

RETHINKING DATA CURATION IN LLM TRAINING: ONLINE REWEIGHTING OFFERS BETTER GENERALIZATION THAN OFFLINE METHODS

Wanru Zhao^{1*}, Yihong Chen², Yuzhi Tang^{3 †}, Wentao Ma^{3 †}, Shengchao Hu⁴,
Shell Xu Hu⁵, Alex Jacob¹, Abhinav Mehrotra⁵, Nicholas D. Lane¹

¹ University of Cambridge, ² OATML, University of Oxford,

³ University of Toronto, ⁴ Shanghai Jiao Tong University, ⁵ Samsung AI Center

ABSTRACT

Data curation is a critical yet underexplored area in large language model (LLM) training. Existing methods, such as data selection and mixing, operate in an *offline* paradigm, detaching themselves from training. This separation introduces engineering overhead and makes the curation brittle: the entire pipeline must be re-run under model/task shifts. Moreover, offline methods alter data size through hard filtering or resampling, often sacrificing data diversity and harming generalization.

We propose to rethink data curation as an *online* reweighting problem, where sample importance is dynamically adjusted during training via loss weighting rather than static pre-processing. Specifically, we introduce **ADAPT** (Adaptive Data reweighting for Pretraining and FineTuning), a dynamic online framework that reweights training samples with adaptive per-sample learning rates guided by similarity-based quality signals, without changing the number of training samples. Unlike offline methods that enforce a static data distribution, ADAPT acts as an implicit curriculum learner, progressively shifting focus from coarse-grained patterns to fine-grained semantic distinctions as the model evolves. Experiments on both instruction tuning and large-scale pretraining show that ADAPT consistently outperforms offline selection/mixing and prior online methods, achieving stronger cross-benchmark generalization under equal FLOPs.

 <https://huggingface.co/DATA-ADAPT>

 <https://github.com/RyanOv0/ADAPT>

1 INTRODUCTION

Large language models (LLMs) have demonstrated remarkable capabilities across diverse tasks (Brown et al., 2020; Chowdhery et al., 2022; Touvron et al., 2023a), yet their generalization usually hinges on the quality of the training data, including its diversity and mixing (Hoffmann et al., 2022; Kaplan et al., 2020a). While some recent efforts seek to circumvent the data quality issue and improve generalization through architectural interventions (Peyrard et al., 2022; Chen et al., 2023c; Bengio et al., 2025), most pragmatic practices still rely on curating the “right” dataset for the model and the task. Data curation – deciding what data to keep and how to weight it – has thus become a cornerstone of large language model development. On the other hand, the process of data curation has largely remained a secret sauce with pipelines relying on opaque heuristics, ad-hoc engineering choices, or costly trial-and-error, rather than principled frameworks.

Current approaches to data curation largely follow an *offline* paradigm, falling into two camps: data selection, which keeps a subset of valuable examples, and data mixing, which adjusts sampling frequencies to rebalance distributions. Despite their differences, both approaches typically follow a multi-stage pipeline: (1) train or extract features with a proxy model, (2) compute quality signals as selection or weighting criteria on a validation set, and (3) retrain the main model on the curated data. While effective in some cases, this paradigm suffers from fundamental drawbacks (Sec. 4) due to its static and subtractive nature. First, it ignores training dynamics: the value of a data sample is not

*Corresponding to: Wanru Zhao (wz341@cam.ac.uk)

† Equal contribution.

static but changes as the model learns. Offline selection freezes this value based on a proxy model, creating a misalignment with the training model’s evolving needs. Second, it compromises diversity: by relying on hard filtering to select a fixed subset, offline methods risk discarding the broad data distribution essential for robust generalization. Consequently, the resulting subsets often overfit to specific benchmarks and are brittle, which requires the entire pipeline to be re-run whenever the model architecture or corpus changes.

In contrast, we propose a paradigm shift to *online reweighting*, which fundamentally differs from offline selection not just in *when* the signal is computed, but also in *how* the data is used. Instead of hard-filtering a subset based on static scores, online reweighting retains the full dataset to preserve diversity, but dynamically modulates each sample’s contribution via loss weighting. This preserves data diversity compared to offline’s “hard cuts” while intelligently controlling each sample’s “step size” in parameter updates to adapt to evolving model states.

To solve the limitation of offline data curation, we propose a unified online reweighting framework, **ADAPT** (Adaptive Data reweighting for Pretraining and FineTuning), which aligns training samples with a validation set via similarity-based weighting. Our method incurs nearly zero additional overhead yet achieves superior cross-benchmark generalization. Under a unified evaluation protocol, it consistently outperforms both offline selection/mixing strategies and prior online reweighting methods, across large-scale pretraining and instruction tuning setups.

Our contributions. This work makes the following contributions:

1. **🗂️ Formalization and unification of dataset curation pipelines.** We present a unified view of data curation, showing that data selection, mixing, and balancing can all be expressed as special cases of *data reweighting*. This perspective shifts the focus from offline pre-processing to dynamic, in-training adjustment of data weights. Which offers a fair protocol we propose to view and evaluate them through a unified lens of *online data reweighting*.
2. **🔧 Online sample-level reweighting method.** We introduce a new online global reweighting algorithm that is effective under this unified formalization and protocol. It employs an adaptive per-sample learning rate guided by sample quality signals, without explicitly altering the dataset size, while dynamically adjusting sample weights during training.
3. **📈 Better performance and generalization across tasks.** On both instruction tuning and large-scale pretraining experiments, our framework achieves up to up to 7.2% higher accuracy per FLOP than offline baselines including LESS (Xia et al., 2024), DoReMi (Xie et al., 2023b), and RegMix (Liu et al., 2024), and other online reweighting method (Sow et al., 2025), and demonstrates 6.1% stronger out-of-domain generalization across benchmarks.

2 BACKGROUND AND RELATED WORK

2.1 DATA CURATION FOR LARGE LANGUAGE MODELS

Researchers and engineers usually rely on web crawls to gather large datasets for training large language models (Brown, 2020; Computer, 2023; Penedo et al., 2024; Tang et al., 2024; Bai et al., 2023; Kandpal et al., 2025). However, these crawls often include a substantial amount of low-quality or irrelevant content, which makes data curation necessary to build high-quality training sets. Most data curation efforts focus on methods for improving model performance (Raffel et al., 2020; Brown, 2020; Rae et al., 2021; Penedo et al., 2023; Soldaini et al., 2024), including filtering by language (Raffel et al., 2020; Xue et al., 2020), heuristic-based filtering (Gao et al., 2020; Rae et al., 2021; Penedo et al., 2023; Soldaini et al., 2024), quality filtering (Du et al., 2022; Xie et al., 2023c), data deduplication (Lee et al., 2021) and mixing (Xie et al., 2023a; Soboleva et al., 2023; Albalak et al., 2023). However, current approaches to curating such datasets are generally ad-hoc. We aim to develop a principled and automated method for data curation that can also unify different processing stages.

Pretraining data curation. Several recent studies (Xie et al., 2023a; Chen et al., 2023b; Fan et al., 2023; Thakkar et al., 2023) have explored various reweighting techniques to enhance the generalization and efficiency of language models pretraining. For instance, Xie et al. (2023a) and Fan et al. (2023) are aiming at finding the optimal mixture of pretraining corpora to enhance performance

across domains. [Chen et al. \(2023b\)](#) propose an ordered skill learning method for data selection measuring how effectively it teaches interdependent skills for continual pretraining and fine-tuning. Although effective, these works are aiming at the group level, whereas our work explores reweighting at the sample level, offering fine-grained control during model training dynamics. Instance-level reweighting has been used in pretraining settings of LLMs ([Chen et al., 2024](#); [Jiang et al., 2024](#)), where each sample per mini-batch is weighted over how individual samples are treated based on their loss values. In contrast, our work studies the effects of various adaptive learning rate with different quality signals considering the model states to enhance both performance and generalization for LLMs pretraining and fine-tuning regimes.

Instruction data curation. Research has demonstrated that prioritizing data quality and diversity over quantity is more helpful for instruction-following capabilities ([Cao et al., 2023](#); [Chen et al., 2023a](#); [Bukharin & Zhao, 2023](#); [Du et al., 2023](#); [Liu et al., 2023](#); [Li et al., 2023](#)). Instruction tuning data includes task-based datasets curated from traditional NLP tasks ([Wang et al., 2022](#); [Sanh et al., 2022](#); [Wei et al., 2022a](#); [Longpre et al., 2023](#)), and open-ended datasets ([Taori et al., 2023](#); [Conover et al., 2023](#); [Köpf et al., 2023](#); [Xu et al., 2023](#); [Mukherjee et al., 2023](#); [Zhou et al., 2023](#); [Ding et al., 2023](#)). In our work, we are focusing on adaptively adjust the data weight during instruction tuning without explicitly reduce the number of instruction data samples and obtain model with generalization.

2.2 COMMONLY USED QUALITY SIGNAL FOR DATA CURATION

- **Lexical Similarity (BM25)** ([Silva & Barbosa, 2024](#)) quantifies term-based overlap between training and validation data through sparse retrieval scoring:

$$s_{\text{BM25}}(x) = \frac{1}{|\mathcal{D}_{\text{val}}|} \sum_{v \in \mathcal{D}_{\text{val}}} \text{BM25}(x, v). \quad (1)$$

This metric captures surface-level textual similarity without semantic understanding.

- **Semantic Similarity (Embedding)** ([Rubin et al., 2021](#)) measures dense representation alignment using pretrained encoders:

$$s_{\text{Embed}}(x) = \frac{1}{|\mathcal{D}_{\text{val}}|} \sum_{v \in \mathcal{D}_{\text{val}}} \cos(\phi(x), \phi(v)), \quad (2)$$

where $\phi(\cdot)$ denotes a frozen embedding model. This approach captures semantic proximity beyond lexical overlap.

