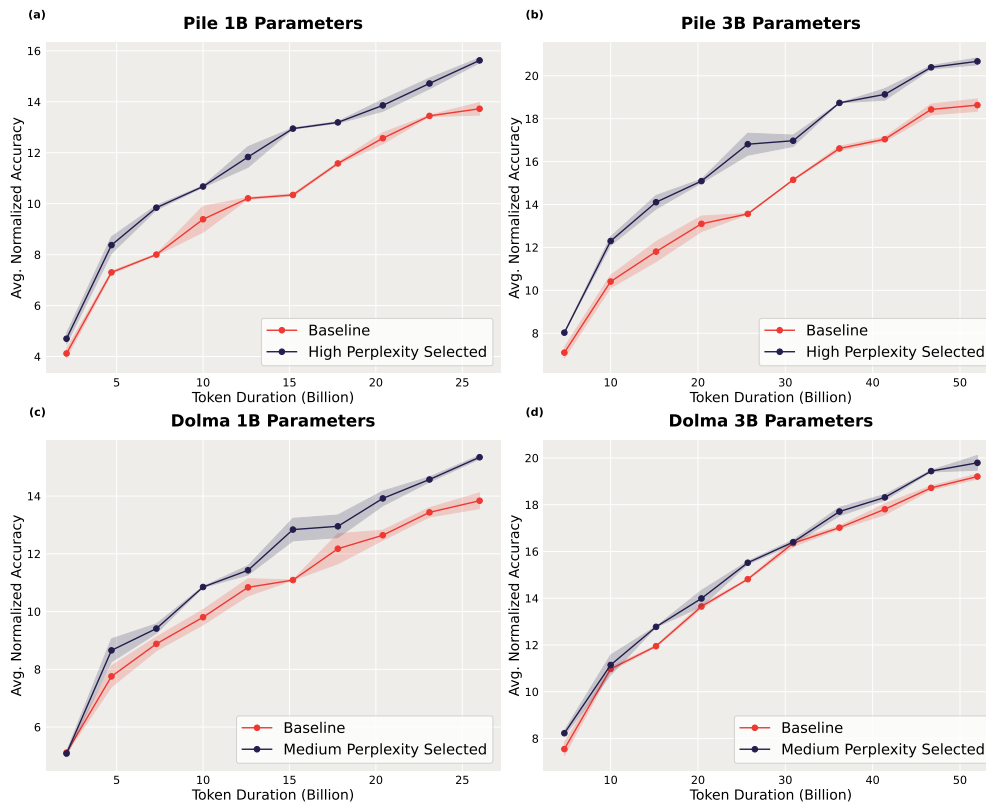


Figure 1: Average normalized task accuracy evaluated intermittently throughout pretraining for each dataset and model size investigated. Perplexity-based data pruning leads to an improvement in performance for all intermediate training steps evaluated.

3.3 PERPLEXITY-BASED DATA PRUNING IMPROVES TRAINING EFFICIENCY

Since perplexity-based data pruning improves the final performance of models, we also investigate how pruned data affects the training dynamics of models. Specifically, we investigate whether training on perplexity pruned data enables models to achieve the same downstream performance as models trained on the unpruned data in training fewer steps. We plot the average downstream performance of partially trained checkpoints from the 1B baseline and perplexity pruned models in Figure 1. Perplexity pruning outperforms the baseline model for all intermediate pretraining durations evaluated. Furthermore, perplexity pruned models reach the same average normalized accuracy as the baseline models in $1.31\times$ and $1.45\times$ fewer steps for Pile 1B and 3B respectively and in $1.29\times$ and $1.14\times$ fewer steps for Dolma 1B and Dolma 3B respectively. These results demonstrate that the resulting high-quality data from perplexity-based data pruning enables faster learning which can be leveraged to achieve the same downstream performance as training on unpruned data with fewer pretraining steps.

3.4 PERPLEXITY-BASED DATA PRUNING FOR OVER-TRAINED MODELS

A recent trend with LLMs has been to over-train models by training them on more tokens than the Chinchilla optimal number of tokens (Touvron et al., 2023; Gadre et al., 2024). As our work targets the data component of LLM pretraining, we investigate the hypothesis that over-training would be more beneficial for models trained on perplexity pruned datasets as the data is of higher quality. We test this hypothesis by training a 1B parameter model for 130B tokens, which is $5\times$ the Chinchilla optimal number of tokens. We evaluate the downstream performance of each over-trained model in Table 2. The main observation is that while the absolute gain in average downstream normalized accuracy from perplexity-based data pruning on the Pile is similar for both compute optimal and

Table 2: Downstream task performance for Chinchilla Optimal and $5\times$ over-trained data budgets. The “Improvement Over Baseline” column refers to the gain observed from perplexity pruning as compared to the baseline trained in the same setting.

Pruning Method	Average	Improvement Over Baseline
1B Parameters Trained on High Perplexity Pile		
Chinchilla Optimal	15.62	1.89
$5\times$ Over-Trained	18.83	1.74
1B Parameters Trained on Medium Perplexity Dolma		
Chinchilla Optimal	15.35	1.51
$5\times$ Over-Trained	18.67	0.84

over-trained models, the gain decreases for Dolma when over-training. On the Pile, we find that the gain from perplexity pruned data is similar in the compute optimal regime and the over-trained regime: we see a gain in average performance of 1.89 when training compute optimal and a gain of 1.74 when over-training. On Dolma, the gain from perplexity pruned data decreases in the over-trained regime: we see a gain of 1.51 when training for a compute optimal duration but this decreases to a gain of 0.84 when over-training. These results show that while the higher quality data resulting from perplexity-based data pruning does still lead to an improvement in downstream performance in the over-trained regime, there is not a relative increase in downstream improvement over the baseline when over-training.

3.5 PERPLEXITY-BASED DATA PRUNING FOR THE DATA CONSTRAINED REGIME

Our experiments so far were conducted in the setting where there exists a sufficient abundance of data such that even after pruning with the desired selection rate there are enough data points to fill the desired token budget without requiring any data to be repeated. However, there are many training settings that do not fall under this data-abundant regime. Consequently, we evaluate how perplexity-based data pruning performs when the number of tokens is constrained, and pruning induces a greater number of repetitions of the data. For each dataset, we vary the available data such that training for a Chinchilla optimal number of tokens requires a different number of repetitions. Specifically, we investigate data budgets that require $\{0.5, 1, 2, 4, 8\}$ repetitions to reach the Chinchilla optimal number of tokens³. As each number of repeats refers to the total number of tokens available, for all pruning experiments the number of repetitions after pruning is actually greater by a factor of $\frac{1}{r_s}$ since we prune the available tokens according to r_s , the selection rate. Since all models use a selection rate of 0.5, the models trained on the pruned data see the data for $2\times$ more repetitions.

We plot the average downstream performance as a function of the number of repetitions in Figure 2. On both the Pile and Dolma, we find that training on perplexity pruned data yields an improvement for up to two repetitions. These results suggest that perplexity-based data pruning can still provide performance gains for some degree of data constraint. Furthermore, our results replicate the findings of Muennighoff et al. (2023) that more than four repetitions yields negligible gains. Specifically, the baseline model without pruning maintains commensurate performance for up to four repetitions. Similarly, models trained on perplexity-pruned data maintain commensurate performance for up to two repetitions through the base data, which corresponds to four repetitions after pruning. That training on repeated perplexity-pruned data leads to diminishing gains after four repetitions post-pruning suggests that the higher quality data resulting from pruning does not change the point for which repeating data yields diminishing improvements in performance.

