

A Retraining-Free Framework for Analyzing Training Order Effects in Large Language Models

Anonymous ACL submission

Abstract

The order of training samples plays a crucial role in large language models (LLMs), significantly impacting both their external performance and internal learning dynamics. Traditional methods for investigating this effect generally require retraining the model with various sample orders, which is computationally infeasible for LLMs. In this work, we improve traditional methods by designing a retraining-free framework. By approximating Adam optimizer updates with first- and second-order Taylor expansions and utilizing random projection methods to store intermediate checkpoints, our framework can efficiently estimate model parameters for arbitrary training sample orders. Next, we apply our framework to two downstream research problems: (1) Training curriculum design for LLMs — We base our retraining-free framework to propose a novel curriculum learning strategy that augments curriculum proposals with estimated model performances, enabling more informed sample scheduling. (2) LLMs’ memorization and generalization effect analysis — We use our retraining-free framework to estimate how the positions of training samples influence LLMs’ capacity for memorization and generalization. We conduct extensive experiments to validate the effectiveness of our retraining-free framework in reproducing the true model performances, and further demonstrate its potential in optimizing LLM training curricula and analyzing the memorization and generalization effects of LLMs.

1 Introduction

The order of training samples is crucial for optimizing large language models (LLMs), primarily due to the inherent nature of batch-based optimization methods (*e.g.*, mini-batch gradient descent) (Xue et al., 2023; Peng et al., 2025). This insight has spurred significant research in areas such as training curriculum design for LLMs, which

strategically schedules training samples to enhance model optimization (Zhang et al., 2025b,a; Campos, 2021; Dai et al., 2025), and LLMs’ memorization and generalization effect analysis (Lesci et al., 2024; Tirumala et al., 2022; Budnikov et al., 2025; Zheng and Jiang, 2022), which investigates how the sequence of sample exposure influences the model’s ability to retain knowledge and generalize effectively. A straightforward strategy to study these problems is to train the target model multiple times with different sample orders, and then observe the results to either select the optimal one or analyze the underlying patterns (Zhang et al., 2018b; Xue et al., 2023; Kim and Lee, 2024). In traditional machine learning, the above strategy is feasible because sample and parameter sizes are typically manageable, and training costs are relatively low (Zhang et al., 2018a; Graves et al., 2017; Mahadevan and Mathioudakis, 2024; Wu et al., 2020). However, in the era of LLMs, this approach becomes impractical due to the massive scale of samples and parameters. This naturally raises a novel and fundamental research question:

Can we estimate the effect of sample ordering on LLM performance without retraining?

Despite its significance, answering this question is challenging. To begin with, a practical strategy for estimating model performance under a target sample order is to first measure the performance for a reference sample order and then infer the target performance by establishing a relationship between these two orders. However, since the target sample order can be arbitrary in an extremely large space, identifying a common basis to effectively bridge the reference and target performances becomes a non-trivial challenge. And then, even if we can successfully identify a common basis for relating different sample orders, efficiently storing this basis also poses a significant challenge, as it may involve a vast number of LLM parameters.

To overcome the above challenges, in this paper, we propose a novel retraining-free framework by approximating the parameter updating process with Taylor expansions (called **FUT** for short). Specifically, we focus on the Adam optimizer and reformulate its update term as a function of the current model parameters. Next, we apply Taylor expansions to derive the relationships between the update terms across different model parameters based on the first- and second-order gradients of the loss function. This formulation establishes the common basis for correlating LLM performance across varying sample orders. Finally, we adopt the Random Projection based on the Johnson-Lindenstrauss (JL) theorem (Venkatasubramanian and Wang, 2011) to store the update terms for all training batches, significantly reducing memory consumption while maintaining accuracy.

Building on the above foundational framework, we further apply it to two specific research problems: (1) Training curriculum design for LLMs. Unlike traditional curriculum learning strategies that rely on human heuristics to determine sample orders, our framework empowers users to select sample orders based on the final model performance. Furthermore, for each sample order, our framework provides performance estimations, enabling users to make more informed decisions. (2) LLMs’ memorization and generalization effect analysis. Unlike previous approaches that assess the impact of sample positioning on memorization and generalization through costly retraining or black-box neural network approximations, our framework offers an efficient and principled method to analyze these capabilities in LLMs.

In summary, the main contributions of this paper can be summarized as follows: (1) We formally define the problem of “estimating the impact of training sample orders on model performance without retraining” for LLMs. (2) To solve the above problem, we propose a principled framework based on Taylor expansions and the Random Projection to efficiently estimate LLM performance for arbitrary sample orders. (3) We apply our framework to two specific applications: training curriculum design for LLMs and LLMs’ memorization and generalization effect analysis to demonstrate its fundamental nature and general applicability. (4) We conduct extensive experiments to demonstrate the effectiveness of our framework in approximating the true performance and validate its potential in addressing the aforementioned applications.

2 Problem Formulation

Suppose we have a training dataset with T batches, denoted as $\mathcal{D}_{tr} = \{B_t\}_{t=0}^{T-1}$, and an LLM M . We begin by training M on \mathcal{D}_{tr} following a reference sample order and obtain the corresponding reference checkpoints¹. Specifically, without loss of generality, we assume the reference sample order is B_0, B_1, \dots, B_{T-1} , with the initial parameters of M represented as θ_0 . After processing each batch B_t , the model parameters are updated from θ_t to θ_{t+1} . Ultimately, we collect the reference checkpoints as $\Theta = \{\theta_t\}_{t=0}^T$. For a new sample order, $B_{l_0}, B_{l_1}, \dots, B_{l_{T-1}}$, where B_{l_t} is the $(t+1)$ th training batch, our problem aims to efficiently derive the model parameters $\{\gamma_t\}_{t=0}^T$, where γ_{t+1} is the model parameter after training batch B_{l_t} , and we set $\gamma_0 = \theta_0$.

Relation with the influence function. The above problem shares similarities with the influence function (Koh and Liang, 2017), as both study the effects of training samples. However, there are fundamental differences: our focus is on understanding the impact of sample ordering, while the influence function primarily examines the effect of removing individual samples. Moreover, our problem is situated within the context of LLMs, demanding efficient storage and management of large-scale model parameters.

Straightforward solutions. To solve the above problem, one is to retrain M using the new sample order $B_{l_0}, B_{l_1}, \dots, B_{l_{T-1}}$ and obtain the model parameters $\{\gamma_t\}_{t=0}^T$ after each batch. Another potential solution treats the sample order as the input to a neural network, with the model parameters as the output. In this way, a neural network could be trained to learn the correlation between the input and output, enabling parameter estimation without full retraining. However, the first solution demands substantial time and computational resources to retrain LLMs, rendering it practically infeasible. For the second solution, the limited availability of input-output pairs makes it difficult for a neural network to accurately learn the correlations, resulting in significantly lower performance.

3 The FUT Framework

To address the limitations of the above straightforward solutions, in this section, we propose a principled retraining-free framework. The core idea of

¹Note that the reference sample order can be arbitrary or chosen based on user preference.

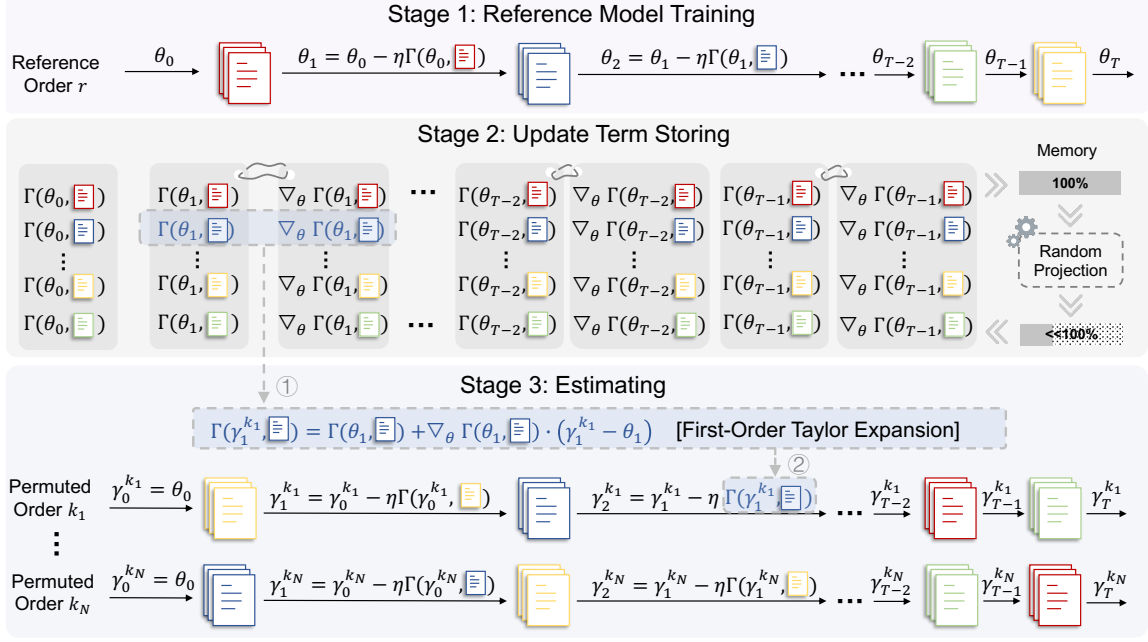


Figure 1: **Overview of the FUT framework.** FUT operates in three stages: **Stage 1:** Compute the reference trajectory $\Theta = \{\theta_t\}_{t=0}^T$ using a fixed data order r . **Stage 2:** Store update and gradient terms for all (θ_t, B_{l_t}) pairs, compressing them via random projection. **Stage 3:** Estimate trajectories $\{\gamma_t^{k_i}\}_{t=0}^T$ under permuted data orders $\{k_i\}_{i=1}^N$ using first-order Taylor expansion based on stored terms. A toy example along the dashed line illustrates: ① retrieving stored terms for expansion, and ② updating parameters along a permuted order.

our approach is to establish a relationship between $\{\gamma_t\}_{t=0}^T$ and $\{\theta_t\}_{t=0}^T$ by delving deeply into their respective generation processes. Then, we derive $\{\gamma_t\}_{t=0}^T$ based on $\{\theta_t\}_{t=0}^T$, which are precomputed as reference checkpoints.

Here, we focus on the Adam optimizer due to its widespread use in LLM optimization. However, our method can be easily extended to other batch-based gradient methods, such as SGD. By applying the updating rule of Adam, we have: ²

$$\theta_{t+1} - \theta_t = -\eta \Gamma(\theta_t, B_t), \forall 0 \leq t \leq T-1 \quad (1)$$

In this equation, $\Gamma(\theta_t, B_t) = m_t / (\sqrt{v_t} + \epsilon)$ is the update term, and

$$\begin{aligned} m_t &= (\beta_1 m_{t-1} + (1 - \beta_1) \nabla_{\theta} \mathcal{L}_t) / (1 - \beta_1^t), \\ v_t &= (\beta_2 v_{t-1} + (1 - \beta_2) \nabla_{\theta} \mathcal{L}_t^2) / (1 - \beta_2^t), \end{aligned} \quad (2)$$

where $\nabla_{\theta} \mathcal{L}_t = \nabla_{\theta} \mathcal{L}(\theta_t, B_t)$ represents the gradient of the loss function \mathcal{L} computed with respect to the model parameters θ_t using the mini-batch B_t . η is the learning rate. m_t and v_t are the first and second momentum statistics, respectively. β_1 and β_2 are both the smoothing coefficients that control

²Without special mention, the updating is applied to each dimension of the parameter separately.

the decay rate of past gradients. ϵ is a small constant to prevent m_t and v_t from being divided by zero.

