

---

# Large Language Model-guided Document Selection

---

Xiang Kong\* Tom Gunter\* Ruoming Pang  
Apple  
{xiang\_kong,tom\_gunter,r\_pang}@apple.com

## Abstract

Large Language Model (LLM) pre-training exhausts an ever growing compute budget, yet recent research has demonstrated that careful document selection enables comparable model quality with only a fraction of the FLOPs. Inspired by efforts suggesting that domain-specific training document selection is in fact an interpretable process [Gunasekar et al., 2023], as well as research showing that instruction-finetuned LLMs are adept zero-shot data labelers [Gilardi et al., 2023], we explore a promising direction for scalable general-domain document selection; employing a prompted LLM as a document grader, we distill quality labels into a classifier model, which is applied at scale to a large, and already heavily-filtered, web-crawl-derived corpus autonomously. Following the guidance of this classifier, we drop 75% of the corpus and train LLMs on the remaining data. Results across multiple benchmarks show that: 1. Filtering allows us to quality-match a model trained on the full corpus across diverse benchmarks with at most 70% of the FLOPs, 2. More capable LLM labelers and classifier models lead to better results that are less sensitive to the labeler’s prompt, 3. In-context learning helps to boost the performance of less-capable labeling models. In all cases we use open-source datasets, models, recipes, and evaluation frameworks, so that results can be reproduced by the community.

## 1 Introduction

Large Language Models (LLMs) (Devlin et al. [2019], Radford et al., 2019], Brown et al. [2020a], Touvron et al. [2023a,b], Achiam et al. [2023], Team et al. [2023]) have demonstrated impressive zero and few-shot generalization across a diverse range of downstream language understanding and generation tasks, after extensive training on large text corpora. Encouraged by the fact that model capabilities seem to continue to improve as more data and compute are put to work, with seemingly no end in sight, scaling law efforts have sought to predict model capability given increased compute and data budgets Kaplan et al. [2020], Hoffmann et al. [2022].

In Sorscher et al. [2023], the authors comment that these power scaling laws show disappointingly slow scaling (a 100x increase in training compute improves cross-entropy loss by only 0.5 nats in Hoffmann et al. [2022], for example), and demonstrate empirically that it may be possible to significantly increase the power-law exponent with careful data selection. Separately, substantial effort typically goes into producing an optimal overall dataset mixture for LLM training Xie et al. [2023a], yet these mixtures always include "high quality" dataset components, which are typically over-sampled (e.g. StackExchange, Wikipedia, or books for the LLaMA model series Touvron et al. [2023a,b]), suggesting that these components make significant—relative to their size—contributions to downstream model quality.

By necessity of volume, the majority of any LLM pre-training dataset will be derived from a filtered bulk web crawl, with a popular source being the CommonCrawl (CC) Authors [2024] corpus. Various

---

\*Equal contribution.

open-source filtered subsets of CC exist, e.g. Gao et al. [2020], Computer [2023], Soboleva et al. [2023], which have already been pruned to only include the highest quality documents. Efforts to improve data pruning, and by doing so the scaling law efficiency, should therefore focus on better selecting documents from the bulk web crawl.

In this work, we explore a promising direction for general-domain web-crawl pruning by introducing a scalable approach for LLM-guided document selection. Our framework utilizes two large language models, a large instruction-finetuned model ( $LM_{large}$ ), and a smaller pre-trained language model ( $LM_{small}$ ). Given a target corpus, we initially leverage the strong zero-shot capabilities of  $LM_{large}$  to assess a number of sampled documents in terms of quality and educational value. Subsequently,  $LM_{small}$  is finetuned on the quality labels from  $LM_{large}$ . Finally, the finetuned  $LM_{small}$  is applied to score the full web-crawl corpus. This distillation step allows us to minimize the total amount of compute required for filtering, as for any given prompt and label-set from  $LM_{large}$  we can establish whether a larger distillation model ( $LM_{small}$ ) is required by simply examining the fine-tuning evaluation metrics.

This paper makes the following contributions:

- We introduce a scalable and general LLM-powered web-crawl document selection framework that does not rely on high-quality reference corpora.
- We apply this pipeline to the filtered RPJ-CC Computer [2023] dataset, as a representative high-quality baseline, and demonstrate significant improvements in downstream LLM quality across multiple parameter and compute budgets.
- We conduct several ablation studies to study the impact of: 1. the labeling prompt, 2. the capability of the labeling model, 3. the capacity of the distillation model.

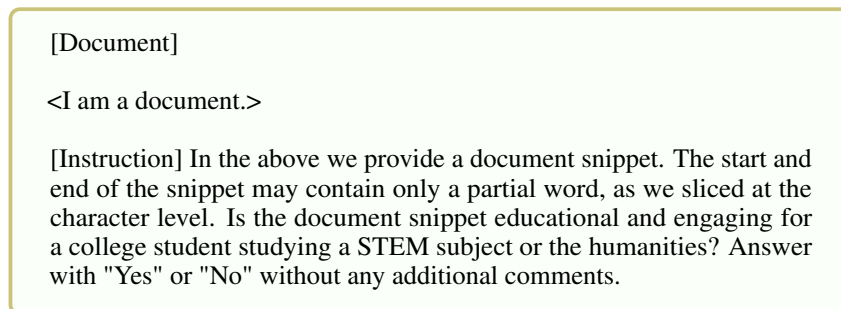


Figure 1: The prompt used to guide the LLM labeler to assess the quality and education value for an input document.