3.6 UPSTREAM PERPLEXITY IS NOT A RELIABLE EVALUATION METRIC FOR DATA PRUNING

As previous works have used the perplexity of the model on a test split of the pretraining dataset as an approximation to downstream performance, we wanted to explore how well such perplexity-based

³Repeat=0.5 means that the available number of tokens is twice the training budget, i.e. the data-abundant setting

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
367
368
369
370
371
372
373
374
375
376
377

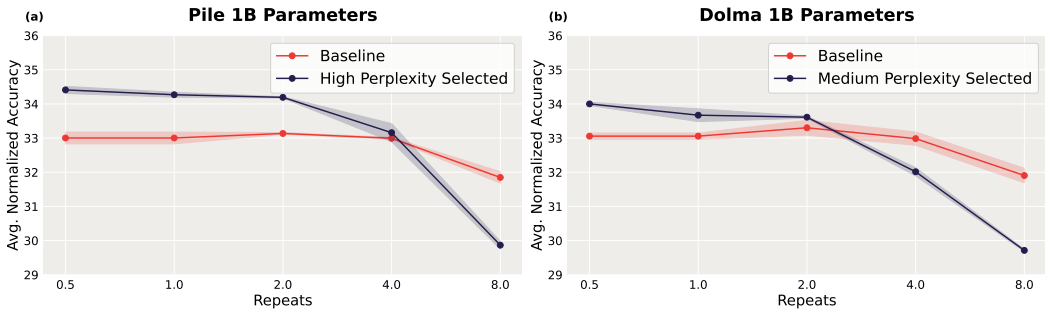


Figure 2: Downstream task performance as a function of available dataset size. The number of repeats denotes the number of repeats over the raw dataset necessary to achieve the Chinchilla optimal number of tokens. Training on perplexity pruned data leads to an improvement for up to two repeats on both the Pile Dolma.

Table 3: Performance as evaluated by perplexity on a test split of the original dataset as well as average normalized task accuracy for 1 billion parameter final models trained on the Pile. The model trained on pruned data has worse pretraining test split perplexity even though it significantly improves average downstream normalized accuracy.

Pruning Method	Test Set Pplx. (\downarrow)	Downstream Task Avg. (\uparrow)
1B Parameters Trained on Pile		
No Pruning (Baseline)	7.83	13.73
High Perplexity Selected	8.51	15.62
1B Parameters Trained on Dolma		
No Pruning (Baseline)	13.53	13.84
Medium Perplexity Selected	14.33	15.35

evaluations agree with downstream performance for data intervention techniques. Pruning performs an intervention on the dataset, making models trained on the pruned dataset biased estimators of the original data distribution. Therefore, it is unlikely that the performance on the original data distribution is a fair evaluation of model quality. We compare the test set perplexity and average downstream performance for 1 billion parameter models trained on the original and pruned version of the Pile and Dolma in Table 3. For both the Pile and Dolma, training on perplexity pruned data significantly worsens perplexity on a test split of the pretraining data, while the average downstream performance is significantly improved. This result suggests that test set perplexity may not always be a sound metric for data pruning work and that researchers should instead directly evaluate on downstream benchmarks.

4 UNDERSTANDING THE EFFECTS OF PERPLEXITY-BASED PRUNING

In this section, we investigate how data pruning works by exploring some of the properties of perplexity-based pruning.

4.1 HOW ARE REFERENCE PERPLEXITIES DISTRIBUTED

In order to better understand how perplexity-based data pruning works, we investigate the distribution of the computed reference model perplexities for each dataset. For each dataset, we randomly sample 10% of the calculated perplexities and perform kernel density estimation to estimate the distribution of log perplexities for a given dataset. We repeat this procedure for the optimal pruned version of the dataset. We plot the resulting estimates of the log perplexity distribution in Figure 3. We find that the log perplexity distribution for the Pile is multimodal and asymmetric, while for Dolma and it is unimodal and symmetric.

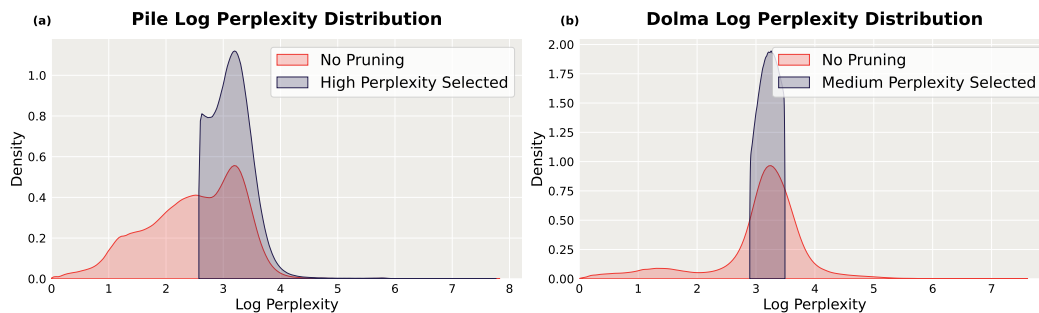


Figure 3: Distribution of sample perplexities as evaluated by the reference model for the Pile and Dolma. We show both the original distribution over the full dataset without pruning as well as the distribution after applying the optimal perplexity-based data pruning technique for a given dataset.

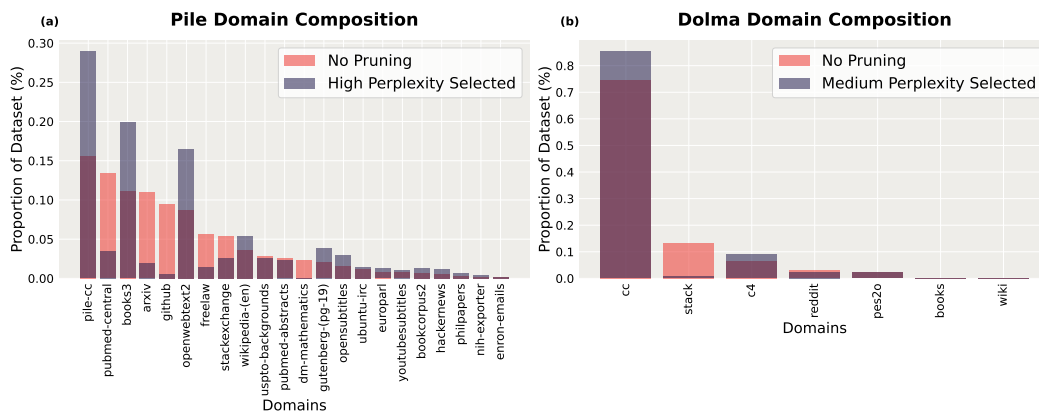


Figure 4: Proportion of the total dataset each domain makes up before and after pruning. For all datasets, pruning tends to select more samples from general web domains while leaving out samples from highly specific domains.

4.2 HOW PRUNING AFFECTS DOMAIN COMPOSITION

We can also interpret the effect that perplexity-based data pruning has on a dataset by examining how pruning affects each domain’s proportion of the total dataset. We plot the pre and post-pruning domain compositions for the Pile and Dolma in Figure 4. Interestingly, for all datasets pruning increases the proportion of data coming from web-scraped domains while decreasing the proportion of data coming from highly specific technical domains such as code or scientific papers. This trend is more pronounced in the Pile, where the proportions of Pile-CC and OpenWebText2 nearly double, while the proportions of domains such as Pubmed Central, ArXiv, and Github are all reduced by at least a factor of three. Future work should investigate how perplexity-based pruning affects a model’s performance on downstream tasks that are in the same category as the highly pruned domains.

5 RELATED WORK

Classical methods for pruning text data. In order to improve the quality of raw web scrapes, which often contain very noisy samples, pruning via quality filtering has become a common practice. Simple rules-based methods have been employed to prune datasets by filtering out low-quality samples according to some hand-crafted heuristic such as whether the text contains prohibited words, is predominantly English, etc. (Bane et al., 2022; Raffel et al., 2020; Rae et al., 2022; Penedo et al., 2023). N-gram perplexity-based methods, in which an n-gram model is first trained on a high quality, curated corpus and then used to score another corpus, have also been applied to filter text

432 data (Moore & Lewis, 2010; Axelrod, 2017; Gao, 2021; Laurençon et al., 2022; Muennighoff et al.,
433 2023). Although our method also uses perplexity to prune data, it does so in a very different manner.
434 In n-gram perplexity pruning, perplexity is used to estimate whether new text is in distribution as
435 compared to the curated text the n-gram was trained on, while in our model-based perplexity pruning,
436 the reference model is trained on the same distribution of text and the perplexity is more akin to an
437 estimate of the difficulty of an example. In this work, the datasets we leverage already have some
438 basic rules-based pruning applied, and as such, the method we investigate is largely complementary
439 to these existing techniques.