Similar to the above updating rule, we have $\gamma_{t+1} - \gamma_t = -\eta \Gamma(\gamma_t, B_{l_t})$ ($0 \leq t \leq T-1$). To compute γ_{t+1} , we regard $\Gamma(\theta, B)$ as a function of the model parameters θ . By using Taylor expansions on $\Gamma(\gamma_t, B_{l_t})$, we have:

$$\Gamma(\gamma_t, B_{l_t}) - \Gamma(\theta_t, B_{l_t}) \approx (\gamma_t - \theta_t) \nabla_{\theta} \Gamma(\theta_t, B_{l_t}) \quad (3)$$

where $\nabla_{\theta} \Gamma(\theta_t, B_{l_t})$ represents the gradient of $\Gamma(\theta_t, B_{l_t})$ with respect to θ . In this equation, since B_{l_t} is one of B_0, B_1, \dots, B_{T-1} , if we can obtain $\Gamma(\theta_t, B_{l_t})$ and $\nabla_{\theta} \Gamma(\theta_t, B_{l_t})$ for all $0 \leq t \leq T-1$, then γ_{t+1} can be recursively computed as follows:

$$\gamma_{t+1} = \gamma_t - \eta \Gamma(\theta_t, B_{l_t}) - \eta (\gamma_t - \theta_t) \nabla_{\theta} \Gamma(\theta_t, B_{l_t}), \quad (4)$$

where all the variables on the right-hand side are known. Here, $\Gamma(\theta_t, B_{l_t})$ and $\nabla_{\theta} \Gamma(\theta_t, B_{l_t})$ form the basis for connecting γ_t and θ_t . According to the Adam computational rules, we have:

$$\nabla_{\theta} \Gamma(\theta_t, B_{l_t}) = \frac{\frac{\partial m_t}{\partial \theta} (\sqrt{v_t} + \epsilon) - \frac{\partial \sqrt{v_t}}{\partial \theta} m_t}{(\sqrt{v_t} + \epsilon)^2} \quad (5)$$

$$\text{where } \frac{\partial m_t}{\partial \theta} = \frac{\beta_1 \frac{\partial m_{t-1}}{\partial \theta} + (1 - \beta_1) \nabla_{\theta}^2 \mathcal{L}_{l_t}}{1 - \beta_1^t}, \quad \frac{\partial \sqrt{v_t}}{\partial \theta} =$$

224 $\frac{\beta_2 \frac{\partial v_{t-1}}{\partial \theta} + 2(1-\beta_2)(\nabla_{\theta} \mathcal{L}_{l_t})(\nabla_{\theta}^2 \mathcal{L}_{l_t})}{2(1-\beta_2^t)\sqrt{v_t}}$, $\nabla_{\theta} \mathcal{L}_{l_t}$ and $\nabla_{\theta}^2 \mathcal{L}_{l_t}$
 225 denote $\nabla_{\theta} \mathcal{L}(\theta_t, B_{l_t})$ and $\nabla_{\theta}^2 \mathcal{L}(\theta_t, B_{l_t})$, respec-
 226 tively.

227 By jointly observing equation (2) and (5), we
 228 can see $\Gamma(\theta_t, B_{l_t})$ and $\nabla_{\theta} \Gamma(\theta_t, B_{l_t})$ only rely on
 229 $\nabla_{\theta} \mathcal{L}(\theta_t, B_{l_t})$ and $\nabla_{\theta}^2 \mathcal{L}(\theta_t, B_{l_t})$. These terms are
 230 the gradients of the loss function with respect to the
 231 reference checkpoint and the training batch. Since
 232 the reference checkpoints $\{\theta_t\}_{t=0}^T$ have already
 233 been collected before, we can efficiently compute
 234 $\nabla_{\theta} \mathcal{L}(\theta_t, B_{l_t})$ and $\nabla_{\theta}^2 \mathcal{L}(\theta_t, B_{l_t})$ simply by bringing
 235 θ_t and B_{l_t} into the gradient functions.

236 The algorithm for deriving $\{\gamma_t\}_{t=0}^T$ is shown in
 237 Algorithm 1 in Appendix A.3. In specific, there
 238 are three stages. In the reference model training
 239 stage, we train M using \mathcal{D}_{tr} based on the reference
 240 sample order. After obtaining $\Theta = \{\theta_t\}_{t=0}^T$, in
 241 the update term storing stage, we derive and store
 242 $\Gamma(\theta_t, B_{l_t})$ and $\nabla_{\theta} \Gamma(\theta_t, B_{l_t})$ for all $0 \leq t \leq T-1$
 243 based on equation (2) and (5). At last, in the esti-
 244 mation stage, for a new sample order $\{l_t\}_{t=0}^{T-1}$, we
 245 compute $\{\gamma_t\}_{t=0}^T$ based on equation (4) in a recur-
 246 sive manner. In practice, the first two stages are
 247 executed only once, after which the performance of
 248 any new sample order can be efficiently estimated.
 249 Figure 1 illustrates the complete FUT framework.
 250 **Enhanced model with the second-order Taylor**
 251 **expansion.** In the above method, we approximate
 252 $\Gamma(\gamma_t, B_{l_t})$ with the first-order Taylor expansion.
 253 To enhance accuracy, we extend our approach by
 254 incorporating the second-order term, resulting in
 255 an updated version of equation (3) as follows:

$$256 \Gamma(\gamma_t, B_{l_t}) - \Gamma(\theta_t, B_{l_t}) \approx (\gamma_t - \theta_t) \nabla_{\theta} \Gamma(\theta_t, B_{l_t}) \\
 257 + \frac{1}{2} (\gamma_t - \theta_t)^2 \nabla_{\theta}^2 \Gamma(\theta_t, B_{l_t}), \quad (6)$$

258 where $\nabla_{\theta}^2 \Gamma(\theta_t, B_{l_t})$ is the second-order gradient
 259 of $\Gamma(\theta_t, B_{l_t})$. By combining this equation with
 260 $\gamma_{t+1} - \gamma_t = -\eta \Gamma(\gamma_t, B_{l_t})$, we have:

$$261 \gamma_{t+1} - \gamma_t = -\eta \left((\gamma_t - \theta_t) \nabla_{\theta} \Gamma(\theta_t, B_{l_t}) \right. \\
 262 \left. + \frac{1}{2} (\gamma_t - \theta_t)^2 \nabla_{\theta}^2 \Gamma(\theta_t, B_{l_t}) + \Gamma(\theta_t, B_{l_t}) \right). \quad (7)$$

263 Please referred to Appendix A.1 for more de-
 264 tails to precompute $\nabla_{\theta}^2 \Gamma(\theta_t, B_{l_t})$. After obtain-
 265 ing $\nabla_{\theta}^2 \Gamma(\theta_t, B_{l_t})$, we can efficiently derive $\{\gamma_t\}_{t=0}^T$
 266 based on equation (7) in a recursive manner.

267 **Efficient storage of the update terms.** Accord-
 268 ing to the above analysis, our framework heavily
 269 rely on $\Gamma(\theta_t, B_{l_t})$, $\nabla_{\theta} \Gamma(\theta_t, B_{l_t})$ and $\nabla_{\theta}^2 \Gamma(\theta_t, B_{l_t})$.

270 However, in the context of LLMs, their dimen-
 271 sions are extremely large, posing significant stor-
 272 age challenges. To address this issue, we lever-
 273 age the Random Projection technique (Chen et al.,
 274 2019a; Zhang et al., 2018b) based on the Johnson-
 275 Lindenstrauss (JL) theorem (Venkatasubramanian
 276 and Wang, 2011) to efficiently reduce their dimen-
 277 sionality. To illustrate this process, consider storing
 278 a 2-dimensional matrix $M \in R^{d_1 \times d_2}$. We first gen-
 279 erate a random matrix $A \in R^{d_2 \times k}$ that follows a
 280 Gaussian distribution $\mathcal{N}(0, 1/k)$, where k is the
 281 target dimension chosen based on the JL theorem.
 282 Next, we perform dimensionality reduction by left-
 283 multiplying A with M , that is, $M' = MA$. Here,
 284 $M' \in R^{d_1 \times k}$ is the compressed representation for
 285 storage. To recover the original matrix M , we
 286 similarly perform a left multiplication using the
 287 Moore-Penrose pseudoinverse of A , denoted as
 288 A^+ , that is, $\widetilde{M} = M'A^+$. This approach effec-
 289 tively reduces the space complexity of M from
 290 $\mathcal{O}(d_1 d_2)$ to $\mathcal{O}(d_1 k)$, where $k \ll d$, significantly
 291 alleviating the storage burden when precomputing
 292 it. For higher-order terms such as $\nabla_{\theta} \Gamma(\theta_t, B_{l_t})$
 293 and $\nabla_{\theta}^2 \Gamma(\theta_t, B_{l_t})$, we similarly apply the random
 294 projection technique to reduce their storage com-
 295 plexity, making the process efficient and scalable.

296 **Comparison between the computational costs**
 297 **of retraining and our method.** Assume that the
 298 time complexity for computing the loss gradient
 299 once is $\mathcal{O}(C)$. Enumerating model parameters
 300 under all possible training orders requires retrain-
 301 ing the model on the original dataset for $T!$ times,
 302 where in each permuted order, we need to perform
 303 $\nabla_{\theta} \mathcal{L}(\theta_t, B_{l_t})$ for T times. Therefore, the total time
 304 complexity of retraining is $\mathcal{O}(T \cdot C \cdot T!)$, which
 305 is computationally prohibitive for LLMs with bil-
 306 lions of parameters. In contrast, our method es-
 307 timates the model updates under different batch
 308 orders without retraining. Its main computational
 309 cost comes from computing the updating terms
 310 $\Gamma(\theta_t, B_{l_t})$, $\nabla_{\theta} \Gamma(\theta_t, B_{l_t})$, and $\nabla_{\theta}^2 \Gamma(\theta_t, B_{l_t})$. In spe-
 311 cific, each of these terms requires a single backward
 312 computation of the model at checkpoint θ_t over
 313 batch B_{l_t} , i.e., $\mathcal{L}(\theta_t, B_{l_t})$. Since there are T^2 such
 314 (θ_t, B_{l_t}) pairs in total, the overall time complexity
 315 of our method is $\mathcal{O}(T^2 \cdot C)$.

316 4 Applications

317 4.1 Training Curriculum Design for LLMs

318 **Problem definition.** Following the notations in
 319 Section 2, suppose we aim to train a model M on

the dataset $\mathcal{D}_{tr} = \{B_t\}_{t=0}^{T-1}$. Let π be a permutation function that maps the standard index set $\{0, 1, \dots, T-1\}$ to $\{\pi(0), \pi(1), \dots, \pi(T-1)\}$, where $\pi(t) \in [0, T-1]$ indicates that batch B_t is placed at the $(\pi(t) + 1)$ -th position in the training sequence. Following common practice, we train the LLM for only one epoch (Xue et al., 2023). The goal is to find an optimal permutation π^* such that the resulting model performs best on a validation set \mathcal{D}_{val} , formally defined as:

$$\pi^* := \arg \max_{\pi \in \Pi} \mathcal{R}(\gamma_T^\pi, \mathcal{D}_{val}), \quad (8)$$

where γ_T^π denotes the final model parameters estimated using our FUT framework, and the training order l_t is induced by π . The performance metric \mathcal{R} is implemented using Perplexity (PPL) (Hu et al., 2024), and Π denotes the space of all possible permutation functions.

Our solution based on FUT. Since objective (8) is non-differentiable, we design a Genetic Algorithm (GA) (Katoch et al., 2020) to obtain π^* . In specific, we maintain a set of candidate sample orders and iteratively apply crossover and mutation operators to generate improved sample orders, aiming to optimize the model performance. For more details, we refer readers to Appendix A.4. Compared to traditional curriculum learning strategies, a key advantage of our method is its ability to estimate model performance for each curriculum proposal, enabling more informed decisions. For instance, by knowing the performance gap between different curricula, users can assess whether the difference is significant. If the gap is small, users can confidently choose one at random.