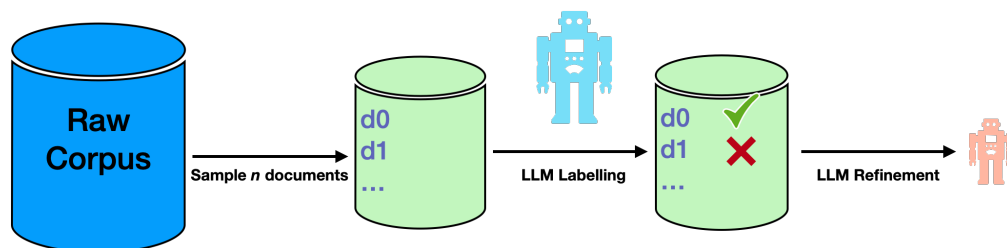


Figure 2: The overall pipeline for our proposed LM-guided data selection pipeline. Given a raw corpus, we first sample  $n$  documents and guide an LLM labeler to assess them in terms of the textual quality and educational value. The resulting (doc, label) pairs could be distilled into an LM-based quality classifier, which will label all documents in the raw corpus.

## 2 Method

In this section we describe the LLM guided web-crawl document selection (LMDS) pipeline.

### 2.1 Proposed Framework

To apply LMDS, we first need access to the following:

- **Target Corpus ( $C$ ):** The source dataset over which document selection is to be applied.  $C$  may contain hundreds of billions of documents sourced from the web, (e.g. plaintext extracted CommonCrawl), with text covering a diverse range of topics, languages, etc.
- $LM_{\text{large}}$ : A capable (the more so the better), instruction-finetuned language model.
- $LM_{\text{small}}$ : A second, low inference cost distillation-target LLM, preferably heavily pretrained (instruction finetuning is optional).

The pipeline is composed of two distinct stages:

**$LM_{\text{large}}$ -Labelling** In this stage, we select a random subset of documents from  $C$  to serve as a representative sample for the whole corpus. The sampled documents are then fed into  $LM_{\text{large}}$  for assessment according to the criteria specified in the prompt (Figure 1), generating high/low-quality document labels. We note that the choice of prompt used to guide  $LM_{\text{large}}$  is important for two reasons: 1. it must be straightforward for the chosen  $LM_{\text{large}}$  to interpret (very capable proprietary models are often able to follow more subtle instructions versus open-source alternatives), 2. it explicitly defines the selection criteria. Assuming a broad definition of "high quality", we find that a concise prompt guiding  $LM_{\text{large}}$  to assess documents for general educational value (inspired by the domain-targeted efforts of Gunasekar et al. [2023]) works well. We leave it to future work to explore more targeted (e.g. subject-specific) prompts, but do believe that the flexibility of the prompt (vs. e.g. sourcing relevant high-quality reference corpora to train a classifier) is a major benefit of LMDS.

**$LM_{\text{small}}$ -Refinement** To enable labelling at very large scale, a cheap-to-run-inference-on (smaller) pretrained language model  $LM_{\text{small}}$  is employed. We first distill the labels sourced from  $LM_{\text{large}}$  on the document assessment task into  $LM_{\text{small}}$ . Once fine-tuning is complete,  $LM_{\text{small}}$  can be deployed to evaluate the entire target corpus  $C$  using your accelerator-backed map-reduce framework of choice.

The overall two-stage pipeline is depicted in Figure 2.

**Rationale for Using  $LM_{\text{small}}$  as well as  $LM_{\text{large}}$**  Given the relative size of a filtered crawl dataset vs the unfiltered corpus, (typically  $\ll 1\%$  of the data remains after selection Computer [2023], Soboleva et al. [2023]), we believe that document selection is a difficult recall problem—it’s probable, given non-perfect filter recall, that a significant number of high quality documents are dropped by typical filter pipelines. For this reason we’d like to minimize the number of stages (asides from LMDS), in the overall filter pipeline, which in turn means that it must be feasible to run LMDS at a hundred-billion document scale. At the same time, we want to use the most powerful  $LM_{\text{large}}$  that we have access to, so that the generated labels adhere to the prompt with high fidelity. To balance this trade-off we introduce  $LM_{\text{small}}$ , and in Table 2 we show that it’s possible to understand the effect of this distillation step (and as a result tune the size of  $LM_{\text{small}}$ ) by examining the F1 score of the classifier at fine-tuning time, with little to no overhead.

## 3 Experiments

### 3.1 Experiment Setup

**Source dataset:** To validate the technique, we will use the filtered Common Crawl (CC) split from the RedPajama v1 data mixture (RPJ-CC) Computer [2023] as  $C$ , because this is a reasonably strong baseline over which we hope to improve. We leave it to future work to demonstrate the effectiveness of LMDS applied to a larger scale unfiltered web-crawl (although we expect this to yield better results still). The RPJ-CC dataset was produced by running five CommonCrawl dumps through the cc\_net pipeline, which de-duplicates (including at a paragraph level) and selects documents using

both perplexity filtering and a linear classifier trained to recognize high-quality Wikipedia reference text. It contains 878 billion tokens.