440
441 **Neural network based methods for pruning text data.** Recently, there has been much interest
442 in using neural networks to compute metrics that can be used to intelligently prune datasets. A
443 common technique in this family of methods is using a model to sample high-quality data from
444 large datasets based on the sample’s similarity to a curated high-quality corpus that serves as a target
445 distribution (Feng et al., 2022; Xie et al., 2023b). Xie et al. (2023a) also consider how to use a
446 small reference model to prune pretraining data for a much larger model, by using a small reference
447 model to learn the optimal weighting of domain proportions to maximize the "learnability" of the
448 resulting dataset. Pruning based on the difficulty or loss of a sample has previously been explored
449 for text data, but the majority of such work focuses on curating data for finetuning (Swayamdipta
450 et al., 2020; Attenu & Corbeil, 2023; Coleman et al., 2020; Mindermann et al., 2022; Mekala et al.,
451 2024). Marion et al. (2023), however, investigate multiple model-based sample difficulty heuristics
452 for pruning pretraining text datasets. Although we use the same method for pruning text pretraining
453 datasets, our analysis differs substantially as we evaluate model quality based on downstream metrics
454 and extend our analysis to multiple different dataset compositions which enables us to conclude that
455 the reference model can be smaller than the final model.

456
457 **Data pruning on vision tasks.** While data pruning is becoming more and more relevant with large
458 amounts of text data, it has also been extensively applied in the vision domain (Paul et al., 2021;
459 Toneva et al., 2018; Park et al., 2023). These works often prune data points based on their loss or
460 gradients during training (Killamsetty et al., 2021; Mirzasoleiman et al., 2020). Model-based methods
461 have also been leveraged for image data pruning (Fang et al., 2024; Schuhmann et al., 2021). Note
462 that in the literature, data pruning is also sometimes referred to as coreset selection (Guo et al., 2022).
463 More recently, Park et al. (2022) show that, somewhat surprisingly, active learning (Castro & Nowak,
464 2008) based algorithms tend to outperform most data subset selection algorithms. In the context of
465 contrastive learning, hard-negative mining has been effective as a data pruning method (Kalantidis
466 et al., 2020; Robinson et al., 2020; Zhang & Stratos, 2021). Recently, Goyal et al. (2024) investigated
467 scaling laws for training on pruned data in the context of vision models.

468 469 6 CONCLUSION

470
471
472 In this work, we conduct an empirical investigation of the impact that perplexity-based data pruning
473 has on model performance. We demonstrate that small reference models can be used to prune the
474 data of models with up to $30\times$ more parameters, leading to both significant downstream performance
475 improvements and increased training efficiency. We then investigate perplexity-based data pruning
476 in two non-standard settings: the over-trained and data-constrained regimes. We find that for both
477 settings, training on perplexity pruned data can outperform training on unpruned data, demonstrating
478 that perplexity-based data pruning is a widely applicable and extensible technique. We also investigate
479 upstream metrics for evaluating data pruning techniques and provide an example where evaluating
480 models based on their perplexity on the test split of the pretraining dataset does not align with
481 evaluating based on downstream model performance. Additionally, we demonstrate that optimal
482 pruning techniques can vary greatly for different dataset compositions. Although we do not present a
483 predictive theory for how pruning parameters should be selected for different datasets, we demonstrate
484 that the optimal pruning parameters for a 1 billion parameter model can successfully transfer to
485 3 billion parameter models, potentially suggesting that empirically determining the optimal pruning
parameters can be done cheaply. Our work takes a key step towards establishing perplexity-based
data pruning as a primary technique in the modern data researcher’s toolkit.

REFERENCES

- 486
487
488 Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh
489 Hajishirzi. MathQA: Towards interpretable math word problem solving with operation-based
490 formalisms. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019*
491 *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 2357–2367, Minneapolis,
492 Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1245.
493 URL <https://aclanthology.org/N19-1245>.
494
- 495 Jean-michel Attendu and Jean-philippe Corbeil. NLU on data diets: Dynamic data subset selection
496 for NLP classification tasks. In Nafise Sadat Moosavi, Iryna Gurevych, Yufang Hou, Gyuwan Kim,
497 Young Jin Kim, Tal Schuster, and Ameeta Agrawal (eds.), *Proceedings of The Fourth Workshop on*
498 *Simple and Efficient Natural Language Processing (SustainNLP)*, pp. 129–146, Toronto, Canada
499 (Hybrid), July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.sustainlp-1.
500 9. URL <https://aclanthology.org/2023.sustainlp-1.9>.
- 501 Amittai Axelrod. Cynical selection of language model training data. *arXiv preprint arXiv:1709.02279*,
502 2017.
503
- 504 Fred Bane, Celia Soler Uguet, Wiktor Stribizew, and Anna Zaretskaya. A comparison of data filtering
505 methods for neural machine translation. In Janice Campbell, Stephen Larocca, Jay Marciano,
506 Konstantin Savenkov, and Alex Yanishevsky (eds.), *Proceedings of the 15th Biennial Conference*
507 *of the Association for Machine Translation in the Americas (Volume 2: Users and Providers Track*
508 *and Government Track)*, pp. 313–325, Orlando, USA, September 2022. Association for Machine
509 Translation in the Americas. URL <https://aclanthology.org/2022.amta-upg.22>.
- 510 Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical
511 commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*,
512 volume 34, pp. 7432–7439, 2020.
- 513 Rui Castro and Robert Nowak. Active learning and sampling. In *Foundations and Applications of*
514 *Sensor Management*, pp. 177–200. Springer, 2008.
515
- 516 Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Hieu Pham, Xuanyi Dong,
517 Thang Luong, Cho-Jui Hsieh, Yifeng Lu, et al. Symbolic discovery of optimization algorithms.
518 *Advances in Neural Information Processing Systems*, 36, 2024.
- 519 Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina
520 Toutanova. BoolQ: Exploring the surprising difficulty of natural yes/no questions. In Jill
521 Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of*
522 *the North American Chapter of the Association for Computational Linguistics: Human Lan-*
523 *guage Technologies, Volume 1 (Long and Short Papers)*, pp. 2924–2936, Minneapolis, Min-
524 nesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1300. URL
525 <https://aclanthology.org/N19-1300>.
- 526 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and
527 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.
528 *arXiv:1803.05457v1*, 2018.
529
- 530 Cody Coleman, Christopher Yeh, Stephen Mussmann, Baharan Mirzasoleiman, Peter Bailis, Percy
531 Liang, Jure Leskovec, and Matei Zaharia. Selection via proxy: Efficient data selection for
532 deep learning. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=HJg2b0VYDr>.
533
- 534 Alex Fang, Albin Madappally Jose, Amit Jain, Ludwig Schmidt, Alexander T Toshev, and Vaishaal
535 Shankar. Data filtering networks. In *The Twelfth International Conference on Learning Represen-*
536 *tations*, 2024. URL <https://openreview.net/forum?id=KAk6ngZ09F>.
537
- 538 Yukun Feng, Patrick Xia, Benjamin Van Durme, and João Sedoc. Automatic document selection
539 for efficient encoder pretraining. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.),
Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pp.