- **Distributional Alignment (Perplexity)** ([Antonello et al., 2020](#)) evaluates likelihood under a reference language model θ_0 :

$$s_{\text{PPL}}(x) = -\log P_{\theta_0}(x). \quad (3)$$

Lower perplexity indicates stronger distributional alignment with the reference corpus.

- **Gradient-Based Influence** ([Xia et al., 2024](#)) estimates training utility through first-order approximation of validation loss reduction:

$$s_{\text{Grad}}(x) = \langle \nabla_{\theta} \ell(x; \theta_0), \nabla_{\theta} \ell(\mathcal{D}_{\text{val}}; \theta_0) \rangle, \quad (4)$$

where θ_0 denotes a proxy model. This metric directly quantifies how training on x influences performance on the validation set \mathcal{D}_{val} .

3 UNIFIED FORMALIZATION OF DATA CURATION

Setup. Given two datasets: the *train* set $\mathcal{D}_{\text{train}}$ and the *validation* set \mathcal{D}_{val} . Usually, the size of the validation set is much smaller than the train set. When \mathcal{D}_{val} is sampled from the distribution of downstream test data, the validation dataset \mathcal{D}_{val} is considered as *in-domain*. The train dataset $\mathcal{D}_{\text{train}}$ on the other hand consists of both *in-domain* and *out-of-domain* samples. This is the case when web crawling is used to collect training data from the whole internet. Our goal is to subsample or weight the train set $\mathcal{D}_{\text{train}}$ under the guidance of the validation set \mathcal{D}_{val} , so that the model training is less affected by the out-of-domain samples.

3.1 DESIGN SPACE OF DATA QUALITY SIGNAL

In Sec. 2, we present a unified framework for data quality assessment where *scoring function* takes the form $s(x) \equiv s(x; \theta, \mathcal{D}_{\text{val}})$ that assigns a quality signal for each data example x , which optionally depends on a (proxy) model θ and the validation set \mathcal{D}_{val} .

3.2 OFFLINE DATA CURATION

The scores are employed once before training (e.g., with a proxy model θ_0 and fixed validation set \mathcal{D}_{val}). The resulting weights $\{w(x)\}_{x \in \mathcal{D}_{\text{train}}}$ are fixed and used to make a decision to pass which training examples to the real training stage.

- **Data Selection (sample-level binary weights)** Data selection (or filtering) removes part of the training corpus before pretraining begins. It keeps data examples with quality signal above a certain threshold τ :

$$\mathcal{D}'_{\text{train}} = \{x \in \mathcal{D}_{\text{train}} \mid s(x) \geq \tau\}.$$

From a data reweighting perspective, filtering data points is equivalent to assigning binary weights of the form $w(x) = \mathbf{1}[s(x) \geq \tau] \in \{0, 1\}$.

- **Data Mixing (domain-level fractional weights)** Data mixing can be seen as a coarse-grained data reweighting method that operates at the domain level. In other words, data points within the same domain receive the same quality score. Pulling all domain scores $\{s_d\}_{d \in \text{domains}}$ and normalize them to obtain data mixing probability for each domain d

$$w_d = \frac{g(s_d)}{\sum_{d'} g(s_{d'})},$$

where g transform the score to be non-negative with a common choice being $g: s \mapsto \exp(s)$. As such, w_d determine how much of a total training budget B is allocated to each domain: $B_d = w_d B$. Alternatively, training can be implemented by sampling domains with probability w_d (probability mixing) or by assigning a fixed quota B_d (quota mixing). Since w_d is decided before training, mixing is an *offline* operation that alters the effective number of examples each domain contributes.

3.3 ONLINE DATA CURATION

In this setting, the scoring function depends on the evolving model state θ_t and, in some cases, on a dynamic validation set. The resulting weights $\{w_t(x)\}_{x \in \mathcal{D}_{\text{train}}}$ evolve throughout training, adjusting the *gradient contributions* of examples while preserving the full volume of the training set $\mathcal{D}_{\text{train}}$.

Data Reweighting (sample-level fractional weights) Data reweighting assigns normalized weights to examples while keeping the training set size unchanged. Similar to Sec. 3.2, the quality score for example x is transformed to a non-negative weight: $w(x) = g(s(x))$, which scales their contribution to the loss:

$$\mathcal{L}^*(\theta) = \frac{1}{Z} \sum_{x \in \mathcal{D}} w(x) \mathcal{L}(\theta; x), \quad Z = \sum_{x \in \mathcal{D}} w(x).$$

Equivalently, from a stochastic gradient descent perspective, the stochastic gradient w.r.t. x is scaled by $\frac{w(x)}{Z}$, which plays a role of *per-sample* learning rate in addition to the *global* learning rate η :

$$\theta_{t+1} = \theta_t - \eta \frac{w(x)}{Z} \nabla_{\theta} \ell(\theta_t; x).$$

Since the total number of training examples remains unchanged, only their relative contribution to parameter updates is modified, data reweighting is inherently an *online* method.

3.4 UNIFIED EVALUATION OF TOTAL FLOPS (OFFLINE VS. ONLINE)

The unification enables us to compare different methods under the same FLOPS calculation framework. We evaluate curation methods under a *cost-aware* metric that accounts for all computation spent to obtain and use the curated data. Let B denote a training budget (tokens or update steps).

Total FLOPs for offline data curation Offline curation modifies the *amount of data that participates in training* before training starts. Its total compute is

$$F_{\text{total}}^{\text{off}} = \underbrace{F_{\text{prep}}^{\text{off}}(\mathcal{D}; f)}_{\text{data scoring as preprocessing}} + \underbrace{F_{\text{train}}(\mathcal{D}', B')}_{\text{model training on the sampled subset}}, \quad (5)$$

where (i) $F_{\text{prep}}^{\text{off}}$ includes any corpus-wide scoring, proxy-model passes, retrieval or filtering necessary to construct \mathcal{D}' ; (ii) $\mathcal{D}' \subseteq \mathcal{D}$ is the retained subset (for selection), or an *effective* subset induced for data mixing; (iii) B' is the effective training budget after curation. By construction, offline selection uses *binary* example weights (keep/drop), and offline domain mixing changes per-domain sampling rates, thereby reducing or reallocating the volume of data seen during training.

Total FLOPs for online data curation Online curation *does not change the amount of training data*; instead, it modulates each example’s contribution *during training*. Its total compute is

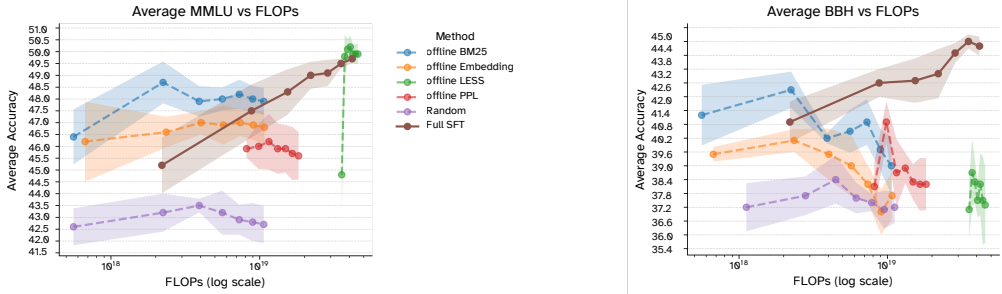
$$F_{\text{total}}^{\text{on}} = \underbrace{F_{\text{train}}(\mathcal{D}, B)}_{\text{model training on the full corpus}} + \underbrace{F_{\text{metrics}}^{\text{on}}}_{\text{on-the-fly features}}, \quad (6)$$

where $F_{\text{metrics}}^{\text{on}}$ denotes lightweight, in-training computations needed to obtain scores (e.g., using current logits, per-example loss, or gradient norms). Because scoring is amortized inside the training loop and no data are removed, $F_{\text{metrics}}^{\text{on}}$ is typically modest relative to F_{train} .

Conversely, online reweighting maintains the full training signal while adding only minimal in-loop metrics computation, providing a clearer accuracy–compute trade-off when measured by *total FLOPs*.

4 REVISITING OF OFFLINE DATA SELECTION

In this section, based on the evaluation protocol derived from the unification above, we revisit the accuracy–FLOPs trade-off of existing offline data selection methods.



(a) In-domain evaluation: MMLU validation set \rightarrow MMLU benchmark. (b) Out-of-domain evaluation: MMLU validation set \rightarrow BBH benchmark.

Figure 1. Efficiency–accuracy tradeoff of different data selection methods and full finetune method under our proposed **total FLOPs** metric. Selection methods (e.g., LESS) appear competitive under in-domain evaluation, but their advantage vanishes or even reverses in out-of-domain settings. In contrast, (full sft) remains consistently strong across domains.

Observation. Figure 1 illustrates the limitation of offline data selection methods: **Overfitting to the validation benchmark**. Model trained on offline selected data with MMLU (Hendrycks et al., 2021) as validation set performs well on the MMLU benchmark (Figure 1a), but the same model generalize poorly to BBH benchmark (Suzgun et al., 2022) (Figure 1b). This issue is especially severe when using LESS (Xia et al., 2024). This reveals that offline methods often *overfit to the chosen validation task*, lacking a true generalization. However, we observed that vanilla full dataset training demonstrate more stable performance in both benchmark. One possible reason is that official data curation which we directly change the number of data involved in the training by repetition would cause the model to replace generalization ability with memorization (Hernandez et al., 2022).

In summary, the offline paradigm has severe limitations, being neither cost-effective nor generalizable. This motivates us to explore an online alternative, which may alleviate memorization while retaining more data to enhance generalization. In the next section, we leverage these insights to design an *online reweighting method*.

5 ADAPTIVE LEARNING RATE FOR ONLINE DATA REWEIGHTING

The limitations of offline data selection suggest that a different paradigm is needed: one that is generalizable, and adaptive to the evolving state of the model. To this end, we propose an **online data reweighting framework** that is also suitable for data selection, data mixing with a unified formalization. Unlike offline methods that commit to a fixed subset before training begins, our framework dynamically adjusts data weights as the model learns, so it naturally adapts to new model states without re-running the data preprocessing and training pipeline from scratch.