**Language Model-guided Document Selection Pipeline:** The data selection pipeline described in Section 2 is applied to these datasets.

In the  $\text{LM}_{\text{large}}$ -**Labelling** phase, we use the Llama-2-chat 70B (Touvron et al. [2023a])<sup>2</sup> as our  $\text{LM}_{\text{large}}$ , as it is a widely accessible open-source model with good language understanding and reasoning capabilities. We sample two million documents from RPJ-CC and feed the middle 1500 tokens of each document into  $\text{LM}_{\text{large}}$ , (along with the prompt in Figure 1), for quality assessment<sup>3</sup>.

During the  $\text{LM}_{\text{small}}$ -**Refinement** stage, we take  $\text{LM}_{\text{small}}$  (a smaller LLM pretrained on Wikipedia, Stackexchange, and Books components from Soboleva et al. [2023]) and finetune this model on the (document, quality label) pairs generated by  $\text{LM}_{\text{large}}$ , in the process distilling the signal from  $\text{LM}_{\text{large}}$ . More specifically, a Sigmoid binary classification head is added to  $\text{LM}_{\text{small}}$ , and the representation produced by this head defines the quality score.

At this point  $\text{LM}_{\text{small}}$  is a distilled representation of the labeling function defined by the prompt and  $\text{LM}_{\text{large}}$ , and is then used to quality score  $C$ .

**Language model training on selected documents:** To evaluate the efficacy of our data selection pipeline, we train language models of varying capacities (1B and 7B) on a high quality subset of  $C$  according to the scores from our  $\text{LM}_{\text{small}}$ . For all documents in the target corpus,  $\text{LM}_{\text{small}}$  assigns a quality score in  $(0, 1)$ , then we choose a cutoff and keep only documents scoring above this threshold for training. We drop around 75% of the data from  $C$ , which is in-line with the fraction of documents assigned to the high-quality bucket by  $\text{LM}_{\text{large}}$  at label generation time.

We compare models trained with the same compute budget (and using sequence packing Brown et al. [2020b]), i.e. all models see exactly the same number of tokens and are trained with the same hyperparameters. For details about the architecture and training processes, please refer to Appendix A.

**Language model evaluation:** We assess the language models across a comprehensive set of benchmarks. For all results reported in our study, higher numbers are better than lower numbers.

- **Core English tasks (CoreEN):** Models are evaluated across various benchmarks focusing on on common sense reasoning, reading comprehension and question answering, including ARC Easy and Challenge (0-shot) (Clark et al. [2018]), HellaSwag (0-shot) (Zellers et al. [2019]), WinoGrande (0-shot) (Sakaguchi et al. [2019]), PIQA (0-shot) (Bisk et al. [2019]), SciQ (0-shot) (Sap et al. [2019]), LAMBADA (0-shot) (Paperno et al. [2016]), TriviaQA (1-shot) (Joshi et al. [2017]) and WebQS (1-shot) (Berant et al. [2013]). We use the EleutherAI evaluation-harness Gao et al. [2021] to run these evaluations. In our experiments, we report both the average performance on the seven 0-shot tasks (denoted as CoreEN (0S)) and the average performance across all nine tasks (denoted as CoreEN (All)).
- **MMLU:** Performance is detailed for the 5-shot setup on the MMLU benchmark (Hendrycks et al. [2021]), a widely-adopted benchmark consisting of exam-style questions from 57 tasks.

### 3.2 Experiment Results

**Model quality with a Fixed Training Budget:** In Table 1, we evaluate models of different scales across multiple benchmarks. Models trained on datasets curated using our proposed data selection pipeline consistently achieve higher performance than their counterparts across all benchmarks. For example, at 7B scale, the model trained on filtered Common Crawl achieves **+5.53** on MMLU score compared to the one trained on entire RPJ-CC dataset, demonstrating the effectiveness of our method to select documents with high educational value. The breakdown performance on CoreEN is listed in Appendix B.

<sup>2</sup><https://huggingface.co/meta-llama/Llama-2-70b-chat-hf>

<sup>3</sup>We set temperature as 0.2 to make the results more deterministic.

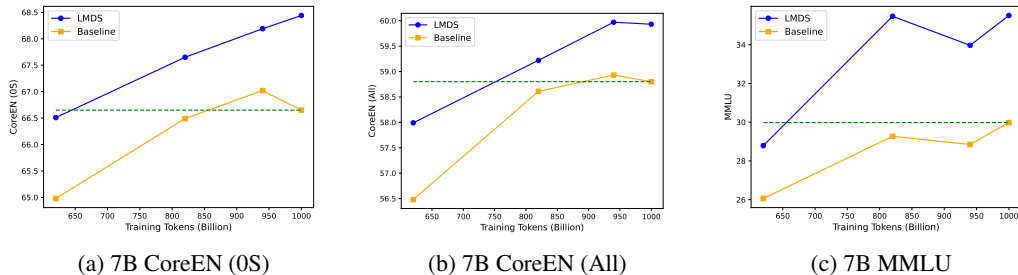


Figure 3: Learning curves on the downstream task for 7B model pretraining on the raw data versus LMDS-based filtered data.