- 540 9522–9530, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational
541 Linguistics. doi: 10.18653/v1/2022.emnlp-main.647. URL [https://aclanthology.org/
542 2022.emnlp-main.647](https://aclanthology.org/2022.emnlp-main.647).
- 543 Samir Yitzhak Gadre, Georgios Smyrnis, Vaishaal Shankar, Suchin Gururangan, Mitchell Wortsman,
544 Rulin Shao, Jean Mercat, Alex Fang, Jeffrey Li, Sedrick Keh, et al. Language models scale reliably
545 with over-training and on downstream tasks. *arXiv preprint arXiv:2403.08540*, 2024.
546
- 547 Leo Gao. An empirical exploration in quality filtering of text data, 2021.
- 548 Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang,
549 Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. The pile: An 800gb
550 dataset of diverse text for language modeling, 2020.
- 551 Sachin Goyal, Pratyush Maini, Zachary C Lipton, Aditi Raghunathan, and J Zico Kolter. Scaling
552 laws for data filtering—data curation cannot be compute agnostic. *arXiv preprint arXiv:2404.07177*,
553 2024.
554
- 555 Chengcheng Guo, Bo Zhao, and Yanbing Bai. Deepcore: A comprehensive library for coreset selec-
556 tion in deep learning. In *International Conference on Database and Expert Systems Applications*,
557 pp. 181–195. Springer, 2022.
- 558 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob
559 Steinhardt. Measuring massive multitask language understanding. In *International Confer-
560 ence on Learning Representations*, 2021. URL [https://openreview.net/forum?id=
561 d7KBjmI3GmQ](https://openreview.net/forum?id=d7KBjmI3GmQ).
- 562 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
563 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Thomas
564 Hennigan, Eric Noland, Katherine Millican, George van den Driessche, Bogdan Damoc, Au-
565 relia Guy, Simon Osindero, Karén Simonyan, Erich Elsen, Oriol Vinyals, Jack Rae, and
566 Laurent Sifre. An empirical analysis of compute-optimal large language model training.
567 In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Ad-
568 vances in Neural Information Processing Systems*, volume 35, pp. 30016–30030. Curran Asso-
569 ciates, Inc., 2022. URL [https://proceedings.neurips.cc/paper_files/paper/
570 2022/file/c1e2faff6f588870935f114ebe04a3e5-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/c1e2faff6f588870935f114ebe04a3e5-Paper-Conference.pdf).
- 571 Aditi Jha, Sam Havens, Jeremy Dohmann, Alexander Trott, and Jacob Portes. LIMIT: Less is more
572 for instruction tuning across evaluation paradigms. In *NeurIPS 2023 Workshop on Instruction
573 Tuning and Instruction Following*, 2023. URL [https://openreview.net/forum?id=
574 QxtL4Q1enz](https://openreview.net/forum?id=QxtL4Q1enz).
- 575 Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. PubMedQA: A dataset
576 for biomedical research question answering. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun
577 Wan (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Language
578 Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-
579 IJCNLP)*, pp. 2567–2577, Hong Kong, China, November 2019. Association for Computational
580 Linguistics. doi: 10.18653/v1/D19-1259. URL <https://aclanthology.org/D19-1259>.
- 581 Yannis Kalantidis, Mert Bulent Sariyildiz, Noe Pion, Philippe Weinzaepfel, and Diane Larlus. Hard
582 negative mixing for contrastive learning. *Advances in Neural Information Processing Systems*, 33:
583 21798–21809, 2020.
584
- 585 Krishnateja Killamsetty, Sivasubramanian Durga, Ganesh Ramakrishnan, Abir De, and Rishabh Iyer.
586 Grad-match: Gradient matching based data subset selection for efficient deep model training. In
587 *International Conference on Machine Learning*, pp. 5464–5474. PMLR, 2021.
- 588 Tomasz Korbak, Kejian Shi, Angelica Chen, Rasika Vinayak Bhalerao, Christopher L. Buckley,
589 Jason Phang, Samuel R. Bowman, and Ethan Perez. Pretraining language models with human
590 preferences. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan
591 Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023,
592 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning
593 Research*, pp. 17506–17533. PMLR, 2023. URL [https://proceedings.mlr.press/
v202/korbak23a.html](https://proceedings.mlr.press/v202/korbak23a.html).

- 594 Hugo Laurençon, Lucile Saulnier, Thomas Wang, Christopher Akiki, Albert Villanova del Moral,
595 Teven Le Scao, Leandro Von Werra, Chenghao Mou, Eduardo González Ponferrada, Huu Nguyen,
596 Jörg Froberg, Mario Šaško, Quentin Lhoest, Angelina McMillan-Major, Gérard Dupont, Stella
597 Biderman, Anna Rogers, Loubna Ben allal, Francesco De Toni, Giada Pistilli, Olivier Nguyen,
598 Somaieh Nikpoor, Maraim Masoud, Pierre Colombo, Javier de la Rosa, Paulo Villegas, Tristan
599 Thrush, Shayne Longpre, Sebastian Nagel, Leon Weber, Manuel Romero Muñoz, Jian Zhu,
600 Daniel Van Strien, Zaid Alyafeai, Khalid Almubarak, Vu Minh Chien, Itziar Gonzalez-Dios, Aitor
601 Soroa, Kyle Lo, Manan Dey, Pedro Ortiz Suarez, Aaron Gokaslan, Shamik Bose, David Ifeoluwa
602 Adelani, Long Phan, Hieu Tran, Ian Yu, Suhas Pai, Jenny Chim, Violette Lepercq, Suzana Ilic,
603 Margaret Mitchell, Sasha Luccioni, and Yacine Jernite. The bigscience ROOTS corpus: A 1.6TB
604 composite multilingual dataset. In *Thirty-sixth Conference on Neural Information Processing
605 Systems Datasets and Benchmarks Track*, 2022. URL [https://openreview.net/forum?
606 id=UoEw6KigkUn](https://openreview.net/forum?id=UoEw6KigkUn).
- 607 Hector Levesque, Ernest Davis, and Leora Morgenstern. The winograd schema challenge. In
608 *Thirteenth International Conference on the Principles of Knowledge Representation and Reasoning*.
609 Citeseer, 2012.
- 610 Xian Li, Ping Yu, Chunting Zhou, Timo Schick, Luke Zettlemoyer, Omer Levy, Jason Weston,
611 and Mike Lewis. Self-alignment with instruction backtranslation. In *The Twelfth International
612 Conference on Learning Representations*, 2024. URL [https://openreview.net/forum?
613 id=loiJHJBRsT](https://openreview.net/forum?id=loiJHJBRsT).
- 614 Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. Logiqa: A
615 challenge dataset for machine reading comprehension with logical reasoning. In Christian Bessiere
616 (ed.), *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence,
617 IJCAI 2020*, pp. 3622–3628. ijcai.org, 2020. doi: 10.24963/IJCAI.2020/501. URL [https:
618 //doi.org/10.24963/ijcai.2020/501](https://doi.org/10.24963/ijcai.2020/501).
- 619 Max Marion, Ahmet Üstün, Luiza Pozzobon, Alex Wang, Marzieh Fadaee, and Sara Hooker. When
620 less is more: Investigating data pruning for pretraining LLMs at scale. In *NeurIPS Workshop on
621 Attributing Model Behavior at Scale*, 2023. URL [https://openreview.net/forum?id=
622 XUIYn3jo5T](https://openreview.net/forum?id=XUIYn3jo5T).
- 623 Dheeraj Mekala, Alex Nguyen, and Jingbo Shang. Smaller language models are capable of selecting
624 instruction-tuning training data for larger language models. *arXiv preprint arXiv:2402.10430*,
625 2024.
- 626 Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct
627 electricity? a new dataset for open book question answering. In *Conference on Empirical Methods
628 in Natural Language Processing*, 2018. URL [https://api.semanticscholar.org/
629 CorpusID:52183757](https://api.semanticscholar.org/CorpusID:52183757).
- 630 Sören Mindermann, Jan M Brauner, Muhammed T Razzak, Mrinank Sharma, Andreas Kirsch,
631 Winnie Xu, Benedikt Hölting, Aidan N Gomez, Adrien Morisot, Sebastian Farquhar, and Yarin
632 Gal. Prioritized training on points that are learnable, worth learning, and not yet learnt. In
633 Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato
634 (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of
635 *Proceedings of Machine Learning Research*, pp. 15630–15649. PMLR, 17–23 Jul 2022. URL
636 <https://proceedings.mlr.press/v162/mindermann22a.html>.
- 637 Baharan Mirzasoleiman, Jeff Bilmes, and Jure Leskovec. Coresets for data-efficient training of
638 machine learning models. In *International Conference on Machine Learning*, pp. 6950–6960.
639 PMLR, 2020.
- 640 Robert C. Moore and William Lewis. Intelligent selection of language model training data. In Jan
641 Hajič, Sandra Carberry, Stephen Clark, and Joakim Nivre (eds.), *Proceedings of the ACL 2010 Con-
642 ference Short Papers*, pp. 220–224, Uppsala, Sweden, July 2010. Association for Computational
643 Linguistics. URL <https://aclanthology.org/P10-2041>.
- 644 MosaicML. Llm evaluation scores, 2023a. URL [https://www.mosaicml.com/
645 llm-evaluation](https://www.mosaicml.com/llm-evaluation).