5.1 PER-SAMPLE LEARNING RATE UPDATE

We cast data selection and mixing as *online data reweighting*. At training step t , given a minibatch $B_t \subset \mathcal{D}$, we update model parameters as

$$\theta_{t+1} = \theta_t - \eta \sum_{i \in B_t} w_t(i) \nabla_{\theta} \ell(f_{\theta}(x_i), y_i), \quad (7)$$

where $w_t(i) \geq 0$ denotes a dynamic weight assigned to sample i . Here, $w_t(i)$ directly scales the gradient contribution of sample i , and thus acts as a *per-sample learning rate multiplier*: larger weights amplify the effective step size on informative examples, while smaller weights downweight less useful ones.

In our framework, weights are derived from scoring functions $s_t(v, i)$ that compare a training example i with validation/query points $v \in \mathcal{D}_{\text{val}}$. These scores are aggregated across v , normalized within the current batch or pool, and transformed by a smooth gating function to produce the final weights $\widehat{w}_t(i)$. To ensure stability, weights are clipped to prevent excessively large effective learning rates and avoid gradient explosion. We now describe our score functions.

5.2 MODEL-AGNOSTIC PER-SAMPLE LEARNING RATE UPDATE

ADAPT-BM25 quantifies term-based overlap between training and validation data through sparse retrieval scoring:

$$s_{\text{BM25}}(x) = \frac{1}{|\mathcal{D}_{\text{val}}|} \sum_{v \in \mathcal{D}_{\text{val}}} \text{BM25}(x, v). \quad (8)$$

This metric captures surface-level textual similarity without semantic understanding. We use the standard BM25 score $s_{\text{BM25}}(v, i)$ over sparse token matches between query v and example i . Aggregated scores are normalized to obtain target weights.

5.3 PER-SAMPLE LEARNING RATE UPDATE WITH MODEL STATES

ADAPT measures alignment using the model’s own dense representations rather than a frozen encoder (Iverson et al., 2025). For an input x with last-layer hidden states $\{h_i\}_{i=1}^L$, we compute a position-weighted mean pooling:

$$w_i = \frac{i}{\sum_{j=1}^L j}, \quad \phi(x) = \sum_{i=1}^L w_i h_i, \quad (9)$$

where later tokens receive higher weights to counteract the causal mask bias of decoder-only models. We then define the similarity score:

$$s_{\text{ADAPT}}(x) = \frac{1}{|\mathcal{D}_{\text{val}}|} \sum_{v \in \mathcal{D}_{\text{val}}} \cos(\phi(x), \phi(v)). \quad (10)$$

Equivalently, we can instantiate a representation-based scorer $s_{\text{ADAPT}}(v, i) = \cos(\phi(v), \phi(i))$, where $\phi(\cdot)$ is the weighted hidden representation.

We use a sigmoid function with temperature scaling to produce *absolute* weights that are independent of the batch composition:

$$w_t(i) = \sigma \left(\frac{s_{\text{ADAPT}}(x_i)}{\max(\tau, \epsilon)} \right) = \frac{1}{1 + \exp(-s_{\text{ADAPT}}(x_i)/\max(\tau, \epsilon))} \quad (11)$$

where $\sigma(\cdot)$ denotes the sigmoid function, $\tau > 0$ is a temperature hyperparameter (default $\tau = 1.0$), and $\epsilon > 0$ is a small constant (e.g., $\epsilon = 10^{-8}$) used for numerical stability. The temperature parameter τ controls the steepness of the sigmoid: larger values produce flatter weight distributions, while smaller values create sharper distinctions between high and low similarity samples.

Global Weighting via Sigmoid: The sigmoid transformation maps similarity scores to absolute weights in the interval $[0, 1]$ without requiring normalization across samples in the batch. This design choice ensures that the weight assigned to each sample x_i depends solely on its similarity score $s_{ADAPT}(x_i)$ relative to the anchor set, rather than its rank within the current mini-batch. Consequently, a sample with a given similarity score receives the same weight regardless of whether it appears in a high-quality or low-quality batch, making the weighting mechanism robust to batch-level variations in data quality. This contrasts with normalized weighting schemes (e.g., softmax normalization) where weights are relative to other samples in the same batch, potentially amplifying or dampening effects based on batch composition.

Embedding Normalization: We apply L2 normalization to embeddings $\phi(x)$ and $\phi(v)$ before computing cosine similarities to ensure scale-invariant similarity measurements. Specifically, for each embedding vector, we compute $\phi(x) \leftarrow \phi(x) / \max(\|\phi(x)\|_2, \epsilon)$ where ϵ prevents division by zero. This normalization step is distinct from weight normalization and serves to standardize the embedding space for similarity computation.

Online Validation Embedding Updates: To maintain alignment with the evolving model representations, anchor embeddings $\{\phi(v)\}_{v \in \mathcal{D}_{val}}$ are refreshed every R training steps via forward passes on the current model parameters θ_t . This online update mechanism ensures that similarity scores reflect the model’s current representation space rather than outdated embeddings from earlier training stages. The refresh interval R is a hyperparameter that balances computational efficiency with representation freshness. The complete algorithm is summarized in Algorithm 1.

While gradient-based methods (e.g., LESS) directly measure influence, gradients significantly fluctuate during early training and are computationally expensive to compute frequently. In contrast, semantic embeddings provide a smoother, more stable signal of relevance that evolves consistently with the model’s representation space, offering a robust guide for reweighting without the noise of instantaneous gradient estimates.

6 EVALUATION

6.1 EXPERIMENTAL DESIGN

Instruction Tuning Models and Data. We use LoRA (Hu et al., 2021) to fine-tune the base model LLAMA-2-7B (Touvron et al., 2023c). Following the experimental setup in Wang et al. (2023), we use the instruction tuning datasets including FLAN V2 (Longpre et al., 2023), CoT (Wei et al., 2022b), DOLLY (Conover et al., 2023) and OPEN ASSISTANT 1 (Köpf et al., 2023). The datasets do not contain any obvious in-domain data for the target queries. We evaluate our method on MMLU (Hendrycks et al., 2020), multiple-choice dataset spanning 57 tasks and BBH (Suzgun et al., 2023) from BIG-Bench selected to evaluate reasoning capabilities. For MMLU, we report 5-shot accuracy; for BBH, we report 3-shot exact match score. Appendix G contains more details on the training hyperparameter, Appendix D and E contain more dataset details.

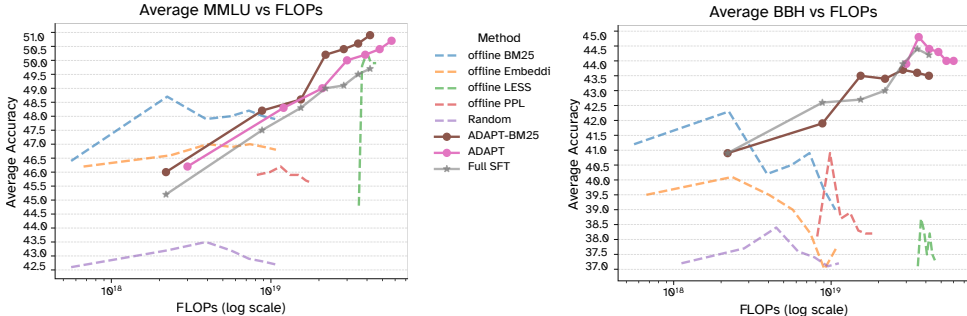
Instruction Tuning Baselines. We introduce instruction finetuning baselines in Sec. 2.2.

Pretraining Models and Data. We adopt TinyLlama architecture (Zhang et al., 2024) with 120M parameters, with FlashAttention (Dao et al., 2022) and Lit-GPT (LightningAI, 2023). More details about the training settings can be found at Appendix F. Following prior works (Touvron et al., 2023b; Zhang et al., 2024; Wettig et al., 2024; Xie et al., 2023b), we employ SlimPajama (Touvron et al., 2023b; Computer, 2023) as the text corpus, which is specifically curated for pre-training LLMs. All selections are performed on about 590M training files of SlimPajama, processed with Llama tokenizer (Touvron et al., 2023b). We evaluate our method on a diverse set of 15 downstream benchmarks, following common practice in prior work such as RegMix (Liu et al., 2024). These tasks span a wide range of realistic settings, including: ARC-E (Clark et al., 2018), ARC-C (Clark et al., 2018), COPA (Sarlin et al., 2020), HellaSwag Zellers et al. (2019), Lambada-S (Paperno et al., 2016), Lambada-O (Radford et al., 2019), LogiQA (Liu et al., 2020), MultiRC (Khashabi et al., 2018), Open-

BookQA (Mihaylov et al., 2018), PiQA (Bisk et al., 2020), QQP (Wang, 2018), RACE (Lai et al., 2017), SciQ (Welbl et al., 2017), Social IQA (Sap et al., 2019), WinoGrande (Sakaguchi et al., 2021). The reported accuracy in table 2 is measured in the 0-shot setting scored using the lm-eval-harness evaluation framework (Gao et al., 2024).

Pretraining Baselines. We compare ADAPT with **Uniform** selection and existing file selection methods for LLM pre-training, including **Doremi** (Xie et al., 2023b) and **RegMix** (Liu et al., 2024). ADAPT requires an anchor set typically consists of examples in the evaluation distribution. To construct this set, we sample 50 validation examples from each of eight evaluation benchmarks: ARC-C, COPA, Lambada, MultiRC, PiQA, RACE, SciQ, and Social IQA. For Doremi and Regmix we use the domain weights in (Lu et al., 2023) as the selection ratio of text samples in different domains in our experiment. For Uniform we use the same ratio to sample from each domain. We also compare **LinUpper** (Sow et al., 2025): an online sample reweighting strategy where the sample weight is proportional to the normalized loss but is capped at a predefined α value, ensuring that outliers do not dominate the training process. Due to the large cost of pretraining, for each method we sample 9B unique tokens from SlimPajama, and train for a total budget of 50B tokens, i.e., train for approximately 5.6 epochs.