4.2 LLMs’ Memorization and Generalization Effect Analysis

Problem definition. We continue to follow the notations introduced in Section 2. For each training batch in \mathcal{D}_{tr} , the memorization problem evaluates model performance when the batch appears at different positions in the training sequence. Specifically, we use the following evaluation method:

$$M_{i,j} = \frac{1}{N} \sum_{k=1}^N \mathcal{R}(\theta_T^{\pi_k^{ij}}, B_i),$$

where π_k^{ij} is a permutation function that fixes B_i at the j -th training position while randomly shuffling all other batches. For each B_i , we generate N such permutations, and the final performance

is computed as the average across these permutations. The generalization problem is defined in a similar manner, with the key distinction that B_i in the above equation is replaced by D_i , a dataset not seen during training, i.e., $D_i \notin \mathcal{D}_{tr}$.

Our solution based on FUT. For each π_k^{ij} , we first generate the sequence l_t and then estimate $\gamma_T^{\pi_k^{ij}}$ using the reference checkpoints $\{\theta_t\}_{t=0}^T$. Finally, we compute $\mathcal{R}(\theta_T^{\pi_k^{ij}}, B_i)$ or $\mathcal{R}(\theta_T^{\pi_k^{ij}}, D_i)$ based on $\gamma_T^{\pi_k^{ij}}$, and average the resulting performances over different values of k . Compared to previous studies that estimate memorization capability using black-box neural network (Zheng and Jiang, 2022; Lesci et al., 2024; Feldman and Zhang, 2020), our method is more principled and grounded in theoretical foundations.

5 Experiments

In this section, we conduct extensive experiments to demonstrate the effectiveness of our framework and its potential applications in designing LLM training curricula and analyzing LLMs’ memorization and generalization capabilities.

5.1 Evaluation on the General Capability of Our Methods

Experimental Setup. To evaluate the effectiveness of our FUT framework, we incorporate the estimated model parameters into the LLM and measure the performance gap between the estimated and actual results. Specifically, we conduct our experiments on the Wikitext dataset (Merity et al., 2018a), a curated collection of high-quality English Wikipedia articles that is widely used for language modeling and evaluation. This dataset is particularly well-suited for assessing model perplexity due to its long-range token dependencies (Merity et al., 2018b). We adopt the architecture of LLaMA (Touvron et al., 2023) to construct a base model with 636 million parameters. The model has a hidden size of 2048 and consists of 10 stacked transformer layers with 10 attention heads. We choose this relatively small architecture because our main experiments involve repeated LLM training to validate that the proposed FUT framework can accurately estimate model parameters under various training orders. In Appendix C.1, we scale the model size up to 6.0 billion parameters to assess the scalability of our approach. Following common practice (Xue et al., 2023), we use the Adam optimizer for LLM

Table 1: Performance estimation accuracy (use AbsDiff as metric) across methods under settings with different numbers of batches.

T	Random	FUT	FUT++
8	0.0205	0.0165	0.0085
16	0.0917	0.0649	0.0703
32	0.0373	0.0290	0.0193
64	0.0644	0.0445	0.0319
128	0.0575	0.0372	0.0284
256	0.0471	0.0205	0.0368

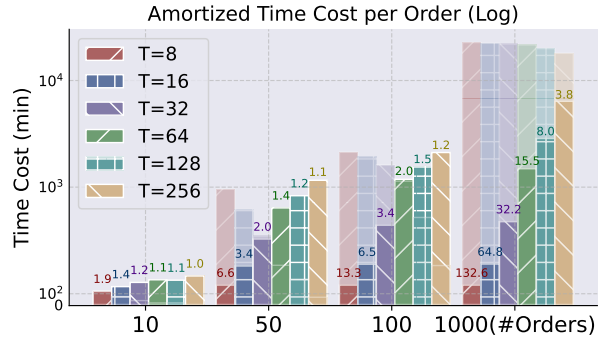


Figure 2: Time cost comparison.

training and train for a single epoch to evaluate performance based on perplexity (Hu et al., 2024). We also report fine-grained performance estimation results at intermediate stages in Appendix C.3.

In our experiments, assuming the dataset consists of T batches, we randomly select N training orders from the total of $T!$ possible permutations. For each selected order, we use our method to estimate the model performance \hat{r} and also train the LLM using that order to obtain the ground-truth performance r . The performance gap is then calculated as: $\text{AbsDiff} = \frac{1}{N} \sum_{k=1}^N |\hat{r}_k - r_k|$, where k indexes the different training orders. We set $N = 10$ to balance evaluation reliability with computational cost. To assess the scalability and robustness of our framework, we vary the number of batches T across the set $\{8, 16, 32, 64, 128, 256\}$. This setup allows us to evaluate how performance estimation behaves under increasing training granularity and longer optimization trajectories.

Baseline. We denote our method using the first- and second-order Taylor expansions as **FUT** and **FUT++**, respectively. The ground-truth performance obtained via actual LLM training is referred to as **Retraining**. Additionally, we introduce a heuristic baseline, named **Random**, where we first obtain all the N ground-truth performances $\{r_k\}_{k=1}^N$, and then randomly estimate the performance within the range $[\min_{k \in [1, N]} r_k, \max_{k \in [1, N]} r_k]$.

Results. Table 1 shows that the Random baseline performs the worst, indicating that estimating LLM performance without retraining itself is a non-trivial task. Both FUT and FUT++ consistently outperform Random across all batch settings with considerable margins, demonstrating their effectiveness. Between our methods, FUT++ performs better than FUT in more cases, suggesting that the inclusion of the second-order term in the Taylor expansion is beneficial for our problem. In addition,

we also compare the efficiency of our method with the Retraining strategy. Here, we vary N over the set $\{10, 50, 100, 1000\}$ to observe the trend as the number of sample orders increases³. We compare different methods with various T 's. The results are presented in Figure 2, where the solid bars represent our method and the dashed bars represent Retraining. We observe that as the total number of orders increases, our method progressively achieves higher time efficiency per order compared to Retraining, with a maximum speedup of 132.6 times. Our methods across all T surpass Retraining, highlighting the advantages of our methods in scalability. We also conduct experiments to demonstrate the necessity of using the random projection for storage in Appendix C.2.

5.2 Evaluation on the Application of Training Curriculum Design for LLMs

Experimental Setup & Baselines. In this experiment, we evaluate whether our methods can assist in designing more effective training curricula for LLMs. Similar to the above section, we use perplexity as the evaluation metric and measure different models by varying T in the range of $\{8, 16, 32, 64, 128, 256\}$. We compare our methods with the following baselines: **Random Order (RO)**, which generates the curriculum by randomly shuffling the training batches. **Sample Length (SL)** (Campos, 2021), which is a difficulty-based curriculum design strategy, and the difficulty score is determined based on the sentence length. **Perplexity (PPL)** (Zhang et al., 2025a), which uses the perplexity from a reference model as a proxy to evaluate sample difficulty and design the curriculum. **Perplexity Difference (PD)** (Zhang et al., 2025b), measuring the perplexity gap between a strong and a weak model, treating samples with

³The costs of the first two stages in our method are amortized across all sample orders, as they are executed only once.

Table 2: Perplexity results across different batch numbers and curriculum design strategies. “Est.” means the performance estimated by our methods, and the colored results represent the best estimation accuracy.

Methods	RO	Len	PPL	PD	FUT (Est.)	FUT++ (Est.)
8	1.4414	1.4392	1.4012	1.4006	1.3996 (1.3963)	1.3998 (1.3962)
16	1.4599	1.5291	1.4531	1.4542	1.4536 (1.4314)	1.4523 (1.4307)
32	1.4109	1.4042	1.3966	1.3933	1.3909 (1.3823)	1.3881 (1.3686)
64	1.4248	1.4079	1.4027	1.4071	1.3785 (1.3838)	1.3804 (1.3856)
128	1.3838	1.3872	1.3790	1.3697	1.3412 (1.3446)	1.3619 (1.3512)
256	1.3696	1.3766	1.3645	1.3660	1.3378 (1.3551)	1.3178 (1.3460)

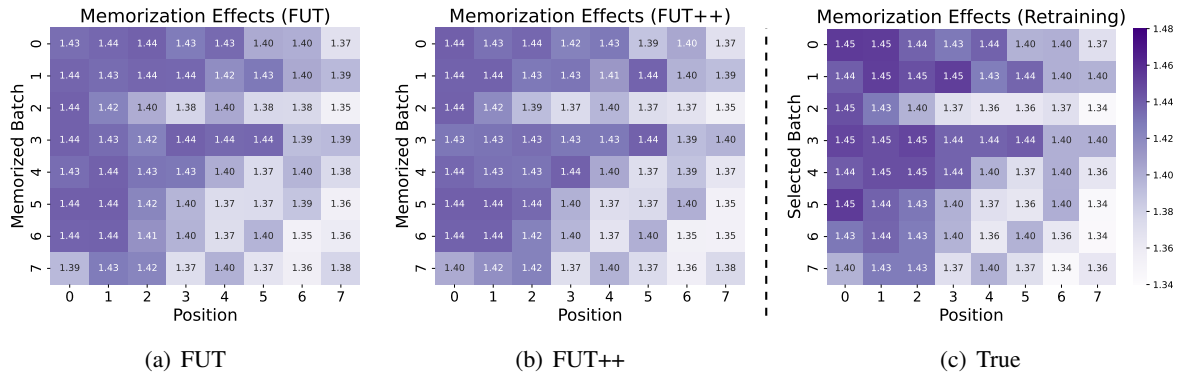


Figure 3: Memorization effects. Heatmaps in (a) and (b) are estimated by our FUT and FUT++ methods, respectively. Heatmap in (c) represents the true memorization effect obtained by retraining.

larger gaps as more difficult to design the curriculum. We use the baseline methods and our proposed approaches (using equation (8)) to generate training curricula, and train the LLM based on them for comparison.

Results. The results are shown in Table 2. We can see: In most cases, **RO** performs the worst, as it lacks any problem-specific design and simply generates the training curriculum randomly. **PPL** and **PD** consistently outperform **SL** across different batch sizes, which is as expected since they both leverage perplexity as a proxy to design the curricula-aligning well with the final evaluation metric. Finally, our methods achieve superior performance compared to all baselines, demonstrating their effectiveness in designing training curricula for LLMs. As shown in the last two columns of Table 2, our methods provide performance estimates that closely match the actual results, enabling more efficient decision-making when selecting optimal training orders.

5.3 Evaluation on the Application of LLM Memorization & Generalization Effect

Experimental Setup. In this experiment, we evaluate the memorization & generalization effects of

LLM when a sample batch is placed at different training positions. In specific, the number of training batches is set as 8 (*i.e.*, $T = 8$). We visualize the value of $M_{i,j}$ in Section 4.2 based on perplexity by setting different (i, j) pairs.

Results. The results are presented in Figure 3, 4 and 5. We can see: Compared to the true memorization effect (in Figure 3(c)), where we retrain the LLM to compute $M_{i,j}$, FUT and FUT++ in Figure 3(a) and (b), can accurately estimate the model’s memorization of different batches at various positions using both first- and second-order approximations, respectively. The results reveal that the model tends to memorize batches appearing later in the training order more effectively, as indicated by lower perplexity.

For generalization analysis, we divide training batches into two groups based on their similarity to the test set D , using the average similarity τ as a threshold. Our metric is the cosine similarity of the sample representations. As shown in Figures 4 and 5, our method (dashed red/blue lines) closely estimates the true performance (black line) and captures the same generalization trend in most cases. In Figure 4, batches similar to the test data generalize better when placed later in training. In contrast,

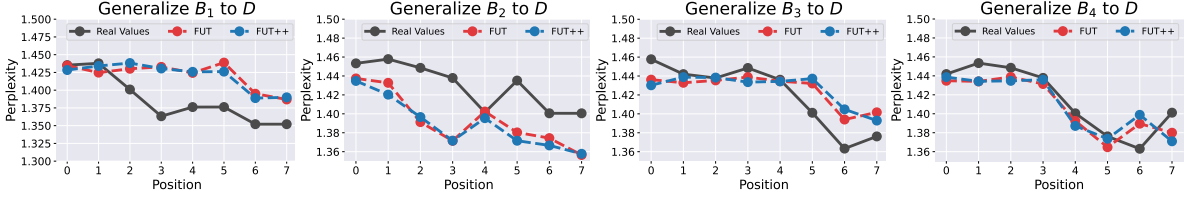


Figure 4: The generalization effect of batch B_i on dataset D , with $\text{sim}(B_i, D) \geq \tau$.