- 648 MosaicML. Llm foundry. <https://github.com/mosaicml/llm-foundry>, 2023b.
- 649
- 650 MosaicML. Introducing mpt-7b: A new standard for open-source, commercially usable llms, 2023c.
- 651 URL <https://www.databricks.com/blog/mpt-7b>.
- 652 Niklas Muennighoff, Alexander M Rush, Boaz Barak, Teven Le Scao, Nouamane Tazi, Aleksandra
- 653 Piktus, Sampo Pyysalo, Thomas Wolf, and Colin Raffel. Scaling data-constrained language
- 654 models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL
- 655 <https://openreview.net/forum?id=j5BuTrEj35>.
- 656
- 657 OpenAI. Tiktoken: A fast bpe tokeniser for use with openai’s models. [https://github.com/](https://github.com/openai/tiktoken)
- 658 [openai/tiktoken](https://github.com/openai/tiktoken), 2022.
- 659 Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Ngoc Quan Pham, Raffaella Bernardi,
- 660 Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. The LAMBADA dataset:
- 661 Word prediction requiring a broad discourse context. In Katrin Erk and Noah A. Smith (eds.),
- 662 *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume*
- 663 *1: Long Papers)*, pp. 1525–1534, Berlin, Germany, August 2016. Association for Computational
- 664 Linguistics. doi: 10.18653/v1/P16-1144. URL <https://aclanthology.org/P16-1144>.
- 665 Dongmin Park, Dimitris Papailiopoulos, and Kangwook Lee. Active learning is a strong baseline for
- 666 data subset selection. In *Has it Trained Yet? NeurIPS 2022 Workshop*, 2022.
- 667
- 668 Dongmin Park, Seola Choi, Doyoung Kim, Hwanjun Song, and Jae-Gil Lee. Robust data pruning
- 669 under label noise via maximizing re-labeling accuracy. *arXiv preprint arXiv:2311.01002*, 2023.
- 670 Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. Deep learning on a data diet: Finding
- 671 important examples early in training. *Advances in Neural Information Processing Systems*, 34:
- 672 20596–20607, 2021.
- 673
- 674 Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Hamza Alobeidli,
- 675 Alessandro Cappelli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. The refinedweb
- 676 dataset for falcon LLM: Outperforming curated corpora with web data only. In *Thirty-seventh*
- 677 *Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023.
- 678 URL <https://openreview.net/forum?id=kM5eGcdCzq>.
- 679 Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John
- 680 Aslanides, Sarah Henderson, Roman Ring, Susannah Young, Eliza Rutherford, Tom Hennigan,
- 681 Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche, Lisa Anne Hendricks,
- 682 Maribeth Rauh, Po-Sen Huang, Amelia Glaese, Johannes Welbl, Sumanth Dathathri, Saffron
- 683 Huang, Jonathan Uesato, John Mellor, Irina Higgins, Antonia Creswell, Nat McAleese, Amy Wu,
- 684 Erich Elsen, Siddhant Jayakumar, Elena Buchatskaya, David Budden, Esme Sutherland, Karen
- 685 Simonyan, Michela Paganini, Laurent Sifre, Lena Martens, Xiang Lorraine Li, Adhiguna Kuncoro,
- 686 Aida Nematzadeh, Elena Gribovskaya, Domenic Donato, Angeliki Lazaridou, Arthur Mensch,
- 687 Jean-Baptiste Lespiau, Maria Tsimpoukelli, Nikolai Grigorev, Doug Fritz, Thibault Sottiaux,
- 688 Mantas Pajarskas, Toby Pohlen, Zhitao Gong, Daniel Toyama, Cyprien de Masson d’Autume,
- 689 Yujia Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin, Aidan Clark, Diego de Las Casas,
- 690 Aurelia Guy, Chris Jones, James Bradbury, Matthew Johnson, Blake Hechtman, Laura Weidinger,
- 691 Iason Gabriel, William Isaac, Ed Lockhart, Simon Osindero, Laura Rimell, Chris Dyer, Oriol
- 692 Vinyals, Kareem Ayoub, Jeff Stanway, Lorraine Bennett, Demis Hassabis, Koray Kavukcuoglu,
- 693 and Geoffrey Irving. Scaling language models: Methods, analysis and insights from training
- 694 gopher, 2022.
- 695
- 696 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,
- 697 Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified
- 698 text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020. URL
- 699 <http://jmlr.org/papers/v21/20-074.html>.
- 700
- 701 Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions
- for machine comprehension of text. In Jian Su, Kevin Duh, and Xavier Carreras (eds.), *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 2383–2392, Austin, Texas, November 2016. Association for Computational Linguistics. doi: 10.18653/v1/D16-1264. URL <https://aclanthology.org/D16-1264>.

