

# SCALING LAWS FOR PREDICTING DOWNSTREAM PERFORMANCE IN LLMs

Anonymous authors

Paper under double-blind review

## ABSTRACT

Precise estimation of downstream performance in large language models (LLMs) prior to training is essential for guiding their development process. Scaling laws analysis utilizes the statistics of a series of significantly smaller sampling language models (LMs) to predict the performance of the target LLM. For downstream performance prediction, the critical challenge lies in the emergent abilities in LLMs that occur beyond task-specific computational thresholds. In this work, we focus on the pre-training loss as a more computation-efficient metric for performance estimation. Our two-stage approach consists of first estimating a function that maps computational resources (*e.g.*, FLOPs) to the pre-training Loss using a series of sampling models, followed by mapping the pre-training loss to downstream task Performance after the critical “emergent phase”. In preliminary experiments, this **FLP** solution accurately predicts the performance of LLMs with 7B and 13B parameters using a series of sampling LMs up to 3B, achieving error margins of 5% and 10%, respectively, and significantly outperforming the FLOPs-to-Performance approach. This motivates **FLP-M**, a fundamental approach for performance prediction that addresses the practical need to integrate datasets from multiple sources during pre-training, specifically blending general corpora with code data to accurately represent the common necessity. **FLP-M** extends the power law analytical function to predict domain-specific pre-training loss based on FLOPs across data sources, and employs a two-layer neural network to model the non-linear relationship between multiple domain-specific loss and downstream performance. By utilizing a 3B LLM trained on a specific ratio and a series of smaller sampling LMs, **FLP-M** can effectively forecast the performance of 3B and 7B LLMs across various data mixtures for most benchmarks within 10% error margins.

## 1 INTRODUCTION

Large language models (LLMs) form the basis for numerous real-world applications (Brown et al., 2020; Jiang et al., 2023; Touvron et al., 2023) and scaling laws analysis serves as the foundation for LLMs development (Kaplan et al., 2020; Bahri et al., 2024). The key idea of scaling laws involves training a sequence of language models (LMs) to gather data (*e.g.*, expended compute and corresponding model performance). This data is then used to build a predictive model that estimates the performance of a substantially larger target LLM (Su et al., 2024; Hoffmann et al., 2022).

Previous efforts focus on predicting the target LLM’s pre-training loss and establish a power-law relation between the computational resource expended (*e.g.*, floating-point operations per second (FLOPs)) and the final loss achieved (Kaplan et al., 2020; Muennighoff et al., 2024; Henighan et al., 2020). Further, we aim to predict the downstream performance in LLMs to more accurately reflect the primary concerns regarding their

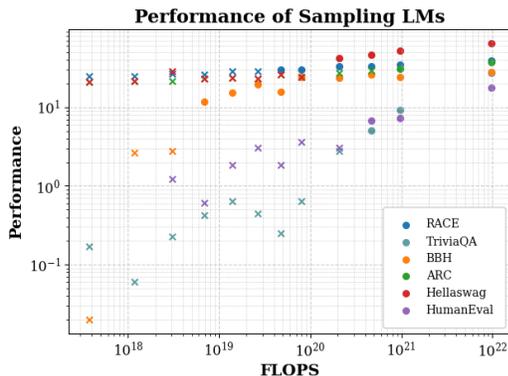


Figure 1: The performance of sampling LMs with increasing compute.  $\times$  represents non-emerged data points, and  $\bullet$  indicates emerged data points that surpass a randomness threshold of 5.