Table 1: Comparisons of models trained on datasets w/ (LMDS) and w/o data selection. We do not report MMLU results for the 1B model, because at this compute budget they are near random.

model	dataset	CoreEN (0S)	CoreEN (All)	MMLU
1B	Baseline	58.10	48.63	N/A
1B	LMDS	<b>61.14</b>	<b>51.41</b>	N/A
7B	Baseline	66.66	58.72	29.98
7B	LMDS	<b>68.44</b>	<b>59.93</b>	<b>35.51</b>

**Data efficiency:** Given the same training budget (FLOPs), training on the LMDS-filtered dataset achieves higher downstream model performance. We further investigate the data efficiency of our approach. As depicted in Figure 3, we show the performance of 7B models trained on w/ and w/o our LMDS pipeline against the number of training tokens. We find that filtering allows us to quality-match a model trained on the full corpus across diverse benchmarks with at most 70% of the FLOPs.

**Downstream Task Accuracy v.s. Selection Ratio:** Using the finetuned  $LM_{small}$  to label the target corpus, each document in the target corpus receives a score ranging from zero to one. One natural question is how to set the cutoff threshold to select documents for model training. In this study, we explore cutoff thresholds to achieve distinct selection ratios ( $\{20\%, 25\%, 30\%, 40\%, 50\%, 100\%\}$ ) and train 1B models on these resulted datasets to understand its impact on the downstream task performance. The results on CoreEN (0S) and CoreEN (All) are shown in Figure 4. Firstly, all selection ratios  $\{20\%, 25\%, 30\%, 40\%, 50\%\}$  show meaningful improvements over the baseline dataset (selection ratio=100%). Secondly, as the data selection ratio is reduced from 55% to 25%, we observe incremental gains in performance, underscoring the importance of data quality to downstream task outcomes, whilst an overly aggressive pruning (reduction from 25% to 20%) begins to reduce performance again. Note that the optimal dataset selection ratio may vary for different corpora, however the raw label-split declared by  $LM_{large}$  is a good rule of thumb when selecting an appropriate threshold.

Table 2: Comparing distillation classifier models across various sizes. Reported F1 score is computed on the validation set. It’s clear that 1. a classifier as small as 85m parameters provides a meaningful improvement over the baseline, 2. larger classifiers are still meaningfully better.

Size	F1-score	CoreEN (0S)	CoreEN (All)
85M	0.78	60.24	50.72
302M	0.81	60.88	51.10
1B	<b>0.84</b>	<b>61.14</b>	<b>51.41</b>

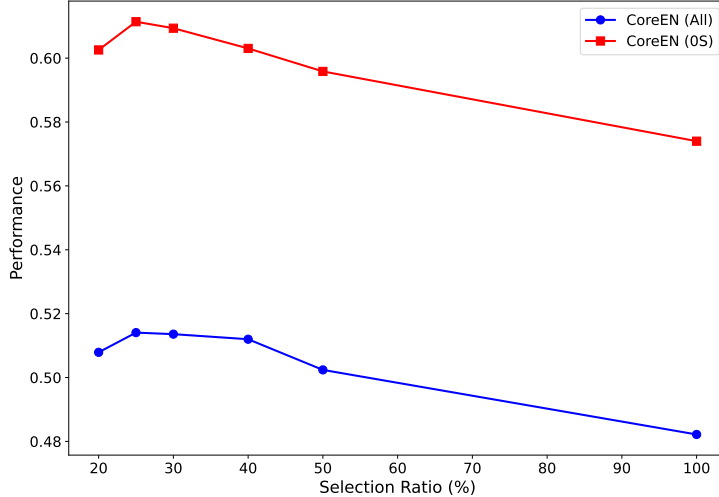


Figure 4: Downstream task accuracy of 1B models trained on LMDS-based filtered datasets with different selection ratios.

Table 3: Comparing the different scale  $LM_{large}$  on downstream task accuracy.

$LM_{large}$	CoreEN (OS)	CoreEN (All)
Baseline	58.10	48.63
Llama-2-chat 7B	57.44	48.13
Llama-2-chat 13B	59.50	50.39
Llama-2-chat 70B	<b>61.14</b>	<b>51.41</b>

**Downstream Task Accuracy v.s. Size of  $LM_{small}$ :** In this study, we examine the impact of different sizes of  $LM_{small}$  on downstream task accuracy. Specifically, we experiment with three sizes of  $LM_{small}$ : 85M, 302M, and 1B parameters. Our first step involves comparing the F1-scores these models achieve on a held-out set when fine-tuned on the (document, quality label) pairs annotated by  $LM_{large}$ . From the results presented in Table 2, it is evident that larger models yield higher F1-scores. This trend suggests that models with greater capacity are more effective at learning from the annotations produced by our large language model,  $LM_{large}$ . To further investigate this phenomenon, we leverage the three fine-tuned  $LM_{small}$  models to label the target corpus. We then train 1B LMs on these datasets. The result is shown in Table 2, we observe that the LM trained on the dataset labeled by the 1B  $LM_{small}$  outperforms those trained on datasets labeled by the 302M and 85M  $LM_{small}$  models, showing that larger capacity  $LM_{small}$  can learn better during the distillation stage. We conclude that models as small as 85m parameters provide a meaningful improvement over the baseline, and that larger models do better still—so suggest that the choice of  $LM_{small}$  capacity is guided by the total inference budget available.