- 702 Joshua Robinson, Ching-Yao Chuang, Suvrit Sra, and Stefanie Jegelka. Contrastive learning with
703 hard negative samples. *arXiv preprint arXiv:2010.04592*, 2020.
704
- 705 Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S. Gordon. Choice of plausible alternatives:
706 An evaluation of commonsense causal reasoning. In *Logical Formalizations of Commonsense*
707 *Reasoning, Papers from the 2011 AAAI Spring Symposium, Technical Report SS-11-06, Stanford,*
708 *California, USA, March 21-23, 2011*. AAAI, 2011. URL [http://www.aaai.org/ocs/](http://www.aaai.org/ocs/index.php/SSS/SSS11/paper/view/2418)
709 [index.php/SSS/SSS11/paper/view/2418](http://www.aaai.org/ocs/index.php/SSS/SSS11/paper/view/2418).
- 710 Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An
711 adversarial winograd schema challenge at scale. In *AAAI*, pp. 8732–8740, 2020. URL [https://](https://aaai.org/ojs/index.php/AAAI/article/view/6399)
712 aaai.org/ojs/index.php/AAAI/article/view/6399.
713
- 714 Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis,
715 Aarush Katta, Theo Coombs, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of
716 clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021.
- 717 Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur,
718 Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, Valentin Hofmann, Ananya Harsh
719 Jha, Sachin Kumar, Li Lucy, Xinxu Lyu, Nathan Lambert, Ian Magnusson, Jacob Morrison, Niklas
720 Muennighoff, Aakanksha Naik, Crystal Nam, Matthew E. Peters, Abhilasha Ravichander, Kyle
721 Richardson, Zejiang Shen, Emma Strubell, Nishant Subramani, Oyvind Tafjord, Pete Walsh, Luke
722 Zettlemoyer, Noah A. Smith, Hannaneh Hajishirzi, Iz Beltagy, Dirk Groeneveld, Jesse Dodge,
723 and Kyle Lo. Dolma: an Open Corpus of Three Trillion Tokens for Language Model Pretraining
724 Research. *arXiv preprint*, 2024.
- 725 Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam
726 Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska,
727 Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W.
728 Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda
729 Askeel, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders Johan
730 Andreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew M. Dai, Andrew La,
731 Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna
732 Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul Menezes,
733 Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut
734 Erdem, Ayla Karakaş, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph, Bartłomiej Bojanowski,
735 Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno Stein, Berk
736 Ekmekci, Bill Yuchen Lin, Blake Howald, Bryan Orinion, Cameron Diao, Cameron Dour, Catherine
737 Stinson, Cedrick Argueta, Cesar Ferri, Chandan Singh, Charles Rathkopf, Chenlin Meng, Chitta
738 Baral, Chiyu Wu, Chris Callison-Burch, Christopher Waites, Christian Voigt, Christopher D
739 Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Siro, Colin Raffel,
740 Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garrette, Dan Hendrycks, Dan Kilman,
741 Dan Roth, C. Daniel Freeman, Daniel Khashabi, Daniel Levy, Daniel Moseguí González, Danielle
742 Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa, David Dohan, David
743 Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin, Denis Kleyko, Deniz
744 Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta Misra, Dilyar Buzan, Dimitri Coelho
745 Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Ekaterina Shutova, Ekin Dogus Cubuk, Elad
746 Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, Ellie Pavlick, Emanuele Rodolà,
747 Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie Chang, Ethan A Chi, Ethan Dyer, Ethan
748 Jerzak, Ethan Kim, Eunice Engefu Manyasi, Evgenii Zheltonozhskii, Fanyue Xia, Fatemeh Siar,
749 Fernando Martínez-Plumed, Francesca Happé, Francois Chollet, Frieda Rong, Gaurav Mishra,
750 Genta Indra Winata, Gerard de Melo, Germán Kruszewski, Giambattista Parascandolo, Giorgio
751 Mariani, Gloria Xinyue Wang, Gonzalo Jaimovitch-Lopez, Gregor Betz, Guy Gur-Ari, Hana
752 Galijasevic, Hannah Kim, Hannah Rashkin, Hannaneh Hajishirzi, Harsh Mehta, Hayden Bogar,
753 Henry Francis Anthony Shevlin, Hinrich Schuetze, Hiromu Yakura, Hongming Zhang, Hugh Mee
754 Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jack Geissinger, Jackson Kernion, Jacob Hilton, Jaehoon
755 Lee, Jaime Fernández Fisac, James B Simon, James Koppel, James Zheng, James Zou, Jan Kocon,
756 Jana Thompson, Janelle Wingfield, Jared Kaplan, Jarema Radom, Jascha Sohl-Dickstein, Jason
757 Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova, Jelle Bosscher, Jennifer Marsh, Jeremy
758 Kim, Jeroen Taal, Jesse Engel, Jesujoba Alabi, Jiacheng Xu, Jiaming Song, Jillian Tang, Joan

- 756 Waweru, John Burden, John Miller, John U. Balis, Jonathan Batchelder, Jonathan Berant, Jörg
757 Frohberg, Jos Rozen, Jose Hernandez-Orallo, Joseph Boudeman, Joseph Guerr, Joseph Jones,
758 Joshua B. Tenenbaum, Joshua S. Rule, Joyce Chua, Kamil Kanclerz, Karen Livescu, Karl Krauth,
759 Karthik Gopalakrishnan, Katerina Ignatyeva, Katja Markert, Kaustubh Dhole, Kevin Gimpel,
760 Kevin Omondi, Kory Wallace Mathewson, Kristen Chiafullo, Ksenia Shkaruta, Kumar Shridhar,
761 Kyle McDonell, Kyle Richardson, Laria Reynolds, Leo Gao, Li Zhang, Liam Dugan, Lianhui Qin,
762 Lidia Contreras-Ochando, Louis-Philippe Morency, Luca Moschella, Lucas Lam, Lucy Noble,
763 Ludwig Schmidt, Luheng He, Luis Oliveros-Colón, Luke Metz, Lütfi Kerem Senel, Maarten
764 Bosma, Maarten Sap, Maartje Ter Hoeve, Maheen Farooqi, Manaal Faruqui, Mantas Mazeika,
765 Marco Baturan, Marco Marelli, Marco Maru, Maria Jose Ramirez-Quintana, Marie Tolkiehn,
766 Mario Giulianelli, Martha Lewis, Martin Potthast, Matthew L Leavitt, Matthias Hagen, Mátyás
767 Schubert, Medina Orduna Baitemirova, Melody Arnaud, Melvin McElrath, Michael Andrew
768 Yee, Michael Cohen, Michael Gu, Michael Ivanitskiy, Michael Starritt, Michael Strube, Michał
769 Śwędrowski, Michele Bevilacqua, Michihiro Yasunaga, Mihir Kale, Mike Cain, Mímee Xu, Mirac
770 Suzgun, Mitch Walker, Mo Tiwari, Mohit Bansal, Moin Aminnaseri, Mor Geva, Mozdeh Gheini,
771 Mukund Varma T, Nanyun Peng, Nathan Andrew Chi, Nayeon Lee, Neta Gur-Ari Krakover,
772 Nicholas Cameron, Nicholas Roberts, Nick Doiron, Nicole Martinez, Nikita Nangia, Niklas
773 Deckers, Niklas Muennighoff, Nitish Shirish Keskar, Niveditha S. Iyer, Noah Constant, Noah
774 Fiedel, Nuan Wen, Oliver Zhang, Omar Agha, Omar Elbaghdadi, Omer Levy, Owain Evans,
775 Pablo Antonio Moreno Casares, Parth Doshi, Pascale Fung, Paul Pu Liang, Paul Vicol, Pegah
776 Alipoormolabashi, Peiyuan Liao, Percy Liang, Peter W Chang, Peter Eckersley, Phu Mon Htut,
777 Pinyu Hwang, Piotr Miłkowski, Piyush Patil, Pouya Pezeshkpour, Priti Oli, Qiaozhu Mei, Qing
778 Lyu, Qinlang Chen, Rabin Banjade, Rachel Etta Rudolph, Raefer Gabriel, Rahel Habacker, Ramon
779 Risco, Raphaël Millière, Rhythm Garg, Richard Barnes, Rif A. Saurous, Riku Arakawa, Robbe
780 Raymaekers, Robert Frank, Rohan Sikand, Roman Novak, Roman Sitelew, Ronan Le Bras, Rosanne
781 Liu, Rowan Jacobs, Rui Zhang, Russ Salakhutdinov, Ryan Andrew Chi, Seungjae Ryan Lee, Ryan
782 Stovall, Ryan Teehan, Rylan Yang, Sahib Singh, Saif M. Mohammad, Sajant Anand, Sam Dillavou,
783 Sam Shleifer, Sam Wiseman, Samuel Gruetter, Samuel R. Bowman, Samuel Stern Schoenholz,
784 Sanghyun Han, Sanjeev Kwatra, Sarah A. Rous, Sarik Ghazarian, Sayan Ghosh, Sean Casey,
785 Sebastian Bischoff, Sebastian Gehrmann, Sebastian Schuster, Sepideh Sadeghi, Shadi Hamdan,
786 Sharon Zhou, Shashank Srivastava, Sherry Shi, Shikhar Singh, Shima Asaadi, Shixiang Shane
787 Gu, Shubh Pachchigar, Shubham Toshniwal, Shyam Upadhyay, Shyamolima Shammie Debnath,
788 Siamak Shakeri, Simon Thormeyer, Simone Melzi, Siva Reddy, Sneha Priscilla Makini, Soo-
789 Hwan Lee, Spencer Torene, Sriharsha Hatwar, Stanislas Dehaene, Stefan Divic, Stefano Ermon,
790 Stella Biderman, Stephanie Lin, Stephen Prasad, Steven Piantadosi, Stuart Shieber, Summer
791 Mishnerghi, Svetlana Kiritchenko, Swaroop Mishra, Tal Linzen, Tal Schuster, Tao Li, Tao Yu,
792 Tariq Ali, Tatsunori Hashimoto, Te-Lin Wu, Théo Desbordes, Theodore Rothschild, Thomas Phan,
793 Tianle Wang, Tiberius Nkinyili, Timo Schick, Timofei Kornev, Titus Tunduny, Tobias Gerstenberg,
794 Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz, Uri Shaham, Vedant Misra, Vera
795 Demberg, Victoria Nyamai, Vikas Raunak, Vinay Venkatesh Ramasesh, vinay uday prabhu,
796 Vishakh Padmakumar, Vivek Srikumar, William Fedus, William Saunders, William Zhang, Wout
797 Vossen, Xiang Ren, Xiaoyu Tong, Xinran Zhao, Xinyi Wu, Xudong Shen, Yadollah Yaghoobzadeh,
798 Yair Lakretz, Yangqiu Song, Yasaman Bahri, Yejin Choi, Yichi Yang, Yiding Hao, Yifu Chen,
799 Yonatan Belinkov, Yu Hou, Yufang Hou, Yuntao Bai, Zachary Seid, Zhuoye Zhao, Zijian Wang,
800 Zijie J. Wang, Zirui Wang, and Ziyi Wu. Beyond the imitation game: Quantifying and extrapolating
801 the capabilities of language models. *Transactions on Machine Learning Research*, 2023. ISSN
802 2835-8856. URL <https://openreview.net/forum?id=uyTL5Bvosj>.
- 803 Swabha Swayamdipta, Roy Schwartz, Nicholas Lourie, Yizhong Wang, Hannaneh Hajishirzi, Noah A.
804 Smith, and Yejin Choi. Dataset cartography: Mapping and diagnosing datasets with training
805 dynamics. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.), *Proceedings of the
806 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 9275–
807 9293, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.
808 emnlp-main.746. URL <https://aclanthology.org/2020.emnlp-main.746>.
- 809 Kushal Tirumala, Daniel Simig, Armen Aghajanyan, and Ari S. Morcos. D4: Improving LLM
810 pretraining via document de-duplication and diversification. In *Thirty-seventh Conference on
811 Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. URL <https://openreview.net/forum?id=CG0L2PFrb1>.