054 capabilities. The critical challenge is the emergent abilities in LLMs, which states that LLMs  
 055 only exceed random performance when the FLOPs expended during training surpass task-specific  
 056 thresholds (Wei et al., 2022). Supposing a task threshold of  $F_c$ , typical methods require training  
 057  $N$  LMs, expending total FLOPs  $F_t = \sum_{i=1}^N \text{FLOPs}_i > N \times F_c$ , to obtain  $N$  effective data points,  
 058 thereby necessitating significant computational resources. Fig. 1 demonstrates that the sampling  
 059 LMs require more than  $5 \times 10^{20}$  FLOPs to perform better than random on most benchmarks, with  
 060 only three data points available to fit the predictive curve across these benchmarks. Hu et al. (2023)  
 061 address this challenge by significantly increasing the sampling times to compute the `PassUntil`  
 062 of a task, basically increasing the “metric resolution” to enable the abilities to emerge earlier (*i.e.*,  
 063 reducing  $F_c$ ). However, this approach faces challenges in translating the `PassUntil` back to the  
 064 original task metric of concerns and requires huge amounts of FLOPs spent on sampling.

065 In this work, we target the actual task performance prediction based on two intuitions: (1) Predicting  
 066 the target pre-training loss is easier and achievable since there is no “emergent phase” in the pre-  
 067 training loss, as extensively verified in Kaplan et al. (2020); Hoffmann et al. (2022); (2) There is an  
 068 observed correlation between the pre-training loss and the downstream task performance after the  
 069 “emergent point” (*i.e.*, the pre-training loss goes below a critical threshold) (Du et al., 2024; Huang  
 070 et al., 2024). Rethinking previous practice, training LM  $i$  to convergence requires expending  $\text{FLOPs}_i$   
 071 for obtaining a single data point. In contrast, our approach stems from a crucial insight: **Collecting**  
 072 **(pre-training loss, performance) data points at intermediate checkpoints prevents the need for**  
 073 **fully training LMs to convergence, thereby enhancing sample efficiency.** Essentially, we can  
 074 collect a huge amount of useful data between the initial point of above-random performance in LMs  
 075 and their convergence point for performance prediction.

076 Thus, our approach consists of two sequential stages: (1) **FLOPs**  $\rightarrow$  **Loss**: Predict the target pre-  
 077 training loss based on the expended FLOPs. Following previous work, we train a series of sampling  
 078 LMs within the same model family to develop a power-law predictive model. For this stage, the  
 079 expended FLOPs are not required to reach above the emergent threshold. (2) **Loss**  $\rightarrow$  **Performance**:  
 080 Predict the downstream performance based on the pre-training loss. We collect data points from inter-  
 081 mediate checkpoints of various sampling LMs that exhibit above-random performance, and develop  
 082 a regression model for prediction. In preliminary experiments with sampling LMs up to 3B, this **FLP**  
 083 solution predicts the performance of 7B and 13B LLMs across various benchmarks with error margins  
 084 of 5% and 10% respectively, significantly outperforming direct FLOPs-to-Performance predictions.

085 Motivated by these findings, we present **FLP-M**, a fundamental solution for performance prediction  
 086 that addresses the growing demand for integrating diverse datasets during LLMs pre-training, focusing  
 087 on integrating the general corpus with code data in this work. **FLP-M** targets fine-grained domain-  
 088 specific pre-training loss to capture the performance changes. Specifically, we extend the power law  
 089 analytical function to predict the domain-specific loss based on FLOPs across multiple data sources.  
 090 Then we employ a two-layer neural network to model the non-linear relationship between multiple  
 091 domain-specific loss and the downstream performance. Through evaluation, we demonstrate that  
 092 **FLP-M** effectively predicts the performance of 3B and 7B LLMs trained on various data mixtures  
 093 (within 10% error margins for most benchmarks). This is achieved by utilizing a 3B LLM trained on  
 094 a specific data mixing ratio along with a series of smaller sampling LMs.