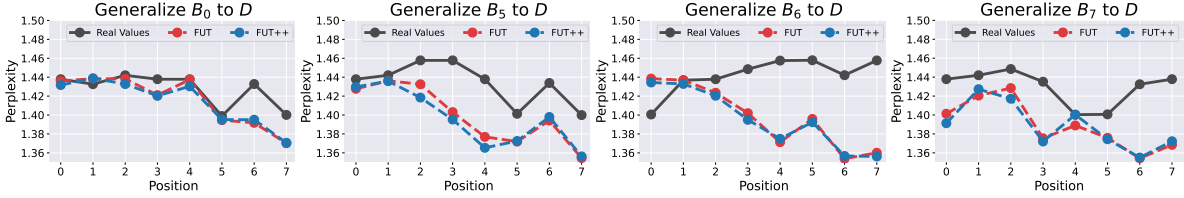


Figure 5: The generalization effect of batch B_i on dataset D , with $\text{sim}(B_i, D) < \tau$.

Figure 5 shows that dissimilar batches have little or random effect on generalization, regardless of their positions in the training sequences.

6 Related Work

Training Dynamics of Language Models. Understanding training dynamics is essential for analyzing how deep models evolve during optimization (Frankle et al., 2020; Raghu et al., 2017; Achille et al., 2018). In the context of language models, early work focused on the evolution of learned representations (Saphra and Lopez, 2018; Saphra, 2021) and the encoding of world knowledge (Liu et al., 2021) during pre-training. These insights have also been extended to downstream tasks such as summarization (Goyal et al., 2022) and speech translation (Savoldi et al., 2022). More recent studies have begun to examine the training dynamics of LLMs (Ren and Sutherland, 2025; Biderman et al., 2023; Teehan et al., 2022; Lesci et al., 2024; Dai et al., 2025), which are harder to analyze due to their scale. For example, Teehan et al. (2022) studies internal representation development and structural changes during training, while Biderman et al. (2023) uses models of varying sizes to study how training behavior shifts with scale. Additionally, Ren and Sutherland (2025) explores how learning certain examples affects the model’s behavior on other inputs.

Influence Function. Influence function is used to estimate the impact of each training sample on a specific test prediction (Koh and Liang, 2017; Bae et al., 2022; Koh et al., 2019). The foundational work by Koh and Liang (2017) applies influence functions by calculating gradients and Hessian-

vector products to measure the contribution of each training example to a test point. However, research in Basu et al. (2021); Guo et al. (2021) has shown that influence functions can be unstable and unreliable in neural network. Additionally, computing the necessary Hessian-vector products is computationally expensive, particularly for LLMs. To address this challenge, a recent study by Lin et al. (2024) introduces a caching mechanism to estimate token-level influences in LLMs. While this method alleviates some computational difficulties, it overlooks the crucial influence of sample order in the training process of LLMs, which plays a significant role in shaping the learning dynamics.

7 Conclusion

In this work, we propose a retraining-free framework FUT for analyzing the effect of training sample order on LLMs, addressing the prohibitive cost of traditional retraining-based approaches. By approximating the optimization dynamics of Adam via Taylor expansion and employing random projection for efficient parameter estimation, our framework enables accurate performance prediction under arbitrary given sample orders. We demonstrate the utility of our FUT framework in two key research problems of LLMs: training curriculum design, and memorization & generalization effect analysis. Extensive experiments show that our framework faithfully approximates true model performance and provides valuable insights into both external performance and internal learning dynamics of LLMs. In summary, our FUT framework offers a practical tool for understanding and optimizing the model behaviors of LLMs.

610 Limitations

611 While our proposed retraining-free framework
612 (FUT) provides a computationally efficient and the-
613oretically grounded method for estimating the ef-
614fects of sample order in LLMs, several limitations
615should be acknowledged.

- 616 1. The accuracy of our estimates relies on the va-
617lidity of Taylor expansions, particularly when
618higher-order nonlinearities dominate the op-
619timization dynamics—scenarios where our
620first- and second-order approximations may
621fall short.
- 622 2. Although the use of random projection signifi-
623cantly reduces memory overhead, it may intro-
624duce approximation noise, especially for mod-
625els with extremely large parameter spaces.
- 626 3. We evaluate the effectiveness of our FUT
627framework solely based on perplexity perfor-
628mance. This is because downstream natural
629language understanding and reasoning tasks
630typically require large-scale models, which
631are infeasible to retrain repeatedly under vary-
632ing conditions. Nevertheless, further valida-
633tion is needed to assess the generalizability of
634our framework in these more complex tasks.

635 Ethical considerations

636 We confirm that this research adheres to the ARR
637Code of Ethics. Our study does not involve human
638subjects, and all datasets used are publicly avail-
639able, with no sensitive or personal data involved.
640We have considered the ethical implications of our
641work, particularly in terms of fairness and potential
642misuse of large language models. We are com-
643mitted to ensuring that our research contributes
644positively to the academic community and society
645at large.

646 References

647 Alessandro Achille, Matteo Rovere, and Stefano Soatto.
648 2018. Critical learning periods in deep networks. In
649 *International Conference on Learning Representa-*
650 *tions*.

651 Juhan Bae, Nathan Ng, Alston Lo, Marzyeh Ghassemi,
652 and Roger B Grosse. 2022. If influence functions are
653 the answer, then what is the question? *Advances in*
654 *Neural Information Processing Systems*, 35:17953–
655 17967.

S Basu, P Pope, and S Feizi. 2021. Influence func-
tions in deep learning are fragile. In *International*
Conference on Learning Representations (ICLR).

Yoshua Bengio, Jérôme Louradour, Ronan Collobert,
and Jason Weston. 2009. Curriculum learning. In
Proceedings of the 26th annual international confer-
ence on machine learning, pages 41–48.

Stella Biderman, Hailey Schoelkopf, Quentin Gregory
Anthony, Herbie Bradley, Kyle O’Brien, Eric Hal-
lahan, Mohammad Aflah Khan, Shivanshu Purohit,
USVSN Sai Prashanth, Edward Raff, and 1 others.
2023. Pythia: A suite for analyzing large language
models across training and scaling. In *International*
Conference on Machine Learning, pages 2397–2430.
PMLR.

Andrei Z. Broder. 1997. On the resemblance and
containment of documents. *Proceedings. Compre-*
sion and Complexity of SEQUENCES 1997 (Cat.
No.97TB100171), pages 21–29.

Mikhail Budnikov, Anna Bykova, and Ivan P
Yamshchikov. 2025. Generalization potential of large
language models. *Neural Computing and Applica-*
tions, 37(4):1973–1997.

Daniel Campos. 2021. Curriculum learning for lan-
guage modeling. *arXiv preprint arXiv:2108.02170*.

Haochen Chen, Syed Fahad Sultan, Yingtao Tian,
Muhao Chen, and Steven Skiena. 2019a. Fast and
accurate network embeddings via very sparse random
projection. In *Proceedings of the 28th ACM inter-*
national conference on information and knowledge
management, pages 399–408.

Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi,
and Cho-Jui Hsieh. 2017. Zoo: Zeroth order op-
timization based black-box attacks to deep neural
networks without training substitute models. In *Pro-*
ceedings of the 10th ACM workshop on artificial
intelligence and security, pages 15–26.

Xiangyi Chen, Sijia Liu, Kaidi Xu, Xingguo Li, Xue
Lin, Mingyi Hong, and David Cox. 2019b. Zo-
adam: Zeroth-order adaptive momentum method
for black-box optimization. *Advances in neural in-*
formation processing systems, 32.

Yalun Dai, Yangyu Huang, Xin Zhang, Wenshan Wu,
Chong Li, Wenhui Lu, Shijie Cao, Li Dong, and
Scarlett Li. 2025. Data efficacy for language model
training. *arXiv preprint arXiv:2506.21545*.

John C Duchi, Michael I Jordan, Martin J Wainwright,
and Andre Wibisono. 2015. Optimal rates for zero-
order convex optimization: The power of two func-
tion evaluations. *IEEE Transactions on Information*
Theory, 61(5):2788–2806.

Vitaly Feldman and Chiyuan Zhang. 2020. What neural
networks memorize and why: Discovering the long
tail via influence estimation. *Advances in Neural*
Information Processing Systems, 33:2881–2891.

711	Abraham D Flaxman, Adam Tauman Kalai, and H Brendan McMahan. 2004. Online convex optimization in the bandit setting: gradient descent without a gradient. <i>arXiv preprint cs/0408007</i> .	Pang Wei Koh and Percy Liang. 2017. Understanding black-box predictions via influence functions. In <i>International conference on machine learning</i> , pages 1885–1894. PMLR.	765
712			766
713			767
714			768
715	Jonathan Frankle, David J Schwab, and Ari S Morcos. 2020. The early phase of neural network training. <i>arXiv preprint arXiv:2002.10365</i> .	Pang Wei W Koh, Kai-Siang Ang, Hubert Teo, and Percy S Liang. 2019. On the accuracy of influence functions for measuring group effects. <i>Advances in neural information processing systems</i> , 32.	769
716			770
717			771
718	Saeed Ghadimi and Guanghui Lan. 2013. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. <i>SIAM journal on optimization</i> , 23(4):2341–2368.	Padmavathi Kora and Priyanka Yadlapalli. 2017. Crossover operators in genetic algorithms: A review. <i>International Journal of Computer Applications</i> , 162(10).	773
719			774
720			775
721			776
722	Henok Ghebrechristos and Gita Alaghband. 2020. Deep curriculum learning optimization. <i>SN Computer Science</i> , 1(5):245.	Pietro Lesci, Clara Meister, Thomas Hofmann, Andreas Vlachos, and Tiago Pimentel. 2024. Causal estimation of memorisation profiles. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15616–15635.	777
723			778
724			779
725	Daniel Golovin, John Karro, Greg Kochanski, Chansoo Lee, Xingyou Song, and Qiuyi Zhang. 2019. Gradientless descent: High-dimensional zeroth-order optimization. <i>arXiv preprint arXiv:1911.06317</i> .	William Liang, Sam Wang, Hung-Ju Wang, Osbert Bastani, Dinesh Jayaraman, and Yecheng Jason Ma. 2024. Environment curriculum generation via large language models. In <i>8th Annual Conference on Robot Learning</i> .	780
726			781
727			782
728			783
729	Tanya Goyal, Jiacheng Xu, Junyi Jessy Li, and Greg Durrett. 2022. Training dynamics for text summarization models. In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 2061–2073.	Huawei Lin, Jikai Long, Zhaozhuo Xu, and Weijie Zhao. 2024. Token-wise influential training data retrieval for large language models. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 841–860.	784
730			785
731			786
732			787
733			788
734	Alex Graves, Marc G Bellemare, Jacob Menick, Remi Munos, and Koray Kavukcuoglu. 2017. Automated curriculum learning for neural networks. In <i>international conference on machine learning</i> , pages 1311–1320. Pmlr.	Sijia Liu, Pin-Yu Chen, Bhavya Kailkhura, Gaoyuan Zhang, Alfred O Hero III, and Pramod K Varshney. 2020. A primer on zeroth-order optimization in signal processing and machine learning: Principals, recent advances, and applications. <i>IEEE Signal Processing Magazine</i> , 37(5):43–54.	789
735			790
736			791
737			792
738			793
739	Liangke Gui, Tadas Baltrušaitis, and Louis-Philippe Morency. 2017. Curriculum learning for facial expression recognition. In <i>2017 12th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2017)</i> , pages 505–511. IEEE.	Zeyu Liu, Yizhong Wang, Jungo Kasai, Hannaneh Hajishirzi, and Noah A Smith. 2021. Probing across time: What does roberta know and when? In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 820–842.	794
740			795
741			796
742			797
743			798
744	Han Guo, Nazneen Rajani, Peter Hase, Mohit Bansal, and Caiming Xiong. 2021. Fastif: Scalable influence functions for efficient model interpretation and debugging. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 10333–10350.	Ananth Mahadevan and Michael Mathioudakis. 2024. Cost-aware retraining for machine learning. <i>Knowledge-Based Systems</i> , 293:111610.	799
745			800
746			801
747			802
748			803
749			804
750	Guy Hacohen and Daphna Weinshall. 2019. On the power of curriculum learning in training deep networks. In <i>International conference on machine learning</i> , pages 2535–2544. PMLR.	Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D Lee, Danqi Chen, and Sanjeev Arora. 2023. Fine-tuning language models with just forward passes. <i>Advances in Neural Information Processing Systems</i> , 36:53038–53075.	805
751			806
752			807
753			808
754	Yutong Hu, Quzhe Huang, Mingxu Tao, Chen Zhang, and Yansong Feng. 2024. Can perplexity reflect large language model’s ability in long text understanding? <i>ArXiv</i> , abs/2405.06105.	Tambet Matiisen, Avital Oliver, Taco Cohen, and John Schulman. 2019. Teacher–student curriculum learning. <i>IEEE transactions on neural networks and learning systems</i> , 31(9):3732–3740.	809
755			810
756			811
757			812
758	Sourabh Katoch, Sumit Singh Chauhan, and Vijay Kumar. 2020. A review on genetic algorithm: past, present, and future . <i>Multimedia Tools and Applications</i> , 80:8091 – 8126.	Stephen Merity, Nitish Shirish Keskar, James Bradbury, and Richard Socher. 2018a. Scalable language modeling: Wikitext-103 on a single gpu in 12 hours. <i>Proceedings of the SYSML</i> , 18.	813
759			814
760			815
761			816
762	Jisu Kim and Juhwan Lee. 2024. Strategic data ordering: Enhancing large language model performance through curriculum learning. <i>ArXiv</i> , abs/2405.07490.		817
763			818
764			819
			820