6.2 RESULTS FOR INSTRUCTION TUNING



(a) In-domain evaluation: MMLU validation set \rightarrow MMLU benchmark. (b) Out-of-domain evaluation: MMLU validation set \rightarrow BBH benchmark.

Figure 2. Efficiency-accuracy tradeoff. (a) While offline BM25 is efficient at low FLOPs, it saturates early. ADAPT dominates in the high-performance regime (higher FLOPs), avoiding the early plateauing of heuristics. (b) Crucially, LESS overfits to the in-domain task, showing regression on OOD tasks (BBH), whereas ADAPT consistently improves generalization on the Pareto frontier.

As shown in Figure 2, our proposed online data reweighting method consistently outperforms existing approaches under comparable computational budgets:

Table 1. Generalization results under different validation-test configurations.

Method	MMLU(val) - MMLU(test)	MMLU(val) - BBH(test)
BM25	48.7 \pm 0.9	42.3 \pm 0.8
Embedding	47.0 \pm 0.6	40.1 \pm 0.5
LESS	50.2 \pm 0.5	38.7 \pm 1.5
PPL	46.2 \pm 1.1	40.9 \pm 0.9
Random	43.5 \pm 0.3	38.4 \pm 1.0
Full Dataset SFT	49.7 \pm 0.2	44.4 \pm 0.3
ADAPT-BM25	50.9 \pm 0.6	43.7 \pm 1.2
ADAPT	50.7 \pm 0.7	44.8 \pm 1.3

Comparison with offline data selection methods. Across all baselines in both Figure 2a and Figure 2b, our approach consistently lies on the Pareto frontier, demonstrating the best trade-off between computational efficiency and final accuracy. At the same FLOPs, our method achieves substantially higher accuracy than the offline data selection SoTA method LESS. To reach the same performance level, our approach requires significantly fewer FLOPs than LESS. Moreover, after reaching this

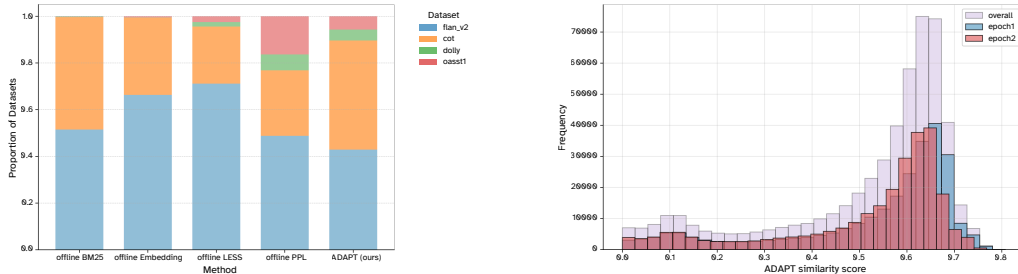
Table 2. Benchmark performance of *TinyLlama-120M* trained on 50B tokens using Uniform, LinUpper, DoReMi, RegMix, and ADAPT.

Tasks	Uniform	LinUpper	DoReMi	RegMix	ADAPT-BM25	ADAPT
ARC-C	17.75 \pm 0.00	18.77 \pm 0.24	18.00 \pm 0.84	18.60 \pm 0.06	18.31 \pm 0.55	19.11 \pm 0.60
ARC-E	39.60 \pm 0.10	38.85 \pm 0.12	40.70 \pm 0.87	41.04 \pm 0.21	39.44 \pm 0.60	39.06 \pm 0.54
COPA	63.00 \pm 0.00	60.00 \pm 1.72	66.00 \pm 1.95	61.00 \pm 1.5	63.20 \pm 0.40	64.00 \pm 1.32
HellaSwag	28.39 \pm 0.00	27.90 \pm 0.13	28.03 \pm 0.03	27.63 \pm 0.03	27.96 \pm 0.07	28.11 \pm 0.00
Lambada-O	24.68 \pm 0.64	23.64 \pm 0.28	22.38 \pm 0.73	24.82 \pm 0.94	26.19 \pm 0.30	24.63 \pm 0.53
Lambada-S	16.98 \pm 0.81	16.79 \pm 0.91	16.30 \pm 1.05	18.30 \pm 0.14	18.43 \pm 0.33	18.07 \pm 0.79
LogiQA	20.43 \pm 0.54	21.35 \pm 0.03	21.51 \pm 0.11	20.89 \pm 0.07	20.52 \pm 0.69	21.66 \pm 0.69
MultiRC	56.68 \pm 0.00	50.74 \pm 0.29	48.99 \pm 2.57	56.70 \pm 0.93	56.53 \pm 0.44	55.67 \pm 1.11
OpenBookQA	15.80 \pm 0.40	14.00 \pm 0.18	15.20 \pm 0.60	18.20 \pm 0.50	15.40 \pm 0.72	14.20 \pm 0.43
PIQA	60.55 \pm 0.29	60.83 \pm 0.45	59.90 \pm 0.39	58.98 \pm 0.29	61.00 \pm 0.19	61.48 \pm 0.23
QQP	36.88 \pm 0.00	37.85 \pm 1.21	36.84 \pm 0.04	36.83 \pm 0.00	36.82 \pm 0.00	36.81 \pm 0.00
RACE	27.85 \pm 0.14	26.41 \pm 0.18	26.41 \pm 0.03	27.85 \pm 0.00	27.39 \pm 0.35	26.60 \pm 0.49
SciQ	71.00 \pm 1.00	70.70 \pm 1.34	72.90 \pm 0.23	71.20 \pm 0.50	72.48 \pm 0.44	72.50 \pm 1.37
SocialIQA	36.90 \pm 0.03	36.85 \pm 0.16	36.95 \pm 0.07	37.41 \pm 0.46	36.52 \pm 0.47	37.05 \pm 0.30
WinoGrande	50.75 \pm 0.17	50.83 \pm 0.09	49.64 \pm 0.79	50.12 \pm 0.72	50.61 \pm 1.25	50.99 \pm 0.00
Average (All) (\uparrow)	37.81 \pm 0.13	37.03 \pm 0.12	37.32 \pm 0.15	37.97 \pm 0.02	38.05 \pm 0.21	38.00 \pm 0.22
Average (Unseen) (\uparrow)	31.98 \pm 0.09	30.56 \pm 0.28	31.95 \pm 0.17	32.46 \pm 0.39	33.49 \pm 0.37	33.73 \pm 0.39
FLOPs overhead (\downarrow)	0	0	4.92×10^{19}	3.072×10^{18}	$\ll 1.0 \times 10^{14}$	$\ll 1.1 \times 10^{15}$

performance, further training with LESS leads to overfitting and accuracy degradation, whereas on-line data reweighting continues to exhibit steady performance gains. Notably, our online reweighting method using BM25 as quality signal (ADAPT-BM25) has better generalization than its offline selection method (BM25) with the other configuration keep the same, which demonstrates the advantage of the online data curation.

Comparison with full-data fine-tuning. In Figure 2a and Table 1, our method also outperforms full-data fine-tuning under equal FLOPs, with the performance gap widening as training progresses. In Fig. 2b and Table. 1, our reweighting strategy that incorporates model state information (ADAPT) outperforms our fixed-state online reweighting (ADAPT-BM25), highlighting the importance of adapting to the evolving model state to better generalisation.

Generalization cross benchmark. As shown in Figure 2b, we further validate the generalization capability of models trained with ADAPT, where we use MMLU (Hendrycks et al., 2021) as the validation set to select the data, and evaluate the model trained on selected data on the BBH (Suzgun et al., 2022) task. The results, also reported in Table 1, show that ADAPT achieves a better generalization performance between benchmark compared to other offline data selection baselines and demonstrates comparable generalization with vanilla full dataset instruction tuning. This highlights a key strength of ADAPT: when new downstream tasks or benchmarks emerge, there is no need to repeat the entire data selection and model training pipeline. Over the long term, this substantially reduces FLOPs costs, making ADAPT highly practical for real-world applications.



(a) Distribution of effective selected examples

(b) Similarity distribution for different epochs

Figure 3. A summary of the distribution of effective selected samples and similarity scores. Both differences and changes show the proposed method’s ability on capturing data features and adaptive learning.

Data mixing and curriculum. We present the distribution of the effective selected examples for different methods in Figure 3a. Different approaches end up selecting very different amounts of data from each training subset. It demonstrates that online reweighting solution can also help with deciding proper data mixture. We also calculate the *effective proportion of used data*, by summing up weights of all data together, with the result of 0.501 in our experimental setting. In practice, this number would automatically adapt according to the overall quality of the training corpus we are selecting from.

We also presented Figure 3b. The shift in similarity distributions from epoch 1 to epoch 2 reflects a transition from collapsed, overly homogeneous representations toward more diverse and fine-grained embeddings. This diversification enhances generalisation by improving the model’s ability to capture subtle distinctions and reducing over-reliance on coarse features. Notably, this process can be viewed as an instance of implicit curriculum learning: the model first clusters samples based on coarse, “easier” patterns and progressively moves toward harder, fine-grained discrimination, thereby mirroring the principles of curriculum learning without explicit scheduling.

6.3 PRETRAINING RESULTS

As shown in Table 2, ADAPT outperforms all baselines in terms of average performance, achieving a 0.19% improvement over Uniform, 0.68% over DoReMi, and a 0.97% over LinUpper. RegMix reaches comparable but slightly lower accuracy on average. Interestingly, LinUpper performs worse than Uniform sampling, suggesting that naive sample-level adjustments may not generalize well. Compared to LinUpper, which conducts normalization and calculates sample-level weights for each batch, our global sample-level reweighting solution demonstrates superior performance in handling unknown ratios of mixed-quality corpora.

To examine the effect of larger training budgets, we extend training to 100B tokens. As shown in Table 3, ADAPT achieves higher benchmark performance than Uniform on 11 out of 15 downstream tasks, with an average improvement of 0.38%. Table 6 further shows that ADAPT consistently yields lower validation perplexity than Uniform at both 50B and 100B training budgets. These findings demonstrate that ADAPT improves both validation perplexity and downstream task performance consistently across compute scales.