## 095 2 RELATED WORK

### 096 2.1 SCALING LAWS

097 Estimating the performance of the target LLM prior to training is essential due to the significant  
 098 resources required for pre-training (Minaee et al., 2024; Wan et al., 2023). The scaling laws of  
 099 LLMs guide the systematic exploration in scaling up computational resources, data, and model  
 100 sizes (Kaplan et al., 2020; Hestness et al., 2017). Previous efforts in this filed demonstrate that  
 101 LLMs’ final pre-training loss on a held-out validation set decreases with an increase in expended  
 102 FLOPs during pre-training (Kaplan et al., 2020; Hoffmann et al., 2022; Yao et al., 2023). The  
 103 following work subsequently establishes the scaling laws for computer vision models (Zhai et al.,  
 104 2022), vision-language models (Henighan et al., 2020; Alabdulmohsin et al., 2022; Li et al., 2024a),  
 105 graph self-supervised learning (Ma et al., 2024), reward modeling (Gao et al., 2023a; Rafailov  
 106 et al., 2024), data filtering (Goyal et al., 2024), knowledge capabilities of LLMs (Allen-Zhu  
 107

108 & Li, 2024), data-constrained LMs (Muennighoff et al., 2024), data poisoning (Bowen et al.,  
 109 2024), LLMs vocabulary size (Tao et al., 2024), retrieval-augmented LLMs (Shao et al., 2024),  
 110 continued pre-training of LLMs (Que et al., 2024), LLMs training steps (Tissue et al., 2024),  
 111 fine-tuning LLMs (Tay et al., 2021; Lin et al., 2024; Hernandez et al., 2021), learning from repeated  
 112 data (Hernandez et al., 2022), the sparse auto-encoders (Gao et al., 2024), hyper-parameters in LLMs  
 113 pre-training (Yang et al., 2022; Lingle, 2024), and the mixture-of-expert LLMs (Clark et al., 2022;  
 114 Frantar et al., 2023; Yun et al., 2024; Krajewski et al., 2024).

115 Despite the efforts, directly estimating the downstream performance of LLMs more accurately  
 116 reflects the models’ capabilities pertinent to our concerns, yet it confronts challenges associated with  
 117 emergent abilities in LLMs (Wei et al., 2022). In general, the compute required for pre-training must  
 118 surpass a task-specific threshold to enable pre-trained LMs to perform better than random chance.  
 119 Previous work addresses this challenge by using the answer loss as an alternative metric (Schaeffer  
 120 et al., 2024) or increasing the metric resolution, such as measuring the average number of attempts to  
 121 solve the task (Hu et al., 2023). However, they encounter difficulties in aligning the proposed metric  
 122 with the original task metric, which is of paramount interest to us. Our research directly predicts  
 123 the task performance metrics of the target LLMs by utilizing readily available intermediate LMs.  
 124 This approach operates independently from and complements existing approaches.

## 125 2.2 DATA MIXTURE

127 Creating the pre-training dataset necessities collecting data from different sources (Liu et al., 2023;  
 128 Shen et al., 2023; Bi et al., 2024; Wei et al., 2023), making the data mixture a critical factor in the  
 129 study of scaling laws. Ye et al. (2024) propose the data mixing laws to predict the pre-training loss  
 130 of the target LLM given the mixing ratios. Liu et al. (2024) build the regression model to predict  
 131 the optimal data mixture regarding the pre-training loss optimization, and Kang et al. (2024) further  
 132 show that the optimal data composition depends on the scale of compute. In this work, we focus on  
 133 integrating the data mixture factor to better predict the downstream performance.

## 135 3 FLP: DOWNSTREAM PERFORMANCE PREDICTION

137 We introduce a *two-stage* approach to predicting downstream performance in LLMs based on two  
 138 established findings: (1) Predicting the target pre-training loss and establishing the power-law relation  
 139 is feasible as it does not involve an emergent phase (Kaplan et al., 2020; Hoffmann et al., 2022). (2)  
 140 When pre-training loss goes below a task-specific threshold, there is an observed correlation between  
 141 pre-training loss and downstream task performance (Du et al., 2024; Huang et al., 2024). In this sec-  
 142 tion, we present FLP as a proof-of-concept for this framework with a straightforward implementation.