821	Stephen Merity, Nitish Shirish Keskar, and Richard Socher. 2018b. An analysis of neural language modeling at multiple scales. <i>arXiv preprint arXiv:1803.08240</i> .	Radu Tudor Ionescu, Bogdan Alexe, Marius Leordeanu, Marius Popescu, Dim P Papadopoulos, and Vittorio Ferrari. 2016. How hard can it be? estimating the difficulty of visual search in an image. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition</i> , pages 2157–2166.	874
822			875
823			876
824			877
825	Marwa Nair, Kamel Yamani, Lynda Said Lhadj, and Riyadh Baghdadi. 2024. Curriculum learning for small code language models. <i>arXiv preprint arXiv:2407.10194</i> .	Suresh Venkatasubramanian and Qiushi Wang. 2011. The johnson-lindenstrauss transform: an empirical study. In <i>2011 Proceedings of the Thirteenth Workshop on Algorithm Engineering and Experiments (ALENEX)</i> , pages 164–173. SIAM.	878
826			879
827			880
828			881
829	Yurii Nesterov and Vladimir Spokoiny. 2017. Random gradient-free minimization of convex functions. <i>Foundations of Computational Mathematics</i> , 17(2):527–566.	Xin Wang, Yuwei Zhou, Hong Chen, and Wenwu Zhu. 2024. Curriculum learning for multimedia in the era of large language models. In <i>Proceedings of the 32nd ACM International Conference on Multimedia</i> , pages 11296–11297.	882
830			883
831			884
832			885
833	Ru Peng, Kexin Yang, Yawen Zeng, Junyang Lin, Dayiheng Liu, and Junbo Zhao. 2025. Dataman: Data manager for pre-training large language models. <i>arXiv preprint arXiv:2502.19363</i> .	Yining Wang, Simon Du, Sivaraman Balakrishnan, and Aarti Singh. 2018. Stochastic zeroth-order optimization in high dimensions. In <i>International conference on artificial intelligence and statistics</i> , pages 1356–1365. PMLR.	886
834			887
835			888
836			889
837	Maithra Raghu, Justin Gilmer, Jason Yosinski, and Jascha Sohl-Dickstein. 2017. Svcca: Singular vector canonical correlation analysis for deep learning dynamics and interpretability. <i>Advances in neural information processing systems</i> , 30.	Daphna Weinshall, Gad Cohen, and Dan Amir. 2018. Curriculum learning by transfer learning: Theory and experiments with deep networks. In <i>International conference on machine learning</i> , pages 5238–5246. PMLR.	890
838			891
839			892
840			893
841			894
842	Yi Ren and Danica J. Sutherland. 2025. Learning dynamics of llm finetuning . <i>Preprint</i> , arXiv:2407.10490.	Yinjun Wu, Edgar Dobriban, and Susan Davidson. 2020. Deltagrad: Rapid retraining of machine learning models. In <i>International Conference on Machine Learning</i> , pages 10355–10366. PMLR.	895
843			896
844			897
845	Naomi Saphra. 2021. Training dynamics of neural language models.	Benfeng Xu, Licheng Zhang, Zhendong Mao, Quan Wang, Hongtao Xie, and Yongdong Zhang. 2020. Curriculum learning for natural language understanding. In <i>Proceedings of the 58th annual meeting of the association for computational linguistics</i> , pages 6095–6104.	898
846			899
847	Naomi Saphra and Adam Lopez. 2018. Understanding learning dynamics of language models with svcca. <i>arXiv preprint arXiv:1811.00225</i> .	Fuzhao Xue, Yao Fu, Wangchunshu Zhou, Zangwei Zheng, and Yang You. 2023. To repeat or not to repeat: Insights from scaling llm under token-crisis . <i>Preprint</i> , arXiv:2305.13230.	900
848			901
849			902
850	Beatrice Savoldi, Marco Gaido, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2022. On the dynamics of gender learning in speech translation. In <i>Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)</i> , pages 94–111. Association for Computational Linguistics.	Liang Zhang, Kiran Koshy Thekumparampil, Sewoong Oh, and Niao He. 2023. Dpzero: dimension-independent and differentially private zeroth-order optimization. In <i>International Workshop on Federated Learning in the Age of Foundation Models in Conjunction with NeurIPS 2023</i> .	903
851			904
852			905
853			906
854			907
855			908
856	Ryan Teehan, Miruna Clinciu, Oleg Serikov, Eliza Szczechla, Natasha Seelam, Shachar Mirkin, and Aaron Gokaslan. 2022. Emergent structures and training dynamics in large language models. In <i>Proceedings of BigScience Episode# 5–Workshop on Challenges & Perspectives in Creating Large Language Models</i> , pages 146–159.	Xuan Zhang, Gaurav Kumar, Huda Khayrallah, Kenton Murray, Jeremy Gwinnup, Marianna J Martindale, Paul McNamee, Kevin Duh, and Marine Carpuat. 2018a. An empirical exploration of curriculum learning for neural machine translation. <i>arXiv preprint arXiv:1811.00739</i> .	909
857			910
858			911
859			912
860			913
861			914
862			915
863	Kushal Tirumala, Aram Markosyan, Luke Zettlemoyer, and Armen Aghajanyan. 2022. Memorization without overfitting: Analyzing the training dynamics of large language models. <i>Advances in Neural Information Processing Systems</i> , 35:38274–38290.	Xuemiao Zhang, Feiyu Duan, Liangyu Xu, Yongwei Zhou, Sirui Wang, Rongxiang Weng, Jingang Wang, and Xunliang Cai. 2025a. Frames: Boosting llms with a four-quadrant multi-stage pretraining strategy. <i>arXiv preprint arXiv:2502.05551</i> .	916
864			917
865			918
866			919
867			920
868	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .		921
869			922
870			923
871			924
872			925
873			926
			927
			928
			929
			930

931 Xuemiao Zhang, Liangyu Xu, Feiyu Duan, Yongwei
932 Zhou, Sirui Wang, Rongxiang Weng, Jingang Wang,
933 and Xunliang Cai. 2025b. Preference curriculum:
934 Lms should always be pretrained on their preferred
935 data. *arXiv preprint arXiv:2501.13126*.

936 Yihua Zhang, Pingzhi Li, Junyuan Hong, Jiayang Li,
937 Yimeng Zhang, Wenqing Zheng, Pin-Yu Chen, Ja-
938 son D Lee, Wotao Yin, Mingyi Hong, and 1 oth-
939 ers. 2024. Revisiting zeroth-order optimization for
940 memory-efficient llm fine-tuning: A benchmark.
941 *arXiv preprint arXiv:2402.11592*.

942 Ziwei Zhang, Peng Cui, Haoyang Li, Xiao Wang, and
943 Wenwu Zhu. 2018b. Billion-scale network embed-
944 ding with iterative random projection. In *2018 IEEE*
945 *international conference on data mining (ICDM)*,
946 pages 787–796. IEEE.

947 Deli Zhao, Jiapeng Zhu, Zhenfang Guo, and Bo Zhang.
948 2019. Curriculum learning for deep genera-
949 tive models with clustering. *arXiv preprint*
950 *arXiv:1906.11594*.

951 Xiaosen Zheng and Jing Jiang. 2022. An empirical
952 study of memorization in nlp. In *Proceedings of the*
953 *60th Annual Meeting of the Association for Compu-*
954 *tational Linguistics (Volume 1: Long Papers)*, pages
955 6265–6278.

956 A Technical Details

957 A.1 Precomputation in Update Term Storing Stage

958 Recall that the proposed FUT framework consists of three stages in Figure 1, where the update term storing (Stage 2) plays an important role to bridge the gap between the learning dynamic of reference order and that of the new order. In specific, in Stage 2, we need to compute three kinds of update terms: $\Gamma(\theta_t, B_{l_t})$ and $\nabla_{\theta}\Gamma(\theta_t, B_{l_t})$ for the first-order Taylor expansion (in equation (3)), and the additional $\nabla_{\theta}^2\Gamma(\theta_t, B_{l_t})$ for the second-order Taylor expansion (in equation (7)). In the following, we describe how we compute these three update terms in detail, respectively.

959 • First, for each $\Gamma(\theta_t, B_{l_t})$ term, we can compute it by directly applying the checkpoint θ_t over batch B_{l_t} following the updating rule of Adam optimizer:

$$974 \Gamma(\theta_t, B_{l_t}) = m_t / (\sqrt{v_t} + \epsilon), \quad (9)$$

975 where

$$976 \begin{aligned} m_t &= (\beta_1 m_{t-1} + (1 - \beta_1) \nabla_{\theta} \mathcal{L}_{l_t}) / (1 - \beta_1^t), \\ v_t &= (\beta_2 v_{t-1} + (1 - \beta_2) \nabla_{\theta} \mathcal{L}_{l_t}^2) / (1 - \beta_2^t), \end{aligned} \quad (10)$$

977 where $\nabla_{\theta} \mathcal{L}_{l_t}$ denote $\nabla_{\theta} \mathcal{L}(\theta_t, B_{l_t})$. The accumulative terms m_{t-1} and v_{t-1} terms in m_t and v_t are constructed by the gradient from the last step in the original training process, *i.e.*, $\nabla_{\theta} \mathcal{L}(B_{t-1}; \theta_{t-1})$.