Table 3. Benchmark performance of *TinyLlama-120M* trained on 100B tokens using Uniform and ADAPT.

Tasks	Uniform	ADAPT
ARC-C	17.75	18.86
ARC-E	40.66	40.45
COPA	63.00	61.00
HellaSwag	28.35	28.75
Lambada-O	25.13	25.91
Lambada-S	18.47	20.14
LogiQA	20.28	21.66
MultiRC	55.98	56.44
OpenBookQA	16.00	15.00
PIQA	61.48	61.97
QQP	36.81	36.82
RACE	27.37	27.46
SciQ	71.50	72.40
SocialIQA	35.77	37.05
WinoGrande	51.14	51.54
Average (All) (↑)	37.98	38.36
Average (Unseen) (↑)	32.51	32.67

7 CONCLUSION AND DISCUSSION

Our work unifies data selection, mixing, and reweighting under a FLOPs-aware framework. We show that offline pipelines often incur significant cost via data preprocessing and faces issues with generalization. In contrast, online reweighting integrates seamlessly into training, adapts to model state, and improves efficiency without affecting effective data size.

Empirically, we demonstrate that ADAPT consistently outperforms both offline and online baselines across instruction tuning and pretraining setups. In instruction tuning, ADAPT not only achieves higher accuracy under equal FLOPs but also exhibits stronger cross-benchmark generalization (e.g., from MMLU to BBH), mitigating the brittleness of offline methods. In pretraining, ADAPT improves both downstream task accuracy and validation perplexity under 50B and 100B token budgets, underscoring its robustness across compute scales. We further discuss the implications and outline future directions in Sec. A.

ETHICS AND REPRODUCIBILITY STATEMENTS

We adhere to the ICLR Code of Ethics. This paper focuses on comparing methodologies for data curation. To ensure fairness, we evaluate online data reweighting and offline data mixing methods under a FLOPs-equivalent setting. We also rigorously test their performance in both instruction fine-tuning and pretraining scenarios to provide a comprehensive and balanced comparison. During the preparation of this manuscript, we utilized large language models (LLMs) to assist with grammar correction and refinement of the writing. We provide all necessary details to ensure reproducibility of our work. The theoretical justification for our FLOPs-equivalent data curation formalization is presented in Section 3, with detailed FLOP calculations in Appendix C. Implementation details and training protocols are provided in Section 6.1, Appendix B, Appendix F, and Appendix G, while descriptions of the training data are given in Appendix D.

ACKNOWLEDGEMENT

We would like to thank Colin Raffel, Zhihao Zhan, Ziqing Fan, Jonathan Hayase, Guilherme Penedo, Hynek Kydlíček and Loubna Ben allal for many stimulating and helpful discussions. This research was supported by the following entities: The Royal Academy of Engineering via DANTE (a RAEng Chair); the European Research Council, specifically the REDIAL project; SPRIND under the composite learning challenge; Google through a Google Academic Research Award. Yihong acknowledges funding from the Horizon Europe grant No. 101213369 (DVPS).

REFERENCES

- Alon Albalak, Liangming Pan, Colin Raffel, and William Yang Wang. Efficient online data mixing for language model pre-training. In *R0-FoMo: Robustness of Few-shot and Zero-shot Learning in Large Foundation Models Workshop*, 2023.
- Zachary Ankner, Cody Blakeney, Kartik Sreenivasan, Max Marion, Matthew L. Leavitt, and Man-sheej Paul. Perplexed by perplexity: Perplexity-based data pruning with small reference models, 2024. URL <https://arxiv.org/abs/2405.20541>.
- Richard Antonello, Nicole Beckage, Javier Turek, and Alexander Huth. Selecting informative contexts improves language model finetuning. *arXiv preprint arXiv:2005.00175*, 2020.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- Samyadeep Basu, Phil Pope, and Soheil Feizi. Influence functions in deep learning are fragile. In *International Conference on Learning Representations*, 2020.
- Yoshua Bengio, Michael Cohen, Damiano Fornasiere, Joumana Ghosn, Pietro Greiner, Matt MacDermott, Sören Mindermann, Adam Oberman, Jesse Richardson, Oliver Richardson, et al. Superintelligent agents pose catastrophic risks: Can scientist ai offer a safer path? *arXiv preprint arXiv:2502.15657*, 2025.
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 7432–7439, 2020.
- 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, 2020.
- Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- Alexander Bukharin and Tuo Zhao. Data diversity matters for robust instruction tuning. *arXiv preprint arXiv:2311.14736*, 2023.
- Yihan Cao, Yanbin Kang, and Lichao Sun. Instruction mining: High-quality instruction data selection for large language models. *arXiv preprint arXiv:2307.06290*, 2023.

- 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*, 2023a.
- Mayee F Chen, Nicholas Roberts, Kush Bhatia, Jue WANG, Ce Zhang, Frederic Sala, and Christopher Re. Skill-it! a data-driven skills framework for understanding and training language models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023b. URL <https://openreview.net/forum?id=Ioizw01NLf>.
- Xuxi Chen, Zhendong Wang, Daouda Sow, Junjie Yang, Tianlong Chen, Yingbin Liang, Mingyuan Zhou, and Zhangyang Wang. Take the bull by the horns: Hard sample-reweighted continual training improves llm generalization. *arXiv preprint arXiv:2402.14270*, 2024.
- Yihong Chen. *Structure and Destructure: Dual Forces in the Making of Knowledge Engines*. PhD thesis, UCL (University College London), 2025.
- Yihong Chen, Kelly Marchisio, Roberta Raileanu, David Ifeoluwa Adelani, Pontus Stenetorp, Sebastian Riedel, and Mikel Artetxe. Improving language plasticity via pretraining with active forgetting. In *Advances in Neural Information Processing Systems 36*, 2023c.
- Yihong Chen, Luca Franceschi, Xiangxiang Xu, Pontus Stenetorp, and Sebastian Riedel. Jet expansions: Restructuring LLM computation for model inspection. In *The Fourteenth International Conference on Learning Representations*, 2026. URL <https://openreview.net/forum?id=u6JLh0B05h>.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. PaLM: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- 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. *arXiv preprint arXiv:1803.05457*, 2018.
- Together Computer. Redpajama: an open dataset for training large language models, 2023. URL <https://github.com/togethercomputer/RedPajama-Data>.
- Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. Free Dolly: Introducing the world’s first truly open instruction-tuned LLM, 2023.
- Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*, 35:16344–16359, 2022.
- Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. Enhancing chat language models by scaling high-quality instructional conversations. *arXiv preprint arXiv:2305.14233*, 2023.
- 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 (ICML)*, pp. 5547–5569. PMLR, 2022.
- Qianlong Du, Chengqing Zong, and Jiajun Zhang. Mods: Model-oriented data selection for instruction tuning, 2023.
- Simin Fan, Matteo Pagliardini, and Martin Jaggi. DOGE: Domain reweighting with generalization estimation. In *Second Agent Learning in Open-Endedness Workshop*, 2023. URL <https://openreview.net/forum?id=qiKqsqwYXm>.
- Yarin Gal and Stephen Casper. Customizable ai systems that anyone can adapt bring big opportunities—and even bigger risks. *Nature*, 646(8084):286–287, 2025.

- 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, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muenighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot language model evaluation, 2024.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring Massive Multitask Language Understanding. 2020.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *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, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Hamish Ivison, Muru Zhang, Faeze Brahman, Pang Wei Koh, and Pradeep Dasigi. Large-scale data selection for instruction tuning. *arXiv preprint arXiv:2503.01807*, 2025.
- Chunyang Jiang, Chi-min Chan, Wei Xue, Qifeng Liu, and Yike Guo. Importance weighting can help large language models self-improve. *arXiv preprint arXiv:2408.09849*, 2024.
- Youhe Jiang, Ran Yan, Xiaozhe Yao, Yang Zhou, Beidi Chen, and Binhang Yuan. Hexgen: Generative inference of large language model over heterogeneous environment. *arXiv preprint arXiv:2311.11514*, 2023.
- Youhe Jiang, Fangcheng Fu, Xiaozhe Yao, Guoliang He, Xupeng Miao, Ana Klimovic, Bin Cui, Binhang Yuan, and Eiko Yoneki. Demystifying cost-efficiency in llm serving over heterogeneous gpus. *arXiv preprint arXiv:2502.00722*, 2025a.
- Youhe Jiang, Fangcheng Fu, Wanru Zhao, Stephan Rabanser, Jintao Zhang, Nicholas D. Lane, and Binhang Yuan. Cascadia: An efficient cascade serving system for large language models, 2025b. URL <https://arxiv.org/abs/2506.04203>.
- Youhe Jiang, Ran Yan, and Binhang Yuan. Hexgen-2: Disaggregated generative inference of llms in heterogeneous environment. *arXiv preprint arXiv:2502.07903*, 2025c.
- Nikhil Kandpal, Brian Lester, Colin Raffel, Sebastian Majstorovic, Stella Biderman, Baber Abbasi, Luca Soldaini, Enrico Shippole, A Feder Cooper, Aviya Skowron, et al. The common pile v0. 1: An 8tb dataset of public domain and openly licensed text. *arXiv preprint arXiv:2506.05209*, 2025.
- 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. *arXiv preprint arXiv:2001.08361*, 2020a.
- 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, 2020b. URL <https://arxiv.org/abs/2001.08361>.