### 143 3.1 FLOPS $\rightarrow$ LOSS

145 We follow the previous practice to use the analytical power law function to characterize the relation  
 146 between expended FLOPs  $C$  and the pre-training loss  $L$ :

$$147 L(C) = \left( \frac{C}{C_N} \right)^{\alpha_N}, \quad (1)$$

149 where  $C_N$  and  $\alpha_N$  are constant terms to be estimated. In FLP, we train a series of  $N$  LLMs within the  
 150 same model family in the same pre-training distribution, progressively increasing model size and  
 151 training tokens to achieve even sampling. Then we measure their pre-training loss in our curated  
 152 validation dataset to obtain  $N$  pairs of  $(C_i, L_i)$  to estimate the constants in Eq. 1.

### 154 3.2 LOSS $\rightarrow$ PERFORMANCE

156 Based on our empirical observation of the scatter plots showing (pre-training loss, performance)  
 157 data points (see §A), we select the analytical linear function to characterize the relation between the  
 158 pre-training loss  $L$  on general validation data and the task performance  $P$ :

$$159 P(L) = w_0 + w_1 * L, \quad (2)$$

160 where  $w_0$  and  $w_1$  are constant terms to be estimated. In FLP, we fetch the intermediate checkpoints  
 161 of each sampling LM, and measure its task performance and pre-training loss. If the performance  $P_i$

Table 1: The configurations of the sampling and target LMs with various sizes. HD denotes the hidden dimension, BS denotes the batch size, and LR denotes the learning rate.

Model Size	#Layer	HD	#Head	FFN	#Tokens	Non-embedding FLOPs	BS	LR
43M	3	384	3	1032	8,021,606,400	3.70504E+17	448	0.0052
64M	4	512	4	1376	11,714,691,072	1.18417E+18	544	0.0042
89M	5	640	5	1720	16,184,770,560	3.03607E+18	576	0.0038
0.12B	6	768	6	2064	21,799,895,040	6.81931E+18	640	0.0040
0.15B	7	896	7	2408	28,846,325,760	1.39581E+19	672	0.0042
0.2B	8	1024	8	2752	37,213,962,240	2.63435E+19	736	0.0036
0.25B	9	1152	9	3096	47,563,407,360	4.71817E+19	768	0.0034
0.32B	10	1280	10	3440	59,674,460,160	8.01571E+19	800	0.0028
0.5B	12	1536	12	4128	90,502,594,560	2.05963E+20	960	0.0023
0.72B	14	1792	14	4816	132,026,204,160	4.70331E+20	1024	0.0019
1B	16	2048	16	5504	185,535,037,440	9.75926E+20	1152	0.0016
3B	24	3072	24	8256	556,793,856,000	9.63212E+21	1536	0.0004
7B	32	4096	32	11008	1,258,291,200,000	5.09208E+22	2048	0.0003
13B	40	5120	40	13824	1,258,291,200,000	9.89592E+22	2048	0.0003

of  $LM_i$  exceeds the random performance, we can obtain one effective data point  $(L_i, P_i)$  to estimate the constants in Eq. 2, where  $L_i$  is the pre-training loss of  $LM_i$ .

## 4 VALIDATION OF FLP FRAMEWORK

### 4.1 SAMPLING AND TARGET LMS

We train a series of 12 sampling LMs up to 3B parameters to predict the performance of target LLMs with 7B and 13B parameters. The configurations of LMs are shown in Tab. 1. We first determine the number of training tokens required for the 7B LLM (approximately 180 times the model size), considering practical needs and inference-time costs. In real-world applications, prioritizing inference efficiency often involves training smaller LMs with a higher token-to-parameter ratio beyond the optimal factor of 20x (Hoffmann et al., 2022). Our preliminary experiments indicate that scaling laws remain applicable even in this over-training regime (within 2.8% error margins). We then proportionally scale down this number to determine the required training tokens for the sampling LMs.