**Impact of  $LM_{large}$ :** Another critical component of our data selection pipeline involves crafting effective prompts that guide large language models  $LM_{large}$  accurately tag documents with respect to their quality and educational value. To achieve this,  $LM_{large}$  needs to have strong instruction following, textual understand and reasoning capabilities. Therefore, in this study, we first try three  $LM_{large}$  with different scales, i.e., Llama-2-chat 7B, 13B and 70B Touvron et al. [2023a]. We first employ these three models to label the same sample documents from the target corpus and finetune  $LM_{small}$  on them. We train 1B LMs on data labeled by these  $LM_{small}$ . From the Table 3, we find that the larger  $LM_{large}$  is, the better accuracy we can obtain. For the Llama-2-chat 7B, it only achieve similar accuracy compared to the baseline and after a closer look, we find that approximately 98% of the sample documents are labeled as *yes*, indicating that the Llama-2-chat 7B model lacks the ability to analyze intricate and nuanced textual content effectively.

Table 4: The robustness of  $LM_{\text{large}}$  with various sizes to prompt instructions. We train models based on three versions of instructions for  $LM_{\text{large}}$  to label sample documents.

$LM_{\text{large}}$	Prompt Version	CoreEN (All)	Avg. (Std)
Llama-2-chat 13B	V1	50.05	$50.27_{\pm 0.57}$
	V2	50.86	
	V3	49.90	
Llama-2-chat 70B	V1	51.41	$51.24_{\pm 0.24}$
	V2	51.34	
	V3	50.96	

Table 5: Comparing the Llama-2-chat 7B w/ and w/o in-context learning on downstream task accuracy.

$LM_{\text{large}}$	CoreEN (0S)	CoreEN (All)
Llama-2-chat 7B	57.44	48.13
+ ICL (5-shot)	58.55	48.80
Llama-2-chat 70B	61.14	51.41

Furthermore, we also try to understand the robustness of Llama-2-chat 13 and 70B models with respect to instructions. Ideally,  $LM_{\text{large}}$  should provide a consistent quality assessment given semantically similar instructions. Therefore, we try three prompts (more details in Appendix C), and compute the agreement of downstream evaluation results across these prompts. As shown in Table 4, the Llama-2-chat 70B model exhibits higher agreement compared to the Llama-2-chat 13B, indicating that larger capability  $LM_{\text{large}}$  is more reliable and less influenced by variations in prompts.

**In-Context Learning for Llama-2-chat 7B:** Due to its limited instruction-following capabilities, Llama-2-chat 7B struggles to distinguish high-quality documents from low-quality ones, as compared to its larger siblings. However, large language models exhibit in-context learning (ICL) abilities (Brown et al. [2020a]), where they can learn from a few examples provided in the context without fine-tuning and by doing so overcome instruction-following shortcomings. To leverage this, we provide 5 examples from the strongest  $LM_{\text{large}}$  Llama-2-chat 70B in this study as a demonstration context for Llama-2-chat 7B. We then trained 1B models based on the ICL-enhanced Llama-2-chat 7B results. As shown in Table 5, incorporating demonstrations from a stronger  $LM_{\text{large}}$  yields better downstream accuracy, showing ICL helps the Llama-2-chat 7B to more effectively identify high-quality documents. However, there still remains a significant gap compared to  $LM_{\text{large}}$ , indicating that raw model capability (especially instruction following) is critical for accurate document labeling.

## 4 Related Work

The goal of data filtering is to determine the optimal subset of training data to train on, where optimal is typically measured via trained model performance on a diverse suite of benchmarks. It is well known that better data selection leads to improved downstream performance and/or a reduced compute budget to achieve the same performance for large-scale model training efforts, e.g. LLMs Sorscher et al. [2023].

We believe that the most closely related work to our own is Sachdeva et al. [2024], (developed concurrently), where the authors employ a T5 Raffel et al. [2020] model to produce quality labels before training new T5 models on the filtered data via a process they term "Ask-LLM". Compared to Ask-LLM, we focus on decoder-only LLMs, demonstrating that the technique works for much larger models and compute budgets (7B parameters trained for 1T tokens), as well as tougher benchmarks (MMLU), and also works when starting with a stronger baseline dataset in RPJ-CC (versus C4Raffel et al. [2020]). Furthermore, we ablate the impact of the prompt as well as the choice of labeler and

distillation model, showing that more capable labelers are both more effective and less sensitive to the choice of prompt.

In what follows we describe more broadly approaches for data selection by quality, with a focus on efforts applied to language model training.

**Effect of data selection on pre-training stage:** Training large language models typically involves utilizing extensive text datasets sourced from massive and diverse origins such as CommonCrawl and GitHub. However, it is widely recognized that pre-training data is very noisy, and often includes boilerplate text, templates, error messages, and other forms of repetitive or not-informative content that does not contribute meaningfully to model quality as measured by downstream benchmarks. This poor quality data can lead models to learn and perpetuate irrelevant or redundant information, ultimately impacting their performance and generalization capabilities across various tasks (Raffel et al. [2020], Touvron et al. [2023a,b]). For instance, during the creation of the Pile dataset (Gao et al. [2020]), authors find that the most common 13-grams are character repetitions, such as strings of dashes, with more than 10 million instances. Removing such undesirable text is crucial but must be done efficiently due to the large number of documents involved. A common approach to filtering data is to employ simple yet computationally efficient heuristics. The goal of these heuristic approaches is to constrain the training distribution along certain dimensions (e.g., sentence length, repetitiveness), with the assumption that the evaluation distribution will exhibit similar characteristics. The heuristics used in past works are extensive but generally fall into one of the following categories: item count, repetition count, ratio, or per-document statistical measures.