810 Mariya Toneva, Alessandro Sordani, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and
 811 Geoffrey J Gordon. An empirical study of example forgetting during deep neural network learning.
 812 *arXiv preprint arXiv:1812.05159*, 2018.

813
 814 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
 815 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand
 816 Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language
 817 models, 2023.

818 Trieu H. Trinh and Quoc V. Le. A simple method for commonsense reasoning, 2019.

819
 820 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
 821 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von
 822 Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Ad-
 823 vances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.,
 824 2017. URL [https://proceedings.neurips.cc/paper_files/paper/2017/
 825 file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).

826
 827 Guillaume Wenzek, Marie-Anne Lachaux, Alexis Conneau, Vishrav Chaudhary, Francisco Guzmán,
 828 Armand Joulin, and Edouard Grave. CCNet: Extracting high quality monolingual datasets
 829 from web crawl data. In Nicoletta Calzolari, Frédéric Béchet, Philippe Blache, Khalid Choukri,
 830 Christopher Cieri, Thierry Declerck, Sara Goggi, Hitoshi Isahara, Bente Maegaard, Joseph Mariani,
 831 Hélène Mazo, Asuncion Moreno, Jan Odijk, and Stelios Piperidis (eds.), *Proceedings of the
 832 Twelfth Language Resources and Evaluation Conference*, pp. 4003–4012, Marseille, France, May
 833 2020. European Language Resources Association. ISBN 979-10-95546-34-4. URL [https://
 834 aclanthology.org/2020.lrec-1.494](https://aclanthology.org/2020.lrec-1.494).

835 Thom Wolfe, Lewis Tunstall, and Patrick von Platen. Jeopardy dataset on hugging face hub. [https://
 836 huggingface.co/datasets/jeopardy](https://huggingface.co/datasets/jeopardy), 2022.

837
 838 Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy Liang, Quoc V
 839 Le, Tengyu Ma, and Adams Wei Yu. Doremi: Optimizing data mixtures speeds up language model
 840 pretraining. In *Thirty-seventh Conference on Neural Information Processing Systems, 2023a*. URL
 841 <https://openreview.net/forum?id=1XuByUeHhd>.

842
 843 Sang Michael Xie, Shibani Santurkar, Tengyu Ma, and Percy Liang. Data selection for language
 844 models via importance resampling. In *Thirty-seventh Conference on Neural Information Processing
 845 Systems, 2023b*. URL <https://openreview.net/forum?id=uPSQv01eAu>.

846
 847 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine
 848 really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for
 849 Computational Linguistics*, 2019.

850 Wenzheng Zhang and Karl Stratos. Understanding hard negatives in noise contrastive estimation.
 851 *arXiv preprint arXiv:2104.06245*, 2021.

852 853 854 7 FULL DATA PRUNING SETTINGS SWEEP 855

856
 857 In this section, we report the results of sweeping over different perplexity-based pruning setting
 858 configurations. In particular, for each dataset, we first sweep over the selection criteria to determine
 859 where from the distribution of perplexities samples should be selected. Then, using the best selection
 860 criteria, we sweep the selection rate to determine how much we should prune.

861 **Setup.** We use the same training and evaluation setup as detailed in Section 3.1. We only perform
 862 the sweep over pruning settings for 1 billion parameter final models for computational budget reasons;
 863 however, we find that the best selection criteria at the 1 billion parameter scale also confers a
 performance improvement at the 3 billion parameter scale, as detailed in 3.2.

Table 4: Results from sweeping different selection criteria. We report the average normalized accuracy for each task grouping as well as across all tasks. While high perplexity selection is optimal for the Pile, medium perplexity selection is optimal for Dolma. Bold results are within one standard error of the highest normalized accuracy.

Pruning Method	World Knowledge	Common Sense Reasoning	Language Understanding	Symbolic Problem Solving	Reading Comprehension	Average
1B Parameters Trained on Pile						
No Pruning (Baseline)	15.51	10.31	28.11	3.53	11.16	13.73
Low Perplexity Selected	11.14	5.76	18.66	3.54	8.72	9.56
Medium Perplexity Selected	16.12	9.01	28.1	3.41	10.86	13.5
High Perplexity Selected	18.18	12.75	33.2	3.36	10.63	15.62
1B Parameters Trained on Dolma						
No Pruning (Baseline)	16.48	12.32	28.86	3.58	7.95	13.84
Low Perplexity Selected	16.13	10.1	27.28	3.45	7.85	12.96
Medium Perplexity Selected	17.98	13.03	31.87	3.44	10.41	15.35
High Perplexity Selected	16.65	13.12	31.14	3.15	8.55	14.52

Table 5: Results from sweeping different selection rates. We report the average normalized accuracy for each task grouping as well as across all tasks. Bold results are within one standard error of the highest normalized accuracy.

Pruning Method	World Knowledge	Common Sense Reasoning	Language Understanding	Symbolic Problem Solving	Reading Comprehension	Average
1B Parameters Trained on Pile						
25% Selection Rate	18.21	12.88	34.44	3.73	9.44	15.74
50% Selection Rate	18.18	12.75	33.2	3.36	10.63	15.62
75% Selection Rate	17.08	10.11	31.37	3.81	9.02	14.28
1B Parameters Trained on Dolma						
25% Selection Rate	17.94	12.16	31.63	3.58	8.91	14.85
50% Selection Rate	17.98	13.03	31.87	3.44	10.41	15.35
75% Selection Rate	18.2	11.78	29.96	3.32	10.82	14.82

7.1 FINDING THE BEST SELECTION CRITERIA

For each dataset, we first sweep the selection criteria while keeping the selection rate fixed at 50%. We report the performance of each selection criteria in Table 4. We find that on the Pile high perplexity selection works the best and on Dolma medium perplexity selection works the best, improving the average downstream performance by 1.89 and 1.51 respectively. An important observation from the sweep is that the best selection criteria from one dataset does not transfer to another dataset and may actually degrade performance compared to the baseline. Although medium-perplexity selection is the best method on Dolma, selecting medium-perplexity samples on the Pile leads to a decrease in the average downstream performance of 0.23 as compared to not performing pruning. These results inform us that high and medium perplexity selection are the optimal selection criteria for the Pile and Dolma respectively, and that the optimal selection criteria does not necessarily transfer between datasets with different domain compositions.