978 • Second, for each first-order update term $\nabla_{\theta}\Gamma(\theta_t, B_{l_t})$, we first expand it as:

$$983 \begin{aligned} \nabla_{\theta}\Gamma(\theta_t, B_{l_t}) &= \frac{\partial \Gamma(\theta_t, B_{l_t})}{\partial \theta} \\ &= \frac{\frac{\partial m_t}{\partial \theta} (\sqrt{v_t} + \epsilon) - \frac{\partial \sqrt{v_t}}{\partial \theta} m_t}{(\sqrt{v_t} + \epsilon)^2} \end{aligned} \quad (11)$$

984 where

$$985 \begin{aligned} \frac{\partial m_t}{\partial \theta} &= \frac{\beta_1 \frac{\partial m_{t-1}}{\partial \theta} + (1 - \beta_1) \nabla_{\theta}^2 \mathcal{L}_{l_t}}{1 - \beta_1^t}, \\ \frac{\partial \sqrt{v_t}}{\partial \theta} &= \frac{\beta_2 \frac{\partial v_{t-1}}{\partial \theta} + 2(1 - \beta_2) \nabla_{\theta} \mathcal{L}_{l_t} \nabla_{\theta}^2 \mathcal{L}_{l_t}}{2(1 - \beta_2^t) \sqrt{v_t}}. \end{aligned} \quad (12)$$

986 To compute the second-order gradient $\nabla_{\theta}^2 \mathcal{L}_{l_t} = \nabla_{\theta}^2 \mathcal{L}(\theta_t, B_{l_t})$, a straightforward approach is to apply the backward operator to $\mathcal{L}(\theta_t, B_{l_t})$ twice. However, this requires computing the Hessian matrix of the parameters, which is prohibitively expensive, especially for LLMs with a large number of parameters.

993 To address this limitation, we approximate the second-order gradient using $(\nabla_{\theta} \mathcal{L}(\theta_t, B_{l_t}) - \nabla_{\theta} \mathcal{L}(\theta_{t-1}, B_{l_t})) / (\theta_t - \theta_{t-1})$, where θ_{t-1} denotes the parameter at step $t-1$ in the original training process. This approximation is justified by the limited variation in parameter updates between adjacent training steps.

994 • At last, for each $\nabla^2\Gamma(\theta_t, B_{l_t})$ term, we can also expand it as:

$$995 \begin{aligned} \nabla^2\Gamma(\theta_t, B_{l_t}) &= \frac{\partial^2 \Gamma(\theta_t, B_{l_t})}{\partial \theta^2} = \frac{1}{(\sqrt{v_t} + \epsilon)^2} \\ &\cdot \left[\left(\frac{\partial^2 m_t}{\partial \theta^2} \right) (\sqrt{v_t} + \epsilon) - \left(\frac{\partial^2 \sqrt{v_t}}{\partial \theta^2} \right) m_t \right. \\ &\left. - 2 \left(\frac{\partial \sqrt{v_t}}{\partial \theta} \right) \left(\frac{\partial m_t}{\partial \theta} \right) + 2 \left(\frac{\partial \sqrt{v_t}}{\partial \theta} \right)^2 \frac{m_t}{\sqrt{v_t} + \epsilon} \right], \end{aligned} \quad (1002)$$

1003 where

$$1004 \begin{aligned} \frac{\partial^2 m_t}{\partial \theta^2} &= \frac{\beta_1 \frac{\partial^2 m_{t-1}}{\partial \theta^2} + (1 - \beta_1) \nabla_{\theta}^3 \mathcal{L}_{l_t}}{1 - \beta_1^t}, \\ \frac{\partial^2 \sqrt{v_t}}{\partial \theta^2} &= \frac{\beta_2 \frac{\partial^2 v_{t-1}}{\partial \theta^2} + 2(1 - \beta_2) [(\nabla_{\theta}^2 \mathcal{L}_{l_t})^2 + \nabla_{\theta} \mathcal{L}_{l_t} \nabla_{\theta}^3 \mathcal{L}_{l_t}]}{2(1 - \beta_2^t) \sqrt{v_t}} \\ &\quad - \frac{\left(\beta_2 \frac{\partial v_{t-1}}{\partial \theta} + 2(1 - \beta_2) \nabla_{\theta} \mathcal{L}_{l_t} \nabla_{\theta}^2 \mathcal{L}_{l_t} \right) \frac{\partial v_t}{\partial \theta}}{4(1 - \beta_2^t) v_t^{3/2}} \end{aligned} \quad (13)$$

1005 Similarly, to compute the third-order gradient $\nabla_{\theta}^3 \mathcal{L}_{l_t} = \nabla_{\theta}^3 \mathcal{L}(\theta_t, B_{l_t})$, we use $(\nabla_{\theta}^2 \mathcal{L}(\theta_t, B_{l_t}) - \nabla_{\theta}^2 \mathcal{L}(\theta_{t-1}, B_{l_t})) / (\theta_t - \theta_{t-1})$ to approximate it.

1006 By computing these update terms for each (θ_t, B_{l_t}) pair, we can access all the update terms we may need in Estimating Stage (Stage 3 in Figure 1). That is, given an arbitrary permuted order, which is different from the reference one, we can recursively execute the first-order Taylor expansion in equation (3) or the second-order Taylor expansion in equation (7) to obtain the new model parameters.

1016 A.2 Random Projection for Storing Update Terms

1017 The update terms $\Gamma(\theta_t, B_{l_t})$, $\nabla_{\theta}\Gamma(\theta_t, B_{l_t})$, and $\nabla_{\theta}^2\Gamma(\theta_t, B_{l_t})$ are essential to our FUT framework. However, in large-scale neural networks such as LLMs, these terms typically have dimensionality comparable to that of the model parameters, making direct precomputation and storage for every

pair (θ_t, B_t) prohibitively expensive in terms of memory.

To mitigate this issue, we adopt a random projection strategy based on the Johnson–Lindenstrauss (JL) theorem (Venkatasubramanian and Wang, 2011), following the well-established compression techniques in Lin et al. (2024). The JL theorem guarantees that high-dimensional vectors can be embedded into a significantly lower-dimensional space with bounded distortion of pairwise distances, which aligns well with our goal of efficiently storing approximate versions of gradient-related terms.

Theorem 1 (Johnson–Lindenstrauss Theorem)

Let $0 < \epsilon < 1$ and let $X = \{x_1, x_2, \dots, x_n\} \subset \mathbb{R}^d$ be a set of n vectors. Then there exists a linear mapping $f : \mathbb{R}^d \rightarrow \mathbb{R}^k$, where $k = \mathcal{O}(\epsilon^{-2} \log n)$, such that for all $x_i, x_j \in X$,

$$(1 - \epsilon)\|x_i - x_j\|_2^2 \leq \|f(x_i) - f(x_j)\|_2^2 \leq (1 + \epsilon)\|x_i - x_j\|_2^2.$$

In our setting, we apply the JL projection to compress each update matrix prior to storage. Formally, for any matrix $M \in \mathbb{R}^{d_1 \times d_2}$ —where M may represent $\Gamma(\theta_t, B_t)$, $\nabla_{\theta}\Gamma(\theta_t, B_t)$, or $\nabla_{\theta}^2\Gamma(\theta_t, B_t)$ —we generate a random projection matrix $A \in \mathbb{R}^{d_2 \times k}$ whose entries are sampled i.i.d. from a Gaussian distribution: $A_{ij} \sim \mathcal{N}(0, 1/k)$. The compressed representation of M is then given by:

$$M' = MA \in \mathbb{R}^{d_1 \times k}.$$

This projection reduces the space complexity from $\mathcal{O}(d_1 d_2)$ to $\mathcal{O}(d_1 k)$ while approximately preserving the geometric structure of the original matrix rows.

To recover these terms for estimating the parameters under a new batch order, an approximate reconstruction can be achieved using the Moore–Penrose pseudoinverse $A^+ \in \mathbb{R}^{k \times d_2}$ as:

$$\widetilde{M} = M'A^+ \approx M.$$

In practice, the target dimension k is selected based on the number of rows d_1 in M , which corresponds to the number of vectors n in Theorem 1. To balance accuracy and memory usage, we empirically choose $k \in \{300, 200, 160, 80, 20, 8\}$ depending on the layer size and update type.

A.3 The algorithm of FUT Framework with First-order Taylor Expansion

The full procedure for deriving the sequence $\{\gamma_t\}_{t=0}^T$ is documented in Algorithm 1. For clarity, we expand here on the three-stage workflow underlying this computation.

(1) Reference model training. We begin by training the reference model M on the training dataset \mathcal{D}_{tr} following a predetermined ordering of mini-batches. This stage produces the trajectory of model parameters $\Theta = \{\theta_t\}_{t=0}^T$, which corresponds to the optimization path under this fixed sample order. This parameter path serves as the backbone for all subsequent computations.

(2) Update-term derivation and storage. Given Θ , we next compute the update-related quantities required for later estimation. Specifically, for each training step $0 \leq t \leq T - 1$, we evaluate $\Gamma(\theta_t, B_t)$ and $\nabla_{\theta}\Gamma(\theta_t, B_t)$ according to equations (2) and (5). These terms describe how the model update behaves at each point along the reference trajectory and are stored for reuse. Since these computations depend only on the reference ordering, they need to be performed only once.

(3) Estimation for a new sample order. Once the above quantities are available, we can estimate the model behavior under any new ordering $\{l_t\}_{t=0}^{T-1}$. For such an order, the sequence $\{\gamma_t\}_{t=0}^T$ is reconstructed recursively using equation (4). This recursive computation relies on the precomputed update terms from Stage (2), allowing each new ordering to be evaluated efficiently without retraining the model.

In practice, the computationally intensive operations—namely the reference training and the derivation of update terms—are executed only once. Afterward, the FUT framework can rapidly estimate the performance associated with arbitrary sample orders. Figure 1 in the main pages provides an overview of the entire FUT pipeline.

A.4 Genetic Algorithm for Training Curriculum Design in FUT Framework

Recall that the objective in equation (9), i.e., $\pi^* := \arg \max_{\pi \in \Pi} \mathcal{R}(\gamma_T^{\pi}, \mathcal{D}_{val})$, is to find the optimal permutation π^* that leads to the best validation performance, where γ_T^{π} represents the final model parameters estimated by FUT framework. However, equation (9) is naturally non-differentiable, hindering its application in finding the optimal curriculum. To address this issue, we design an optimization

Algorithm 1 FUT Framework for Deriving $\{\gamma_t\}_{t=0}^T$ with First-order Taylor Expansion

Require: Initialized model parameter θ_0 , reference training batches $\{B_t\}_{t=0}^{T-1}$, learning rate η , and ϵ .