- Daniel Khashabi, Snigdha Chaturvedi, Michael Roth, Shyam Upadhyay, and Dan Roth. Looking beyond the surface: A challenge set for reading comprehension over multiple sentences. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2018.
- Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi-Rui Tam, Keith Stevens, Abdullah Barhoum, Nguyen Minh Duc, Oliver Stanley, Richárd Nagyfi, et al. Openassistant conversations—democratizing large language model alignment. *arXiv preprint arXiv:2304.07327*, 2023.
- Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. Race: Large-scale reading comprehension dataset from examinations. *arXiv preprint arXiv:1704.04683*, 2017.
- Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. Deduplicating training data makes language models better. In *Annual Meeting of the Association for Computational Linguistics*, 2021. URL <https://api.semanticscholar.org/CorpusID:235829052>.
- Yunshui Li, Binyuan Hui, Xiaobo Xia, Jiayi Yang, Min Yang, Lei Zhang, Shuzheng Si, Junhao Liu, Tongliang Liu, Fei Huang, et al. One shot learning as instruction data prospector for large language models. *arXiv preprint arXiv:2312.10302*, 2023.
- LightningAI. Litgpt. <https://github.com/Lightning-AI/litgpt>, 2023.
- Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. Logiqa: A challenge dataset for machine reading comprehension with logical reasoning. *arXiv preprint arXiv:2007.08124*, 2020.
- Qian Liu, Xiaosen Zheng, Niklas Muennighoff, Guangtao Zeng, Longxu Dou, Tianyu Pang, Jing Jiang, and Min Lin. Regmix: Data mixture as regression for language model pre-training. *arXiv preprint arXiv:2407.01492*, 2024.
- Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. What makes good data for alignment? a comprehensive study of automatic data selection in instruction tuning. *arXiv preprint arXiv:2312.15685*, 2023.
- Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, et al. The flan collection: Designing data and methods for effective instruction tuning. *arXiv preprint arXiv:2301.13688*, 2023.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.
- Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. Chameleon: Plug-and-play compositional reasoning with large language models. In *Proceedings of Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=HtqnVSCj3q>.
- Matthew Macfarlane, Minseon Kim, Nebojsa Jojic, Weijia Xu, Lucas Caccia, Xingdi Yuan, Wanru Zhao, Zhengyan Shi, and Alessandro Sordani. Instilling parallel reasoning into language models. In *2nd AI for Math Workshop @ ICML 2025*, 2025. URL <https://openreview.net/forum?id=a3o4b3hkwp>.
- Max Marion, Ahmet Üstün, Luiza Pozzobon, Alex Wang, Marzieh Fadaee, and Sara Hooker. When less is more: Investigating data pruning for pretraining llms at scale, 2023. URL <https://arxiv.org/abs/2309.04564>.
- Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. *arXiv preprint arXiv:1809.02789*, 2018.

- Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. Orca: Progressive learning from complex explanation traces of gpt-4. *arXiv preprint arXiv:2306.02707*, 2023.
- 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. *arXiv preprint arXiv:1606.06031*, 2016.
- Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. The refinedweb dataset for falcon llm: outperforming curated corpora with web data, and web data only. *arXiv preprint arXiv:2306.01116*, 2023.
- Guilherme Penedo, Hynek Kydlíček, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell, Colin Raffel, Leandro Von Werra, and Thomas Wolf. The fineweb datasets: Decanting the web for the finest text data at scale, 2024. URL <https://arxiv.org/abs/2406.17557>.
- Maxime Peyrard, Sarvjeet Ghotra, Martin Josifoski, Vidhan Agarwal, Barun Patra, Dean Carignan, Emre Kiciman, Saurabh Tiwary, and Robert West. Invariant language modeling. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 5728–5743, 2022.
- 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.
- Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*, 2021.
- 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. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- Stephen Robertson, Hugo Zaragoza, et al. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389, 2009.
- Ohad Rubin, Jonathan Herzig, and Jonathan Berant. Learning to retrieve prompts for in-context learning. *arXiv preprint arXiv:2112.08633*, 2021.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Fevry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M Rush. Multitask prompted training enables zero-shot task generalization. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=9Vrb9D0WI4>.
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. Socialliqa: Commonsense reasoning about social interactions. *arXiv preprint arXiv:1904.09728*, 2019.
- Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superglue: Learning feature matching with graph neural networks. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- Levy Silva and Luciano Barbosa. Improving dense retrieval models with llm augmented data for dataset search. *Knowledge-Based Systems*, 294:111740, 2024.

- 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://cerebras.ai/blog/slimpajama-a-627b-token-cleaned-and-deduplicated-version-of-redpajama>, June 2023. URL <https://huggingface.co/datasets/cerebras/SlimPajama-627B>.
- Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur, Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, Valentin Hofmann, Ananya Harsh Jha, Sachin Kumar, Li Lucy, Xinxu Lyu, Nathan Lambert, Ian Magnusson, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew E. Peters, Abhilasha Ravichander, Kyle Richardson, Zejiang Shen, Emma Strubell, Nishant Subramani, Oyvind Tafjord, Pete Walsh, Luke Zettlemoyer, Noah A. Smith, Hannaneh Hajishirzi, Iz Beltagy, Dirk Groeneveld, Jesse Dodge, and Kyle Lo. Dolma: An Open Corpus of Three Trillion Tokens for Language Model Pretraining Research. *arXiv preprint*, 2024. URL <https://arxiv.org/abs/2402.00159>.
- Daouda Sow, Herbert Woiseschläger, Saikiran Bulusu, Shiqiang Wang, Hans Arno Jacobsen, and Yingbin Liang. Dynamic loss-based sample reweighting for improved large language model pre-training. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=gU4ZgQNs0C>.
- Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. Challenging big-bench tasks and whether chain-of-thought can solve them. *arXiv preprint arXiv:2210.09261*, 2022.
- Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc Le, Ed Chi, Denny Zhou, et al. Challenging big-bench tasks and whether chain-of-thought can solve them. In *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 13003–13051, 2023.
- Liping Tang, Nikhil Ranjan, Omkar Pangarkar, Xuezhi Liang, Zhen Wang, Li An, Bhaskar Rao, Linghao Jin, Huijuan Wang, Zhoujun Cheng, et al. Txt360: A top-quality llm pre-training dataset requires the perfect blend, 2024.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
- Megh Thakkar, Tolga Bolukbasi, Sriram Ganapathy, Shikhar Vashishth, Sarath Chandar, and Partha Talukdar. Self-influence guided data reweighting for language model pre-training. *arXiv preprint arXiv:2311.00913*, 2023.
- 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, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023b.
- 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*, 2023c.
- Alex Wang. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*, 2018.
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, et al. Super-NaturalInstructions: Generalization via declarative instructions on 1600+ NLP tasks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 5085–5109, 2022.

- Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Raghavi Chandu, David Wadden, Kelsey MacMillan, Noah A Smith, Iz Beltagy, et al. How far can camels go? exploring the state of instruction tuning on open resources. In *Neural Information Processing Systems Track on Datasets and Benchmarks (NeurIPS Datasets and Benchmarks)*, 2023.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*, 2022a. URL <https://openreview.net/forum?id=gEzrGCozdqR>.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*, 2022b.
- Johannes Welbl, Nelson F Liu, and Matt Gardner. Crowdsourcing multiple choice science questions. *arXiv preprint arXiv:1707.06209*, 2017.
- Alexander Wettig, Aatmik Gupta, Saumya Malik, and Danqi Chen. Qurating: Selecting high-quality data for training language models. In *Forty-first International Conference on Machine Learning*, 2024.
- Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. LESS: Selecting influential data for targeted instruction tuning. In *International Conference on Machine Learning (ICML)*, 2024.
- 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. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023a. URL <https://openreview.net/forum?id=1XuByUeHhd>.
- 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, 2023b.
- Sang Michael Xie, Shibani Santurkar, Tengyu Ma, and Percy Liang. Data selection for language models via importance resampling. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023c.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. Wizardlm: Empowering large language models to follow complex instructions. *arXiv preprint arXiv:2304.12244*, 2023.
- Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. mt5: A massively multilingual pre-trained text-to-text transformer. *arXiv preprint arXiv:2010.11934*, 2020.
- Ran Yan, Youhe Jiang, Xiaonan Nie, Fangcheng Fu, Bin Cui, and Binhang Yuan. Hexiscale: Accommodating large language model training over heterogeneous environment. *arXiv preprint arXiv:2409.01143*, 2024.
- Ran Yan, Youhe Jiang, Zhuoming Chen, Haohui Mai, Beidi Chen, and Binhang Yuan. Fsa: An alternative efficient implementation of native sparse attention kernel. *arXiv preprint arXiv:2508.18224*, 2025.
- Junjie Oscar Yin and Alexander M. Rush. Compute-constrained data selection, 2024. URL <https://arxiv.org/abs/2410.16208>.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.
- Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. Tinyllama: An open-source small language model. *arXiv preprint arXiv:2401.02385*, 2024.

Wanru Zhao, Vidit Khazanchi, Haodi Xing, Xuanli He, Qionikai Xu, and Nicholas Donald Lane. Attacks on third-party apis of large language models, 2024. URL <https://arxiv.org/abs/2404.16891>.

Wanru Zhao, Yihong Chen, Royson Lee, Xinchu Qiu, Yan Gao, Hongxiang Fan, and Nicholas D. Lane. Breaking physical and linguistic borders: Multilingual federated prompt tuning for low-resource languages, 2025a. URL <https://arxiv.org/abs/2507.03003>.

Wanru Zhao, Hongxiang Fan, Shell Xu Hu, Wangchunshu Zhou, Bofan Chen, and Nicholas D. Lane. Clues: Collaborative high-quality data selection for llms via training dynamics, 2025b. URL <https://arxiv.org/abs/2507.03004>.

Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, LILI YU, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. LIMA: Less is more for alignment. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=KBMOkmX2he>.

APPENDIX

Contents

A Implication and Future Work.	20
B Instruction Tuning Baselines	20
C Data-Selection FLOPs	20
D Training Datasets in Instruction Tuning	21
E Evaluation Datasets in Instruction Tuning	21
F Pretraining Details	21
G Instruction Tuning Details	21
H Pretraining Perplexity	21
I Pseudo Code	22
J Quantitative Analysis	22

A IMPLICATION AND FUTURE WORK.