In addition to heuristic data selection methods, many works focus on developing more advanced model-based techniques to identify and select "high quality" text. Brown et al. [2020a], Du et al. [2022], Touvron et al. [2023a], Computer [2023] try to select data in a noisy dataset based on the similarity to a pre-defined high quality reference corpus, e.g., Wikipedia, to improve the model quality. Gao et al. [2020] train GPT-2 (Achiam et al. [2023]) and GPT-3 (Brown et al. [2020a]) models from scratch to compare the "Pile" dataset with CommonCrawl to demonstrate the effectiveness of data filtering. Raffel et al. [2020] train 1B parameter models to show the positive impact of the de-duplication and quality filters. Hernandez et al. [2022] train multiple language models with various epoch counts but fixed compute budgets to understand the effect of repeated data. Muennighoff et al. [2024] explore the effect of dataset repetition on scaling laws, finding that repeating the highest quality data for up to four epochs outperforms fresh, less-heavily filtered data seen only once. In Tirumala et al. [2024], the authors conduct a large-scale data-efficient pre-training evaluation, showing that a 6.7B OPT model (Zhang et al. [2022]) can converge up to 20% faster on data curated by a technique based on pre-trained model embeddings on top of de-duplicated data. Instead of hard-selection, Xie et al. [2023b] propose a data selection algorithm with importance resampling to estimate importance weights and select data according to these weights during the LLM training. One line of research focuses on the data mixture between multiple components. Longpre et al. [2023] discuss the effect of several factors including toxicity, age and domain distribution on the downstream performance. Xie et al. [2024] propose to seek data mixture ratios through training a small proxy model using Group DRO, which could be used to train large scales models. Another series of studies focuses on training models using a selection of "textbook quality" data to demonstrate the importance of data quality (Javaheripi et al. [2023], Gunasekar et al. [2023], Li et al. [2023b], Eldan and Li [2023]). In, Maini et al. [2024], the authors use off-the-shelf instruction-finetuned language models to paraphrase documents on the web to improve the data quality without selection, and as previously mentioned, Sachdeva et al. [2024] employs T5 to choose the high quality data with density sampling to preserve the coverage and diversity of the training data at the same time. Zhang et al. [2024] employs a language model to help identify domain specific (Math) documents to boost math-related benchmark performance.

**Effect of data selection on post-training stage:** In contrast to the pre-training, the post-training stage aims to align LLMs with a downstream user's objective to make the model outputs more controllable and helpful for users (Ouyang et al. [2022]). Many works aim to improve the data quality through focusing on the difficulty, complexity and diversity as components of the utility. Zhou et al. [2024] show that only 1,000 carefully curated data with high quality and diversity can achieve superior instruction following capabilities. Li et al. [2023a] design a difficulty metric from instruction-following perspective to let model itself to select the training samples. Kung et al. [2023] propose *activate instruction tuning* to identify informative tasks to continuously improve its cross-task



generalization ability. Although human guidance may be more reliable, recently, the capability of models to precisely identify undesirable data points is also effective. For example, Chen et al. [2023] leverage ChatGPT Achiam et al. [2023] to assess training examples and select a subset of Alpaca dataset. Gunasekar et al. [2023] build an LM-based classifier to filter textbook quality data from the target corpus and this classifier is trained on data generated from LLMs. Some off-the-shelf LLMs such as GPT-4 (Achiam et al. [2023]) are used to filter and rank responses to provide high-quality feedback (Cui et al. [2023], Zhu et al. [2023]).

## 5 Conclusion

This paper introduces a scalable and general large language model-guided document pruning pipeline for the autonomous selection of high-quality training data for large language models (LLMs) from a web-scale corpus. Two language models with different scales,  $LM_{\text{large}}$  and  $LM_{\text{small}}$  achieve a good balance of selection precision/recall and inference efficiency/scalability. We demonstrate the effectiveness of this approach by applying it to the already-filtered RPJ-CC dataset, at multiple model size and training FLOP budgets. Results showed that models trained on the filtered dataset achieve significantly better quality across multiple benchmarks, in most cases doing so with a reduced compute budget. Several ablation studies highlight the importance of key components in our framework. In future work we hope to demonstrate the effectiveness of this pipeline applied at full scale to a (near) unfiltered web-corpus, and expect that domain targeted  $LM_{\text{large}}$  prompts will also allow for the targeted selection of domain-specific data for domains which do not have an existing high-quality and high-coverage reference corpus.