Ensure: Derived sequence $\{\gamma_t\}_{t=0}^T$

- 1: [Reference Model Training Stage](#):
 - 2: **for** $t = 0$ to $T - 1$ **do**
 - 3: Compute the $(t + 1)$ th checkpoint:
 - 4: $\theta_{t+1} \leftarrow \theta_t - \eta\Gamma(\theta_t, B_t)$ (Eq. 1)
 - 5: **end for**
 - 6: Obtain $\Theta = \{\theta_t\}_{t=0}^T$
 - 7: [Update Term Storing Stage](#):
 - 8: **for** $t = 0$ to $T - 1$ **do**
 - 9: Compute first- and second-order update terms:
 - 10: $\Gamma(\theta_t, B_{l_t}) \leftarrow$ calculate $\nabla_{\theta}\mathcal{L}(\theta_t, B_{l_t})$ with checkpoint θ_t on batch B_{l_t}
 - 11: $\nabla_{\theta}\Gamma(\theta_t, B_{l_t}) \leftarrow$ calculate $\nabla_{\theta}\mathcal{L}(\theta_t, B_{l_t}), \nabla_{\theta}^2\mathcal{L}(\theta_t, B_{l_t})$ with checkpoint θ_t on batch B_{l_t}
 - 12: **end for**
 - 13: [Estimation Stage](#):
 - 14: $\gamma_0 \leftarrow \theta_0$
 - 15: **for** $t = 0$ to $T - 1$ **do**
 - 16: First-order expansion for $\Gamma(\gamma_t, B_{l_t})$:
 - 17: $\Gamma(\gamma_t, B_{l_t}) \leftarrow \Gamma(\theta_t, B_{l_t}) + (\gamma_t - \theta_t)\nabla_{\theta}\Gamma(\theta_t, B_{l_t})$ (Eq. 3)
 - 18: Update γ_{t+1} :
 - 19: $\gamma_{t+1} \leftarrow \gamma_t - \eta\Gamma(\gamma_t, B_{l_t})$ (Eq. 4)
 - 20: **end for**
 - 21: **Return** $\{\gamma_t\}_{t=0}^T$
-

algorithm based on Genetic Algorithm (GA) (Katoch et al., 2020). GA is a well-established meta-heuristic algorithm inspired by Darwinian evolution, which iteratively evolves a population of candidate solutions based on the principle of survival of the fittest. In our context, each candidate represents a specific sample order π , and the fitness of each individual is evaluated by the model’s performance $r^{\pi} = \mathcal{R}(\gamma_T^{\pi}, \mathcal{D}_{val})$. By leveraging crossover, mutation, and selection operators, GA enables us to efficiently explore the exponentially large permutation space without exhaustive enumeration. We describe the detailed design of our GA-based search strategy as follows:

1. **Population Initialization:** Randomly select N sample orders $\text{POP} = \{\pi^i\}_{i=1}^N$ from \mathcal{S}_T as the initial populations, where $\mathcal{S}_T = \{\pi \mid \pi \text{ is a permutation of } \{1, \dots, T\}\}$, with

Algorithm 2 Genetic Algorithm for Finding Optimal Training Curriculum

Require: Validation set \mathcal{D}_{val} , number of batches T , population size N , number of generations K , mutation probability p_m

Ensure: Optimal sample order π^{GA^*}

- 1: Initialize permutation space $\mathcal{S}_T = \{\pi \mid \pi \text{ is a permutation of } \{1, \dots, T\}\}$
 - 2: Randomly sample N permutations as initial population: $\text{POP} = \{\pi^i\}_{i=1}^N \subset \mathcal{S}_T$
 - 3: **for** $k = 1$ to K **do**
 - 4: **for all** $\pi^i \in \text{POP}$ **do**
 - 5: Compute $\gamma_T^{\pi^i}$ using FUT with sample order π^i
 - 6: Evaluate fitness $r^i = \mathcal{R}(\gamma_T^{\pi^i}, \mathcal{D}_{val})$
 - 7: **end for**
 - 8: Retain top 50% individuals with highest fitness to form $\text{POP}_{\text{survive}}$
 - 9: **while** Size of new children $< N/2$ **do**
 - 10: Randomly select two parents π^a, π^b from $\text{POP}_{\text{survive}}$
 - 11: Randomly choose crossover points l, r such that $1 \leq l < r \leq T$
 - 12: Generate child $\pi^c = \text{PMX}(\pi^a, \pi^b, l, r)$
 - 13: **if** $\text{random}() < p_m$ **then**
 - 14: Randomly select positions i, j and swap π_i^c and π_j^c
 - 15: **end if**
 - 16: Add π^c to new children
 - 17: **end while**
 - 18: Replace discarded individuals in POP with new children
 - 19: **end for**
 - 20: **return** $\pi^{\text{GA}^*} = \arg \max_{\pi \in \text{POP}} \mathcal{R}(\gamma_T^{\pi}, \mathcal{D}_{val})$
-

$$|\mathcal{S}_T| = T!.$$

2. **Fitness Selection:** For each $\pi^i \in \text{POP}$, evaluate the model performance $\mathcal{R}(\gamma_T^{\pi^i}, \mathcal{D}_{val})$ as its fitness, where $\gamma_T^{\pi^i}$ is estimated via the FUT method. Retain the top 50% individuals with the highest fitness scores for reproduction, and discard the rest.
3. **Crossover:** Generate new children by applying the partially matched crossover (PMX) (Kora and Yadlapalli, 2017) to randomly selected parent pairs π^a and π^b from the surviving population. Specifically, randomly choose two crossover points l and r such that $1 \leq l < r \leq T$, then exchange

the subsequences $\pi_{l:r}^a$ and $\pi_{l:r}^b$ between the parents. At last, resolve conflicts using the mapping induced by the swapped segments to produce a valid permutation child $\pi^c = \text{PMX}(\pi^a, \pi^b, l, r) \in \mathcal{S}_T$.

- Mutation:** With a predefined mutation probability p_m , randomly select two indices i and j in π^c and swap their values: $\pi^c \leftarrow \pi_{i \leftrightarrow j}^c$. This operation introduces diversity and prevents premature convergence.
- Replacement:** Insert the newly generated children into the population, replacing the discarded individuals. The updated population then forms the basis for the next generation.

By iteratively performing 2-5 steps over a fixed number of generations K , or until a convergence criterion is met (e.g., no improvement in validation performance over several generations), the algorithm ultimately returns the best sample order π^{GA^*} with the highest validation performance. Therefore, this GA-based optimization reduces the inference time complexity of FUT from $\mathcal{O}(T!)$ to $\mathcal{O}(K \cdot N)$, significantly accelerating the search for the optimal sample order.

B Experimental Details

B.1 General Capability

In this section, we introduce more details for the experiments to test the general capability of our FUT framework in Section 5.1.

B.1.1 Base Model

We conduct all of our experiments on a language model that follows the LLaMA architecture (Touvron et al., 2023), but with a reduced number of parameters—specifically, a hidden size of 2048 and 10 stacked transformer layers, resulting in approximately 636 million parameters.

We choose this relatively small model to enable repeated training under varying experimental conditions, which is essential for rigorously evaluating the effectiveness of our proposed FUT framework in both training curriculum design and the analysis of memorization and generalization behaviors. In contrast, training large-scale models typically takes tens or even hundreds of days, making such extensive experimentation prohibitively time-consuming and computationally expensive.

B.1.2 Dataset

WikiText-103 (Merity et al., 2018a) is a widely used benchmark dataset for evaluating language models, particularly in long-range dependency modeling. It consists of over 100 million tokens extracted from high-quality Wikipedia articles, specifically curated to preserve coherent paragraph and document-level structures. Unlike other common datasets that contain shuffled or sentence-level data, WikiText-103 maintains the original article formatting and ordering, enabling models to better learn contextual and discourse-level information. The vocabulary is relatively large and diverse, making it a challenging and realistic corpus for testing the generalization and memorization capabilities of large-scale language models. In our experiments, we partition the dataset into 80% for training, 10% for validation, and the remaining 10% for testing.

To preprocess the WikiText-103 dataset, we first remove short texts with fewer than five characters to eliminate noise. Then, we apply MinHash-based deduplication (Broder, 1997) to efficiently identify and discard near-duplicate samples. Specifically, each text is tokenized into a set of words, and a MinHash signature is computed using 128 permutations. Texts with identical MinHash digests are considered duplicates, and only one representative is retained. This process effectively reduces redundancy while preserving semantically diverse content.

B.1.3 Training and Evaluation Protocols

Training Protocol. For the preprocessed WikiText dataset, we split the data into 80%, 10%, and 10% for training, validation, and testing, respectively. The learning rate is selected from the range $[0.0001, 0.005]$ based on validation performance, and we choose the number of batches from $\{8, 16, 32, 64, 128, 256\}$. Since it is not feasible to process very large batch sizes directly due to memory constraints, we apply gradient accumulation over multiple smaller mini-batches to effectively simulate the desired larger batch size. For the Adam optimizer, we fix the hyperparameters β_1 and β_2 to 0.9 and 0.95, respectively. Within our FUT framework, to stabilize parameter estimation and mitigate the influence of outliers, we apply parameter clipping. Specifically, the parameters are constrained within a tunable range, with the clipping threshold selected from the interval $[-1.1, -0.3] \cup [0.3, 1.1]$ to ensure numerical stability and prevent extreme values from dominating the

update dynamics. The experiments were conducted on a computing platform equipped with NVIDIA A800-SXM GPUs, with a total of 4 GPUs each providing 80GB of memory.

Evaluation Protocol. We adopt Perplexity (PPL) (Hu et al., 2024) as the evaluation metric to assess language modeling performance. Given a token sequence $x = (x_1, x_2, \dots, x_N)$, the perplexity is defined as:

$$\begin{aligned} \text{PPL}(x) &= P(x_1, \dots, x_N)^{-\frac{1}{N}} \\ &= \exp\left(-\frac{1}{N} \sum_{t=1}^N \log P(x_t | x_{<t})\right). \end{aligned} \quad (14)$$

This is equivalent to the exponential of the average cross-entropy loss. Thus, for a given validation set \mathcal{D}_{val} and final model parameters θ_T , we compute:

$$\text{PPL}(\mathcal{D}_{val}) = \exp(\mathcal{L}(\mathcal{D}_{val}; \theta_T)), \quad (15)$$

where $\mathcal{L}(\mathcal{D}_{val}; \theta_T)$ denotes the average cross-entropy loss over the validation set.

B.2 Training Curriculum Design for LLMs

B.2.1 Baselines

Although curriculum learning largely depends on human heuristics or empirical findings, there are still many works that make efforts to design a rational curriculum in the field of LLMs, primarily based on either the characteristics of the dataset (Campos, 2021), or the quantitative criteria (Zhang et al., 2025a,b) that are perceptible to the model. In this section, we introduce all the baselines used in training curriculum design in detail. For better understanding, we define ρ_{B_i} as the difficulty score for batch B_i .

- **Random Order (RO).** RO is a naive baseline, which randomly assigns the difficulty score ρ_{B_i} to each batch B_i in the range of $[0, 1]$.

- **Sample Length (SL) (Campos, 2021).** SL is a purely statistical method based on the intuition that longer sentences are inherently more difficult to model. This is because they require more effective tracking of dependencies, making the learning process more challenging. Therefore, the difficulty score of each batch B_i is defined as the total number of tokens in the batch, computed as $\rho_{B_i} = \sum_{x \in B_i} |x|$, where $|x|$ denotes the length of sample x .

- **Perplexity (PPL) (Zhang et al., 2025a).** PPL metric closely aligns with the self-supervised learning objective of LLMs and effectively measures model-data fit, making it appropriate for data organization. Recent studies (Zhang et al., 2025a) empirically show that training on high-PPL data followed by low-PPL data can significantly reduce loss and boost performance. Following this finding, we introduce a reference model M_{ref} with parameter θ_R to compute PPL for each batch as the difficulty score, *i.e.*, $\rho_{B_i} = -\mathcal{R}(\theta_R, B_i)$.

- **Perplexity Difference (PD) (Zhang et al., 2025b).** Building on the idea in (Zhang et al., 2025b), PD between strong and weak models can serve as an indicator of how difficult a batch is for the model. Specifically, a low PD implies that both models perform similarly in terms of learning efficiency, while a high PD suggests that the batch presents greater difficulty for the weaker model. Consider two reference models, M_{str} and M_{weak} , with parameters θ_S and θ_W , respectively, both trained on the same dataset. In practice, we train two models: M_{str} with 636 million parameters and M_{weak} with 167 million parameters, using their perplexity differences to guide batch rescheduling. For each batch B_i , we define PD as the difficulty score, given by $\rho_{B_i} = (\mathcal{R}(\theta_W, B_i) - \mathcal{R}(\theta_S, B_i)) / \mathcal{R}(\theta_W, B_i)$.

B.2.2 Genetic Algorithm Configuration

To effectively search the optimal sample order within the exponentially large permutation space, we employ a Genetic Algorithm (GA) tailored to our FUT framework. The key design choices focus on maintaining a balance between exploration and exploitation: a moderately sized population ensures sufficient diversity, while elitist selection preserves high-quality solutions across generations. The complete set of hyperparameters and their configurations are summarized in Table 3.