Implication. ADAPT could be effectively used for getting better control under constrained budgets (both *data budget* and *compute budget*). Our empirical observation provides a principled view and reference for automatically allocating computational resources to the most valuable samples while maintaining end-to-end efficiency. Additionally, we expect our method to perform even better where the data quality is often low, which is the common case in practical setup. where data is crude and is not ideal for direct use for LLM training, sources of generally low quality, and biases inherent to the distribution of content on the web.

Future Work. We didn't explicitly discuss and evaluate deduplication, which is a key step for data preprocessing. For example, we will treat deduplication as reweighting related documents in accordance with their frequency. Based on this framework, we will explore different trade-offs inherent in each stage, such as quality-aware deduplication where high-quality documents are allowed to be duplicated more than once.

In future work, we aim to explore ADAPT for safety-oriented data curation (Gal & Casper, 2025) and to study how dynamic data reweighting influences memorization and internal knowledge structure formation in large language models (Chen, 2025). In particular, combining ADAPT with interpretability techniques such as influence functions (Basu et al., 2020) and Jet Expansions (Chen et al., 2026) may enable us to trace how safety-critical data shapes learned representations and to examine whether models preserve underlying data distributions under dynamic reweighting. More broadly, we view ADAPT as a mechanism for dynamically modulating what information is reinforced during training. Such *data-centric* control of model memory opens avenues not only for improving model quality and safety but also for gaining insights into LLM internal knowledge formation when paired with model-side interpretability methods.

B INSTRUCTION TUNING BASELINES

We compare ADAPT with the following baselines: **1) Random Selection:** We randomly sample data from the instruction tuning dataset. **2) BM25** (Robertson et al., 2009): We assign weights to training samples based on textual statistical features (i.e., TF-IDF), and select the top k data points with the highest scores. **4) PPL** (Yin & Rush, 2024; Marion et al., 2023; Ankner et al., 2024): We compute the loss of each training sample on our original base model and use it as its score. **6) LESS** Xia et al. (2024): We train LoRAs on a random subset of the data, and then selecting data by computing the gradient-based influence of each training sample to validation samples.

C DATA-SELECTION FLOPS

To estimate computational costs throughout our paper, we adopt the methodology of Kaplan et al. (2020b), which approximates the training step computation as approximately $6N$ FLOPs per processed token, where N represents the model's parameter count (approximately 7B). According to Kaplan et al. (2020b), the forward pass consumes roughly half the computational resources of the backward pass, yielding an estimate of $2N$ FLOPs per token during sample processing. We employ an approximation of 2,048 tokens per sample, as we limit all samples to this maximum length during both the training and selection phases. Note that in all experimental configurations, we conduct full fine-tuning of models over two complete epochs. If we define N as model size, P as the data pool magnitude (measured in sample quantity), and D as the number of samples chosen for training, we can calculate the computational expense for each methodology as follows:

1. **Random Selection:** $2k \times 6N \times D \times E$
2. **BM25:** $2k \times 6N \times D \times E$
3. **Embedding:** $2k \times 2N \times P + 2k * 6N \times D \times E$ (embed model N)
4. **PPL:** $2k \times 2NP + 2k \times 6N \times D \times E$
5. **LESS:** $1.53 \times 2k \times 6N \times P + 2k \times 6N \times D \times E$ (LESS computes gradients for three checkpoints over the entire pool.)

D TRAINING DATASETS IN INSTRUCTION TUNING

Table 4 contains information about the training sets used in instruction tuning.

Table 4. Details of training dataset from Wang et al. (2023). Len. is short for token length.

Dataset	# Instance	Sourced from	# Rounds	Prompt Len.	Completion Len.
FLAN V2	100,000	NLP datasets and human-written instructions	1	355.7	31.2
CoT	100,000	NLP datasets and human-written CoTs	1	266	53.2
DOLLY	15,011	Human-written from scratch	1	118.1	91.3
OPEN ASSISTANT 1	55,668	Human-written from scratch	1.6	34.8	212.5

E EVALUATION DATASETS IN INSTRUCTION TUNING

Table 5 contains detailed statistics of the evaluation datasets used in instruction finetuning.

Table 5. Statistics of evaluation datasets. The selection of evaluation tasks cover different kinds of answer types.

Dataset	# Shot	# Tasks	$ \mathcal{D}_{\text{val}} $	$ \mathcal{D}_{\text{test}} $	Answer Type
MMLU	5	57	285	18,721	Letter options
BBH	3	23	69	920	COT and answer

F PRETRAINING DETAILS

We follow all settings in TinyLlama (Zhang et al., 2024). The optimizer is AdamW (Loshchilov & Hutter, 2019), setting parameters β_1 at 0.9 and β_2 at 0.95. We adopt the cosine learning rate schedule with a maximum learning rate of $4e-4$ and the minimum of $4e-5$, the batch size of 2M tokens, the weight decay of 0.1, and the gradient clipping threshold of 1.

G INSTRUCTION TUNING DETAILS

All experiments were conducted with parameter-efficient finetuning method LoRA (Hu et al., 2021). For the LoRA adapter, we specified a rank of 128, an α value of 512, and a dropout rate of 0.1 and applied it across all attention matrices. Adding the LoRA adapter introduce minimal FLOPs overhead during training—having no impact on our FLOPS analysis—and mainly reduce memory requirements for more accessible training.

H PRETRAINING PERPLEXITY

In Table 6, we evaluate validation perplexity on SlimPajama for TinyLlama-120M under two compute budgets (50B and 100B tokens). Across both settings, ADAPT achieves consistently lower perplexity than the Uniform baseline, highlighting its effectiveness in improving data efficiency.

Table 6. Validation perplexity of TinyLlama-120M on SlimPajama under different training budgets.

Method	50B Train Budget	100B Train Budget
Uniform	17.15	15.94
ADAPT	16.55	15.36

I PSEUDO CODE

Algorithm 1: Online Data Reweighting via Embedding Similarity

Input : Training set $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$, validation set $\mathcal{D}_{val} = \{v_j\}_{j=1}^{|\mathcal{D}_{val}|}$, model f_θ , temperature τ , learning rate η , refresh interval R , training steps T , small constant $\epsilon > 0$

Output : Updated parameters θ_T

- 1 Initialize model parameters θ_0
- 2 **for** $t = 1$ **to** T **do**
- 3 Sample mini-batch $B_t \subset \mathcal{D}$
- 4 \triangleright Refresh anchor embeddings every R steps via forward pass on current model θ_t (online update)
- 5 **if** $t \bmod R = 1$ **or** $t = 1$ **then**
- 6 **foreach** $v \in \mathcal{D}_{val}$ **do**
- 7 Forward pass: $(h_1^{(v)}, \dots, h_{L_v}^{(v)}) \leftarrow f_{\theta_t}(v)$
- 8 Position weights: $w_\ell \leftarrow \ell / \sum_{j=1}^{L_v} j$ for $\ell = 1, \dots, L_v$
- 9 $\phi(v) \leftarrow \sum_{\ell=1}^{L_v} w_\ell h_\ell^{(v)}$
- 10 $\phi(v) \leftarrow \phi(v) / \max(\|\phi(v)\|_2, \epsilon)$
- 11 **end**
- 12 \triangleright Forward propagation for training samples
- 13 **foreach** $(x_i, y_i) \in B_t$ **do**
- 14 Forward pass: $(\hat{y}_i, h_1^{(i)}, \dots, h_{L_i}^{(i)}) \leftarrow f_{\theta_t}(x_i)$
- 15 Position weights: $w_\ell^{(i)} \leftarrow \ell / \sum_{j=1}^{L_i} j$ for $\ell = 1, \dots, L_i$
- 16 $\phi(x_i) \leftarrow \sum_{\ell=1}^{L_i} w_\ell^{(i)} h_\ell^{(i)}$
- 17 $\phi(x_i) \leftarrow \phi(x_i) / \max(\|\phi(x_i)\|_2, \epsilon)$
- 18 **end**
- 19 \triangleright Compute quality scores
- 20 **foreach** $i \in B_t$ **do**
- 21 $s_{ADAPT}(x_i) \leftarrow \frac{1}{|\mathcal{D}_{val}|} \sum_{v \in \mathcal{D}_{val}} \cos(\phi(x_i), \phi(v))$
- 22 **end**
- 23 \triangleright Compute per-sample weights
- 24 **foreach** $i \in B_t$ **do**
- 25 $w_t(i) \leftarrow \sigma(s_{ADAPT}(x_i) / \max(\tau, \epsilon))$
- 26 **end**
- 27 \triangleright Compute weighted loss
- 28 $\mathcal{L}_t \leftarrow \sum_{i \in B_t} w_t(i) \ell(\hat{y}_i, y_i)$
- 29 \triangleright Backward propagation and parameter update
- 30 $\theta_{t+1} \leftarrow \theta_t - \eta \sum_{i \in B_t} w_t(i) \nabla_{\theta_t} \ell(\hat{y}_i, y_i)$
- 31 **end**
- 32 **return** θ_T

J QUANTITATIVE ANALYSIS

We provide a detailed qualitative comparison in Table 7, Table 8 and Table 9.

Table 7. Examples with highest and lowest scores retrieved by BM25 for the MMLU validation set.