918 7.2 FINDING THE BEST SELECTION RATE

919
920 Using the optimal selection criteria that we found for each dataset, we next investigate the best
921 selection rate for each dataset. We investigate three different selection rates: 25%, 50%, and 75%. We
922 present the results for each selection rate in Table 5. On the Pile, we find that there is no significant
923 difference in downstream performance for selection rates of 25% and 50%; on Dolma we find that a
924 selection rate of 50% achieves the best average downstream performance. For simplicity, we chose
925 to conduct the rest of the experiments in the paper using a selection rate of 50% on both datasets.
926 Furthermore, we find that all the selection rates tested outperform the baseline of no data pruning as
927 measured by average downstream performance. This suggests that the selection criteria has a greater
928 impact on the performance of a pruning configuration than the selection rate.

929 8 DETAILED EVALUATION SETUP

930
931 Jha et al. (2023) also use the MosaicML evaluation gauntlet to perform evaluations in their work. As
932 such, with explicit permission from the authors, we exactly reproduce their text describing the tasks
933 and tasks categories in the evaluation gauntlet. The following is from Section D of their paper:
934

935 The **World Knowledge** category includes the following datasets:

- 936 • Jeopardy (2,117 questions that are a custom subset of the dataset originally obtained from
- 937 Wolfe et al. (2022))
- 938 • MMLU (14,042 four-choice multiple choice questions distributed across 57 categories
- 939 Hendrycks et al. (2021))
- 940 • BIG-bench wikidata (20,321 questions regarding factual information pulled from
- 941 wikipedia) Srivastava et al. (2023)
- 942 • ARC easy (2,376 easy multiple choice middle school science questions) Clark et al. (2018)
- 943 • ARC challenge (1,172 hard multiple choice science questions) Clark et al. (2018)
- 944 • BIG-bench: misconceptions (219 true or false questions regarding common misconceptions)
- 945 Srivastava et al. (2023)
- 946
- 947

948 The **Commonsense Reasoning** category loosely assesses a model’s ability to do basic reasoning
949 tasks that require commonsense knowledge of objects, their properties, and their behavior. It includes
950 the following datasets:

- 951 • BIG-bench Strategy QA (2,289 very eclectic yes/no questions on a wide range of common-
- 952 sense subjects e.g “Can fish get Tonsilitis?”)Srivastava et al. (2023)
- 953 • BIG-bench Strange Stories (174 short stories followed by questions about the charac-
- 954 ters)Srivastava et al. (2023)
- 955 • BIG-bench Novel Concepts (32 find-the-common-concept problems)Srivastava et al. (2023)
- 956 • COPA (100 cause/effect multiple choice questions) Roemmele et al. (2011)
- 957 • PIQA (1,838 commonsense physical intuition 2-choice questions) Bisk et al. (2020)
- 958 • OpenBook QA (500 questions that rely on basic physical and scientific intuition about
- 959 common objects and entities) Mihaylov et al. (2018).
- 960
- 961

962 **Language Understanding** tasks evaluate the model’s ability to understand the structure and properties
963 of languages, and include the following datasets:

- 964 • LAMBADA (6,153 passages take from books - we use the formatting adopted by OpenAI’s
- 965 version)Paperno et al. (2016)
- 966 • HellaSwag (10,042 multiple choice scenarios in which the model is prompted with a scenario
- 967 and choose the most likely conclusion to the scenario from four possible options)Zellers
- 968 et al. (2019)
- 969 • Winograd Schema Challenge (273 scenarios in which the model must use semantics to
- 970 correctly resolve the anaphora in a sentence. The Eval Gauntlet uses the partial evaluation
- 971 technique technique introduced in Trinh & Le (2019)) Levesque et al. (2012)

- 972 • Winogrande (1,267 scenarios in which two possible beginnings of a sentence are presented
973 along with a single ending) Sakaguchi et al. (2020)
- 974 • BIG-bench language identification (10,000 questions on multiple choice language identifica-
975 tion) Srivastava et al. (2023)
- 976 • BIG-bench conceptual combinations (103 questions using made up words) Srivastava et al.
977 (2023)
- 978 • BIG-bench conlang translation (164 example problems in which the model is given transla-
979 tions of simple sentences between English and some fake constructed language) Srivastava
980 et al. (2023)

982 **Symbolic problem solving** tasks test the model’s ability to solve a diverse range of symbolic tasks
983 including arithmetic, logical reasoning, algorithms, and algebra. These datasets include:

- 984 • BIG-bench elementary math QA (38,160 four-choice multiple choice arithmetic word
985 problems) Srivastava et al. (2023)
- 986 • BIG-bench dyck languages (1000 complete-the-sequence questions) Srivastava et al. (2023)
- 987 • BIG-bench algorithms (1,320 questions) Srivastava et al. (2023)
- 988 • BIG-bench logical deduction (1500 four-choice multiple choice questions relating to relative
989 ordering of objects) Srivastava et al. (2023)
- 990 • BIG-bench operators (210 questions involving mathematical operators) Srivastava et al.
991 (2023)
- 992 • BIG-bench repeat copy logic (32 samples in which the model is required to follow some
993 instructions for copying words/symbols)
- 994 • Simple arithmetic with spaces (1000 arithmetic problems consisting of up to 3 operations
995 and using numbers of up to 3 digits, developed by MosaicML)
- 996 • Simple arithmetic without spaces (1000 arithmetic problems consisting of up to 3 operations
997 and using numbers of up to 3 digits, developed by MosaicML)
- 998 • Math QA (2,983 four-choice multiple choice math word problems) Amini et al. (2019)
- 999 • LogiQA (651 four-logical word problems) Liu et al. (2020)

1000 The **Reading comprehension** benchmarks test a model’s ability to answer questions based on the
1001 information in a passage of text. The datasets include:

- 1002 • BIG-bench Understanding fables (189 short stories) Srivastava et al. (2023)
- 1003 • Pubmed QA Labeled (1000 hand-labeled medical documents followed by a related question
1004 for which the model must respond yes/no/maybe) Jin et al. (2019)
- 1005 • SQuAD (10,570 short documents followed by a related question. The model is expected to
1006 output the exact correct answer) Rajpurkar et al. (2016)
- 1007 • BoolQ (3,270 short passages on a diverse range of subjects followed by a yes/no questions)
1008 Clark et al. (2019)

1013 8.1 EVALUATION PROCEDURE

1014 To compute model performance on the above datasets, the Eval Gauntlet uses one of the following
1015 three ICL metrics for each dataset (from MosaicML’s composer library).

- 1016 1. InContextLearningQAAccuracy: This metric uses the query, the corresponding correct
1017 answer and a list of alternative answers to measure a model’s prediction. If the model’s
1018 response conditioned on the query starts with either the correct answer or with one of the
1019 alternative answers, it is considered correct. This is used for question-answering tasks such
1020 as TriviaQA.
- 1021 2. InContextLearningLMAccuracy: This metric tests a model’s ability to output a precise set
1022 of tokens. A model’s output conditioned on a given query is judged to be correct only if the
1023 model’s highest probability tokens match the correct sequence of tokens. This is used for
1024 language modeling tasks such as LAMBADA.

1025

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

3. InContextLearningMultipleChoiceAccuracy: This metric is used for testing a model's ability to answer multiple choice questions accurately. It compares the respective perplexity of the query prepended to each of the possible choices, according to the model. If the query-choice pair with the lowest per token perplexity is indeed the correct choice, then the model's output is judged to be correct. This is used for multiple choice tasks such as HellaSwag, Winograd etc.