## References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- The Common Crawl Authors, 2024. URL <https://commoncrawl.org/overview>.
- Jonathan Berant, Andrew K. Chou, Roy Frostig, and Percy Liang. Semantic parsing on freebase from question-answer pairs. In *Conference on Empirical Methods in Natural Language Processing*, 2013. URL <https://api.semanticscholar.org/CorpusID:6401679>.
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about physical commonsense in natural language, 2019.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020a.
- 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. Language models are few-shot learners, 2020b.
- Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa Gunaratna, Vikas Yadav, Zheng Tang, Vijay Srinivasan, Tianyi Zhou, Heng Huang, et al. Alpapasus: Training a better alpaca with fewer data. *arXiv preprint arXiv:2307.08701*, 2023.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge, 2018.
- Together Computer. Redpajama: An open source recipe to reproduce llama training dataset, 2023. URL <https://github.com/togethercomputer/RedPajama-Data>.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. Ultrafeedback: Boosting language models with high-quality feedback. *arXiv preprint arXiv:2310.01377*, 2023.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Tamar Solorio, editors, *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, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.
- Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. Glam: Efficient scaling of language models with mixture-of-experts. In *International Conference on Machine Learning*, pages 5547–5569. PMLR, 2022.
- Ronen Eldan and Yuanzhi Li. Tinstories: How small can language models be and still speak coherent english? *arXiv preprint arXiv:2305.07759*, 2023.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, Jason Phang, Laria Reynolds, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot language model evaluation, September 2021. URL <https://doi.org/10.5281/zenodo.5371628>.

- Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli. Chatgpt outperforms crowd workers for text-annotation tasks. *Proceedings of the National Academy of Sciences*, 120(30), July 2023. ISSN 1091-6490. doi: 10.1073/pnas.2305016120. URL <http://dx.doi.org/10.1073/pnas.2305016120>.
- Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, Adil Salim, Shital Shah, Harkirat Singh Behl, Xin Wang, Sébastien Bubeck, Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and Yuanzhi Li. Textbooks are all you need, 2023.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- Danny Hernandez, Tom Brown, Tom Conerly, Nova DasSarma, Dawn Drain, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Tom Henighan, Tristan Hume, et al. Scaling laws and interpretability of learning from repeated data. *arXiv preprint arXiv:2205.10487*, 2022.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training compute-optimal large language models, 2022.
- Mojan Javaheripi, Sébastien Bubeck, Marah Abdin, Jyoti Aneja, Sebastien Bubeck, Caio César Teodoro Mendes, Weizhu Chen, Allie Del Giorno, Ronen Eldan, Sivakanth Gopi, et al. Phi-2: The surprising power of small language models. *Microsoft Research Blog*, 2023.
- Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension, 2017.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *CoRR*, abs/2001.08361, 2020. URL <https://arxiv.org/abs/2001.08361>.
- Taku Kudo and John Richardson. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing, 2018.
- Po-Nien Kung, Fan Yin, Di Wu, Kai-Wei Chang, and Nanyun Peng. Active instruction tuning: Improving cross-task generalization by training on prompt sensitive tasks. *arXiv preprint arXiv:2311.00288*, 2023.
- Ming Li, Yong Zhang, Zhitao Li, Jiuhai Chen, Lichang Chen, Ning Cheng, Jianzong Wang, Tianyi Zhou, and Jing Xiao. From quantity to quality: Boosting llm performance with self-guided data selection for instruction tuning. *arXiv preprint arXiv:2308.12032*, 2023a.
- Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. Textbooks are all you need ii: phi-1.5 technical report. *arXiv preprint arXiv:2309.05463*, 2023b.
- Shayne Longpre, Gregory Yauney, Emily Reif, Katherine Lee, Adam Roberts, Barret Zoph, Denny Zhou, Jason Wei, Kevin Robinson, David Mimno, et al. A pretrainer’s guide to training data: Measuring the effects of data age. *Domain Coverage, Quality, & Toxicity*, 2023.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019.
- Pratyush Maini, Skyler Seto, He Bai, David Grangier, Yizhe Zhang, and Navdeep Jaitly. Rephrasing the web: A recipe for compute and data-efficient language modeling. *arXiv preprint arXiv:2401.16380*, 2024.
- Niklas Muennighoff, Alexander Rush, Boaz Barak, Teven Le Scao, Nouamane Tazi, Aleksandra Piktus, Sampo Pyysalo, Thomas Wolf, and Colin A Raffel. Scaling data-constrained language models. *Advances in Neural Information Processing Systems*, 36, 2024.

- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Quan Ngoc Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. The lambada dataset: Word prediction requiring a broad discourse context, 2016.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- Noveen Sachdeva, Benjamin Coleman, Wang-Cheng Kang, Jianmo Ni, Lichan Hong, Ed H Chi, James Caverlee, Julian McAuley, and Derek Zhiyuan Cheng. How to train data-efficient llms. *arXiv preprint arXiv:2402.09668*, 2024.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale, 2019.
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. Socialiqa: Commonsense reasoning about social interactions, 2019.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units, 2016.
- Daria Soboleva, Faisal Al-Khateeb, Robert Myers, Jacob R Steeves, Joel Hestness, and Nolan Dey. SlimPajama: A 627B token cleaned and deduplicated version of RedPajama. <https://www.cerebras.net/blog/slimpajama-a-627b-token-cleaned-and-deduplicated-version-of-redpajama>, 2023. URL <https://huggingface.co/datasets/cerebras/SlimPajama-627B>.
- Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari S. Morcos. Beyond neural scaling laws: beating power law scaling via data pruning, 2023.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- Kushal Tirumala, Daniel Simig, Armen Aghajanyan, and Ari Morcos. D4: Improving llm pretraining via document de-duplication and diversification. *Advances in Neural Information Processing Systems*, 36, 2024.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- Mitchell Wortsman, Peter J Liu, Lechao Xiao, Katie Everett, Alex Alemi, Ben Adlam, John D Co-Reyes, Izzeddin Gur, Abhishek Kumar, Roman Novak, et al. Small-scale proxies for large-scale transformer training instabilities. *arXiv preprint arXiv:2309.14322*, 2023.
- Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy Liang, Quoc V. Le, Tengyu Ma, and Adams Wei Yu. Doremi: Optimizing data mixtures speeds up language model pretraining, 2023a.