MMLU Validation Set Examples	
<p>User: Which expression is equivalent to 5×9? A. $(5 \times 4) \times (6 \times 5)$ B. $(5 \times 5) + (5 \times 4)$ C. $(5 \times 5) + (5 \times 9)$ D. $(5 \times 9) \times (6 \times 9)$. Answer: Assistant: B</p> <p>User: Let A be the set of all ordered pairs of integers (m, n) such that $7m + 12n = 22$. What is the greatest negative number in the set $B = \{m + n : (m, n) \in A\}$? A. -5 B. -4 C. -3 D. -2 Answer: Assistant: B</p> <p>User: Ten students take a biology test and receive the following scores: 45, 55, 50, 70, 65, 80, 40, 90, 70, 85. What is the mean of the students' test scores? A. 55 B. 60 C. 62 D. 65 Answer: Assistant: D</p>	
Retrieved by BM25 (Highest similarity)	Retrieved by BM25 (Lowest similarity)
<p>Sample Index: cot 60736 Score: 0.9403 User: Question: A multiple choice test consists of 4 questions, and each question has 5 answer choices. In how many ways can the test be completed if every question is unanswered? Options: (A) 24 (B) 120 (C) 625 (D) 720 (E) 1024 Answer: 5 choices for each of the 4 questions, thus total of $5 \times 5 \times 5 \times 5 = 5^4 = 625$ ways. Assistant: The answer is (C).</p> <p>Sample Index: cot 60736 (same entry, continued) Score: 0.9403 User: Question: In how many different ways can the letters of the word "CLAIM" be rearranged? Options: (A) 110 (B) 120 (C) 125 (D) 130 (E) 140 Answer: The total number of arrangements is $5! = 120$. Assistant: The answer is (B).</p> <p>Sample Index: cot 60736 (same entry, continued) Score: 0.9403 User: Question: A pipe can fill a cistern in 20 minutes whereas the cistern when full can be emptied by a leak in 70 minutes. When both are opened, find when the cistern will be full? Options: (A) 17 minutes (B) 87 minutes (C) 70 minutes (D) 28 minutes (E) 11 minutes Answer: $1/20 - 1/70 = (7 - 2)/140 = 5/140 = 1/28$ 28 minutes. Assistant: The answer is (D).</p> <p>Sample Index: cot 60736 (same entry, continued) Score: 0.9403 User: Question: What is the area of a rectangle whose length is twice its width and whose perimeter equals that of a square with area 1? Options: (A) $1/9$ sq units. (B) $4/9$ sq units. (C) $8/9$ sq units. (D) $7/9$ sq units. (E) $6/9$ sq units. Answer: Perimeter of square = 4, so perimeter of rectangle = 4 $\rightarrow 2(l + w) = 4 \rightarrow l + w = 2$. Since $l = 2w \rightarrow 3w = 2 \rightarrow w = 2/3, l = 4/3$. Area = $(4/3) \times (2/3) = 8/9$. Assistant: The answer is (C).</p>	<p>Sample Index: oasst1_41616 Score: 0.0 User: ทำไมภาษาลาวคล้ายกับภาษาไทย Assistant: ภาษาลาวและภาษาไทยมีความคล้ายคลึงกันเนื่องจากเป็นภาษาในกลุ่มภาษาไทย-ลาว ซึ่งเป็นกลุ่มภาษาในภูมิภาคอินโดจีน ซึ่งประกอบด้วยภาษาที่มีลักษณะทางการเขียนเหมือนกันและมีวามคล้ายคลึงกันในระบบไวยากรณ์และการออกเสียง ารที่ภาษาลาวและภาษาไทยมีความคล้ายคลึงกันเป็นเรื่องธรรมชาติเนื่องจากมีประวัติศาสตร์และวัฒนธรรมที่คล้ายคลึงกัน โดยทั้งสองภาษามีผลมาจากการได้รับอิทธิพลจากภูมิภาคและประวัติศาสตร์เดียวกัน นอกจากนี้ การมีการติดต่อกันและการแลกเปลี่ยนทางวัฒนธรรมในประวัติศาสตร์ระหว่างประเทศลาวและไทยก็ส่งผลให้ภาษาลาวและภาษาไทยมีความคล้ายคลึงกันในด้านของภาษาด้วย</p> <p>Sample Index: oasst1_4573 Score: 0.0 Question: Buenos días Answer: Hola buenos días. ¿En qué te puedo ayudar?</p> <p>Sample Index: oasst1_42983 Score: 0.0 Question: 第一次去深圳旅游, 给我安排一下行程, 大概3天的时间。 Answer: 深圳是中国广东省的一座现代化大城市, 以其高科技、美丽的海滨风光和丰富的购物娱乐场所而闻名。这里为您安排一个3天的深圳旅游行程: [...] (完整3天行程)</p> <p>Sample Index: oasst1_43212 Score: 0.0 Question: 私は猫である名前はまだない Answer: 夏目漱石の「吾輩は猫である」という作品の最初の一節だね。</p>

Table 8. Examples with highest and lowest scores retrieved by LESS (gradient-based) for the MMLU validation set.

MMLU Validation Set Examples	
<p>User: Which expression is equivalent to 5×9? A. $(5 \times 4) \times (6 \times 5)$ B. $(5 \times 5) + (5 \times 4)$ C. $(5 \times 5) + (5 \times 9)$ D. $(5 \times 9) \times (6 \times 9)$. Answer: Assistant: B</p> <p>User: Let A be the set of all ordered pairs of integers (m, n) such that $7m + 12n = 22$. What is the greatest negative number in the set $B = \{m + n : (m, n) \in A\}$? A. -5 B. -4 C. -3 D. -2 Answer: Assistant: B</p> <p>User: Ten students take a biology test and receive the following scores: 45, 55, 50, 70, 65, 80, 40, 90, 70, 85. What is the mean of the students' test scores? A. 55 B. 60 C. 62 D. 65 Answer: Assistant: D</p>	
Retrieved by LESS (Highest similarity)	Retrieved by LESS (Lowest similarity)
<p>Sample Index: flan_v2_78815 Score: 0.9922</p> <p>User: the compound ratio of $2/3$, $6/7$, $4/3$ and $1/8$ is given by? Options: a) $2/21$, b) $1/48$, c) $1/98$, d) $1/96$, e) $1/68$ Assistant: a</p> <p>Sample Index: flan_v2_40013 Score: 1.0</p> <p>User: Look, this is a photo of Mr. White's family. The man in the photo is Mr. White, the father. The woman is the mother. They have two daughters. The girl behind Mr. White is Mary. She's 14. The girl behind Mrs. White is Jane. She is 11. Mary and Jane are in the same school, but not the same grade. Mary is in Grade 2. Jane is in Grade 1. They are good students. How old are Mary and Jane? Options: A) Eleven and thirteen B) Fourteen and eleven C) Eleven D) Thirteen Assistant: B</p> <p>Sample Index: flan_v2_40832 Score: 0.946</p> <p>User: Please answer the following question: Dogs have an understanding of fair play and become angry if they feel that another dog is getting a better deal, a new study has found. The study looked at how dogs react when a companion is rewarded for the same trick in an unequal way. Friederike Range, a researcher at the University of Vienna in Austria, and her colleagues did a series of experiments with dogs who knew how to respond to the command "give the paw". The dogs were normally happy to repeatedly give the paw, whether they got a reward or not. Range found that, compared with monkeys, dogs</p> <p>Options: A care more about whether they are rewarded. B care less about what they are rewarded with. C care more about what they are ordered to do. D care less about who gives them orders. Assistant: B</p>	<p>Sample Index: flan_v2_81227 Score: 0.0</p> <p>User: Personality: I only sleep 3 hours a day. I have no parents or siblings. I'm incredibly smart. [...] I don't blame you. What is your favorite food? I like well done steak. Candidates 1) Hi, how are you doing? 2) Hi, I am Kathie! I don't like eating fast food. 4) I like Mac N cheese. Assistant: I like Mac N cheese.</p> <p>Sample Index: flan_v2_40281 Score: 0.0187</p> <p>User: Where will a book being read by one person be kept? Options: school room, university, bedside table, pantry, discount store I believe the correct choice is "bedside table", here's why: Assistant: a book being read by one person be kept bedside table.</p> <p>Sample Index: flan_v2_78179 Score: 0.0153</p> <p>User: sbjerg forenede Boldklubber (Efb) is a professional Danish football club from the port city of Esbjerg in West Jutland. The club was formed in 1924, as a merger between Esbjerg Boldklub af 1898 and Esbjerg Amatørklub af 1911. The club has training facilities and stadium at Gl. Vardevej in Esbjerg, and plays in blue and white striped shirts. [...] Subject: morten karlsen, Relation: place_of_birth. Options: A 1955 B central C copenhagen D denmark E esbjerg F farum G island (H) jutland (I) of (J) once (K) randers (L) time (M) uerdingen (N) union (O) viborg (P) zealand Assistant: copenhagen</p>

Table 9. Examples with highest and lowest scores retrieved by ADAPT for the MMLU validation set.

MMLU Validation Set Examples	
<p>User: Which expression is equivalent to 5×9? A. $(5 \times 4) \times (6 \times 5)$ B. $(5 \times 5) + (5 \times 4)$ C. $(5 \times 5) + (5 \times 9)$ D. $(5 \times 9) \times (6 \times 9)$. Answer: Assistant: B</p> <p>User: Let A be the set of all ordered pairs of integers (m, n) such that $7m + 12n = 22$. What is the greatest negative number in the set $B = \{m + n : (m, n) \in A\}$? A. -5 B. -4 C. -3 D. -2 Answer: Assistant: B</p> <p>User: Ten students take a biology test and receive the following scores: 45, 55, 50, 70, 65, 80, 40, 90, 70, 85. What is the mean of the students’ test scores? A. 55 B. 60 C. 62 D. 65 Answer: Assistant: D</p>	
Retrieved by ADAPT (Highest similarity)	Retrieved by ADAPT (Lowest similarity)
<p>User: Write something based on this summary: canada foreign minister regrets handling of visa denial for haitian prime minister [...] Assistant: the maximum disturbing and unraveling word for an author’s agent to hear is “no.” so in the ceremony of book publishing, the relationship that is most culturally challenging and polarizing is the one between writers’ agents and publishers. [...]</p> <p>User: Generate a short summary: an airplane en route to algiers from paris made an unscheduled stop thursday in central france after a passenger reported mistakenly bringing on board a suitcase that belonged to someone else [...] (several summarization tasks) Assistant: paris-algiers flight makes unexpected stop in central france [...] (several summarization answers)</p>	<p>User: Vale Assistant: Dime Tio, que necessitas?</p> <p>User: Write a sentence not in English. Assistant: MDN DOM 195 Chef - Service d’examen Vérification Description:</p> <p>User: Kiiltoraunioinen Translate to English Assistant: Asplenium adiantum-nigrum</p>