C Additional Experimental Results

C.1 Scalability of FUT Framework

Experimental Setup. In this section, we conduct additional experiments to evaluate whether our proposed FUT framework remains effective in estimating model performance as the base model size increases. Specifically, we scale the original 0.6B model to $\{1.2\text{B}, 2.0\text{B}, 4.0\text{B}, 6.0\text{B}\}$. In these experiments, the number of training batches is set to $T = 8$ and $T = 16$. We adopt perplexity as the

Table 3: Genetic Algorithm hyperparameters used in our framework

Hyperparameter	Notation	Description	Scope
Population size	N	Number of candidates per generation	[16, 12, 8, 4, 2]
Max generations	K	Total evolution rounds	[16, 12, 8, 4, 2, 1]
Number of batches	T	Total number of Batches	[256, 128, 64, 32, 16, 8]
Crossover points	l, r	Random crossover segment indices	$1 \leq l < r \leq T$
Mutation probability	p_m	Swap probability per child	0.1
Selection rate	–	Top individuals retained	50%

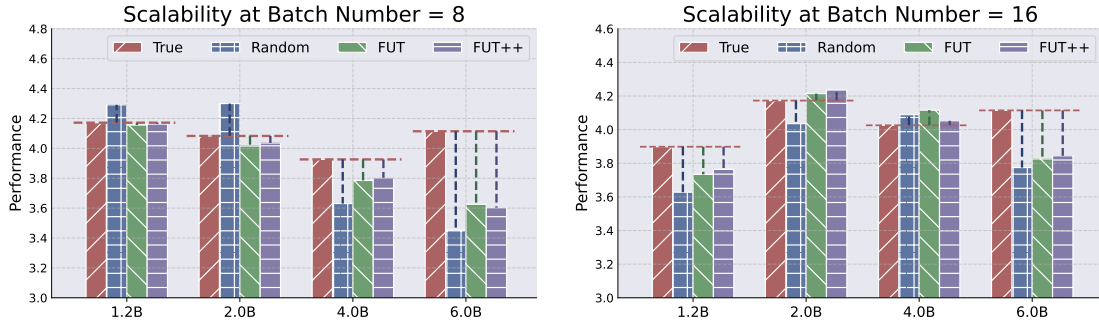


Figure 6: **Scalable estimation performance across model sizes.** We evaluate the estimation accuracy of FUT and FUT++ across model scales $\{1.2\text{B}, 2.0\text{B}, 4.0\text{B}, 6.0\text{B}\}$ under training batch numbers $T = 8$ (left) and $T = 16$ (right). FUT and FUT++ consistently outperform the Random baseline, with FUT++ showing improved accuracy for larger models.

1339 evaluation metric and measure the performance gap
 1340 between the true values and the estimates produced
 1341 by our FUT framework.

1342 **Results.** The results are illustrated in Figure 6.
 1343 Across both batch settings ($T = 8$ and $T = 16$),
 1344 our proposed FUT and FUT++ methods consistently
 1345 outperform the Random baseline in estimating
 1346 model performance, achieving smaller performance
 1347 gaps to the ground truth. This trend holds
 1348 true as we scale the base model size from 1.2B to
 1349 6.0B, validating the scalability of our framework.
 1350 Importantly, we observe that FUT++—which in-
 1351 corporates second-order information—yields even
 1352 more accurate estimates compared to the original
 1353 FUT, particularly for larger models. This suggests
 1354 that higher-order approximations are more effective
 1355 at capturing complex parameter updates in large-
 1356 scale language models. The Random baseline, by
 1357 contrast, lacks theoretical grounding and exhibits
 1358 less consistent performance as model size grows.

1359 C.2 Impact of Random Projection Technique

1360 **Experimental Setup.** Random projection is used
 1361 in our FUT framework to squeeze the storage com-
 1362 plexity, which makes the performance estimation of
 1363 large model with massive parameters become possi-
 1364 ble. In this section, we aim to explore how random

1365 projection can affect the final estimation accuracy.
 1366 Since the entire storage of pairs of update terms
 1367 is unaffordable for a single machine, we conduct
 1368 this ablation study in a relatively small model with
 1369 200M parameters. All other experimental setup is
 1370 consistent with the main experiment.

1371 **Results.** The results are presented in Tale 4,
 1372 where we can see that the use of random projection
 1373 may have negative impact on estimation accuracy,
 1374 because the variant without using random projec-
 1375 tion have a more accurate estimation to the true
 1376 perplexity performance. Nevertheless, we can ob-
 1377 serve that the use of random projection can largely
 1378 save the storage memory, especially when the num-
 1379 ber of batches increases. Therefore in practice,
 1380 random projection can still have great function in
 1381 implementing our framework.

1382 C.3 Batch-wise Analysis of Performance 1383 Estimation Accuracy

1384 **Experimental Setup.** In this section, we conduct
 1385 a fine-grained evaluation of our FUT framework
 1386 by comparing estimated and true model perfor-
 1387 mance at intermediate stages of training. Specifi-
 1388 cally, we consider batch numbers $T \in \{8, 16, 32\}$
 1389 and evaluate performance after each training batch.
 1390 For each time step $1 \leq t \leq T$, we replace the

Table 4: Impact of random projection. We experiment on all different number of batches from 8 to 256 in model with 200M parameters. The results show that the use of random projection can inevitably hurt the estimation accuracy, but it can largely save the storage memory.

Batch Num	Accuracy			Memory	
	True	With RP	W/O RP	With RP	W/O RP
8	1.4885	1.5158	1.5016	~ 0.3G	~ 12G
16	1.4569	1.4793	1.4789	~ 1.6G	~ 52G
32	1.4512	1.4839	1.4783	~ 6.3G	~ 200G
64	1.4485	1.4468	1.4647	~ 26G	~ 820G
128	1.4359	1.4269	1.4394	~ 102G	~ 3.2T
256	1.4268	1.4258	1.4358	~ 410G	~ 13.1T

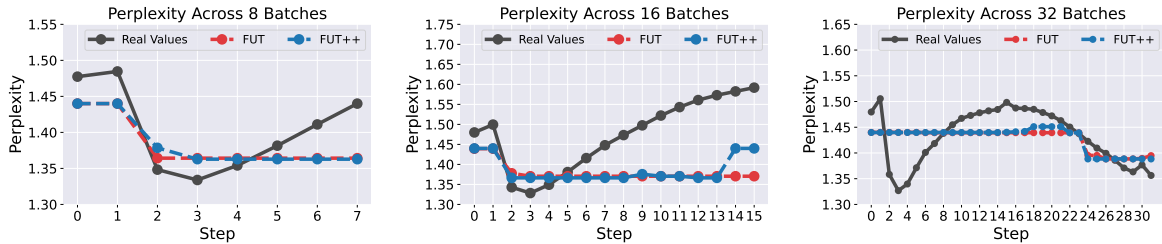


Figure 7: **Perplexity estimation at intermediate training steps.** We visualize the validation perplexity estimated by FUT and FUT++ compared to the real validation perplexity after each batch, for training schedules with $T \in \{8, 16, 32\}$ total batches. FUT and FUT++ both closely follow the true performance trends, with FUT++ consistently providing more accurate estimates—especially when T is larger. These results demonstrate the effectiveness of our methods in tracking training progress in a fine-grained manner.

1391 final-step performance comparison $\mathcal{R}(\gamma_T^\pi, \mathcal{D}_{val})$ and $\mathcal{R}(\theta_T^\pi, \mathcal{D}_{val})$ with the intermediate-step comparison $\mathcal{R}(\gamma_t^\pi, \mathcal{D}_{val})$ and $\mathcal{R}(\theta_t^\pi, \mathcal{D}_{val})$. We use perplexity on the validation set \mathcal{D}_{val} as the evaluation metric to assess how well the FUT-estimated parameters align with those obtained from actual training at each step.

1398 **Results.** As shown in Figure 7, both FUT and FUT++ generate accurate perplexity estimates across different training stages. While FUT performs well in general, FUT++ shows higher fidelity—especially as the number of batches increases. This is most evident in the $T = 32$ case, where FUT++ remains close to the true perplexity throughout, whereas FUT slightly deviates in later stages. These findings affirm the utility of incorporating higher-order dynamics in FUT++, and highlight the robustness of our framework in real-time model monitoring, dynamic training adaptation, and early stopping decisions. In addition, we observe that in certain training stages, particularly under small batch sizes or early steps, the estimated perplexity remains unchanged over multiple steps, forming plateau-like segments. This phenomenon arises from the Taylor-based approxi-

1416 mation mechanism in our framework. Specifically, when the update gradients are small (*e.g.*, due to flat regions in the loss landscape), the computed updates become negligible. Consequently, FUT and FUT++ produce nearly identical estimates across consecutive steps.

1422 D Further Discussions of Related Work

1423 D.1 Curriculum Learning for LLMs

1424 Curriculum learning is a training paradigm that organizes training data in an easy-to-hard manner to facilitate more effective learning (Bengio et al., 2009; Graves et al., 2017; Hachohen and Weinshall, 2019; Xu et al., 2020). In deep learning tasks, sample difficulty is typically defined using either surface-level heuristics or model-based metrics (Matiisen et al., 2019; Hachohen and Weinshall, 2019; Gui et al., 2017; Ghebrechristos and Alagband, 2020; Weinshall et al., 2018). For instance, in sequence modeling, easier examples are often shorter or contain more frequent tokens (Zhang et al., 2018a). In the generative modeling domain, difficulty can be measured by how well a sample aligns with human cognitive expectations or

its deviation from the data distribution center (Tudor Ionescu et al., 2016; Zhao et al., 2019). In the context of LLMs, several empirical studies have explored strategies to score training samples (Nair et al., 2024; Liang et al., 2024; Wang et al., 2024; Matiisen et al., 2019; Campos, 2021; Zhang et al., 2025a,b). Specifically, Campos (2021) reorganizes samples based on their sequence length to progressively improve the model’s ability to capture long-range dependencies. Furthermore, some researchers (Zhang et al., 2025a,b) propose curriculum schemes guided by model-based metrics such as perplexity and perplexity difference, motivated by their empirical observations. **In contrast to conventional curriculum learning approaches that depend on human-designed heuristics for determining sample order, our proposed FUT framework offers an efficient and reliable means of estimating final performance across arbitrary curricula.** This allows practitioners to make well-informed decisions among diverse curriculum strategies without incurring the cost of repeated retraining.

D.2 Zeroth-Order Optimization

Zeroth-order (ZO) optimization refers to a class of derivative-free methods that estimate gradients using only function evaluations, making them suitable for black-box or simulation-based scenarios where gradients are inaccessible or costly (Flaxman et al., 2004; Ghadimi and Lan, 2013; Nesterov and Spokoiny, 2017; Duchi et al., 2015; Wang et al., 2018; Chen et al., 2017; Liu et al., 2020). Classical approaches include finite-difference methods (Flaxman et al., 2004), random gradient estimators (Nesterov and Spokoiny, 2017; Duchi et al., 2015), and ZO-SGD (Ghadimi and Lan, 2013). Recently, ZO has been applied to LLM fine-tuning to reduce the memory burden of back-propagation. Notably, MeZO (Malladi et al., 2023) introduced a forward-only ZO-SGD variant, while Zhang et al. (Zhang et al., 2024) benchmarked and extended ZO techniques—such as ZO-Adam (Chen et al., 2019b) and block-wise estimation—for scalable LLM fine-tuning. **However, applying ZO to pre-training remains impractical due to the extreme dimensionality of LLMs, high variance of estimators, and computational overhead from repeated forward passes** (Zhang et al., 2023; Golovin et al., 2019; Wang et al., 2018). Moreover, most of ZO methods rely on dynamic random perturbations, limiting result reproducibility and reuse. **In contrast, our**

FUT framework is a performance estimation tool—not an optimizer—that precomputes all necessary update terms using Taylor expansions. This enables efficient, deterministic evaluation of arbitrary curricula without retraining, making FUT quite suitable for analyzing training dynamics and guiding curriculum design.

1490
1491
1492
1493
1494
1495
1496