- Sang Michael Xie, Shibani Santurkar, Tengyu Ma, and Percy S Liang. Data selection for language models via importance resampling. *Advances in Neural Information Processing Systems*, 36: 34201–34227, 2023b.
- Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy S Liang, Quoc V Le, Tengyu Ma, and Adams Wei Yu. Doremi: Optimizing data mixtures speeds up language model pretraining. *Advances in Neural Information Processing Systems*, 36, 2024.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence?, 2019.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.
- Yifan Zhang, Yifan Luo, Yang Yuan, and Andrew C Yao. Autonomous data selection with language models for mathematical texts. In *ICLR 2024 Workshop on Navigating and Addressing Data Problems for Foundation Models*, 2024.
- Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36, 2024.
- Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu, and Jiantao Jiao. Starling-7b: Improving llm helpfulness & harmlessness with rlaif, 2023.

Table 6: All model specifications used in this paper.

model	num layers	hidden dim	num heads	batch size (# tokens)
85M	12	768	12	1M
302M	24	1024	16	1M
1B	24	2048	32	1M
7B	32	4096	32	4M

## A Modeling Training Details

**LM<sub>small</sub> pretraining:** Before finetuning LM<sub>small</sub> on (document, label) pair labelled by LM<sub>large</sub>, we first pretrain it on Wikipedia, Stackexchange, and Books components from Soboleva et al. [2023] to build a good initialization for finetuning.

The architecture of all language models including LM<sub>small</sub> and LM<sub>large</sub> is listed in the Table 6. All models are pretrained with the AdamW optimizer (Loshchilov and Hutter [2019]), employing parameters  $\beta_1 = 0.9$ ,  $\beta_2 = 0.95$ , decoupled  $wd = 1e - 4$  and  $\epsilon = 10^{-8}$ . We implement a cosine learning rate schedule, starting with 2000 warmup steps and decaying to 10% of peak learning rate and a gradient clipping at 1.0. A SentencePiece tokenizer (Kudo and Richardson [2018]), employing a vocabulary of 32k and built on the Byte-Pair Encoding algorithm (Sennrich et al. [2016]), is used on C4. All the model are trained on TPU-v4 and v5 chips. All LM<sub>small</sub> are trained with 400B tokens. 1B and 7B main models are trained with 300B and 1T tokens respectively. We use the  $\mu$ Param (simple) (Wortsman et al. [2023]) where the peak learning rate is set to 1e-2.

## B Breakdown performance on CoreEN tasks

We list performance of different models on CoreEN tasks in Table 7.

Table 7: Breakdown results of 1B models on CoreEN in Table 1

	Baseline	LMDS
arc_c	26.37	35.07
arc_e	62.21	70.54
hellaswag	43.75	43.16
lambada	56.44	58.88
piqa	71.22	69.59
sciq	86.80	90.20
winogrande	59.91	60.54
triviaqa	17.68	20.99
webqs	13.29	13.68

## C Instructions exploration for LM<sub>large</sub>

We list the prompts we used to test the robustness of different LM<sub>large</sub>. They are semantically similar to each other and a strong language model should be less sensitive to these prompts.

Table 8: Breakdown results of 7B models on CoreEN in Table 1

	Baseline	LMDS
arc_c	40.02	44.45
arc_e	72.06	77.10
hellaswag	50.88	52.04
lambada	70.25	69.75
piqa	73.78	74.65
sciq	92.50	94.30
winogrande	66.06	66.77
triviaqa	40.79	39.13
webqs	21.90	21.16

V1	<b>[Instruction]</b> In the above we provide a document snippet. The start and end of the snippet may contain only a partial word, as we sliced at the character level. Is the document snippet educational and engaging for a college student studying a STEM subject or the humanities? Answer with "Yes" or "No" without any additional comments.
V2	<b>[Instruction]</b> In the above we provide a document snippet. The start and end of the snippet may contain only a partial word, as we sliced at the character level. Does the document look like it would be helpful for a STEM or Humanities student who is struggling with their course? Answer with "Yes" or "No" without any additional comments.
V3	<b>[Instruction]</b> In the above we provide a document snippet. The start and end of the snippet may contain only a partial word, as we sliced at the character level. Does the document look like it would be educational and helpful for a STEM or Humanities student to help understanding material from their course? Answer with "Yes" or "No" without any additional comments.

Table 9: Prompts used to test the robustness of different  $\text{LM}_{\text{large}}$ .