### 4.2 DATA: PRE-TRAINING, VALIDATION, EVALUATION

**Pre-Training** We use the RedPajama v1 (Computer, 2023), which consists of 1.2T tokens in total, and the data is sourced from Arxiv, C4, Common Crawl, GitHub, Stack Exchange, and Wikipedia.

**Validation** We curate a validation dataset to measure the final pre-training loss, which includes 5 distinct domains: math, code, scientific paper, Wikipedia, and general language corpus. Specifically, we utilize subsets from GitHub, ArXiv, Wikipedia, and the English portion of C4, all from the RedPajama validation sets, along with Proof Pile (Touvron et al., 2023) for the math domain.

**Evaluation** We select the following tasks for evaluation, covering fundamental capabilities in LLMs (e.g., knowledge, reasoning, coding): RACE (Lai et al., 2017), TriviaQA (Joshi et al., 2017), BigBench-Challenge (BBH) (Suzgun et al., 2022), ARC-Challenge (ARC) (Clark et al., 2018), Hellaswag (Zellers et al., 2019), and HumanEval (Chen et al., 2021). The evaluation settings for these benchmarks are listed in Tab. 2. We adopt lm-evaluation-harness (Gao et al., 2023b) for unified evaluation.

Table 2: The evaluation settings of the benchmarks.

Dataset	Evaluation Type	Evaluation Method	Metric	Random Performance
ARC	Multiple Choice	10-shot	Accuracy	25
BBH	Generation	CoT-3-shot	ExactMatch	0
Hellaswag	Multiple Choice	10-shot	Accuracy	25
HumanEval	Generation	0-shot	Pass@100	0
RACE	Multiple Choice	0-shot	Accuracy	25
TriviaQA	Generation	0-shot	ExactMatch	0

### 4.3 EXPERIMENTAL SETTING

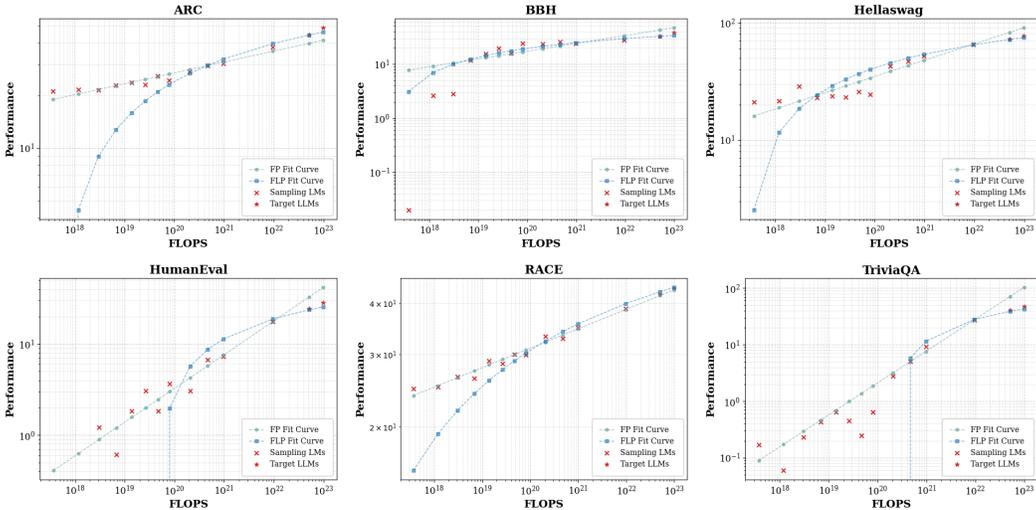
**Baseline** We consider directly using the expended FLOPs  $C$  to predict the downstream performance  $P$ , and experiment with the following analytical form for comparison:

$$P(C) = \left(\frac{C}{C_M}\right)^{\alpha_M}, \quad (3)$$

where  $C_M$  and  $\alpha_M$  are constant terms to be estimated. We denote this approach as **FP**.

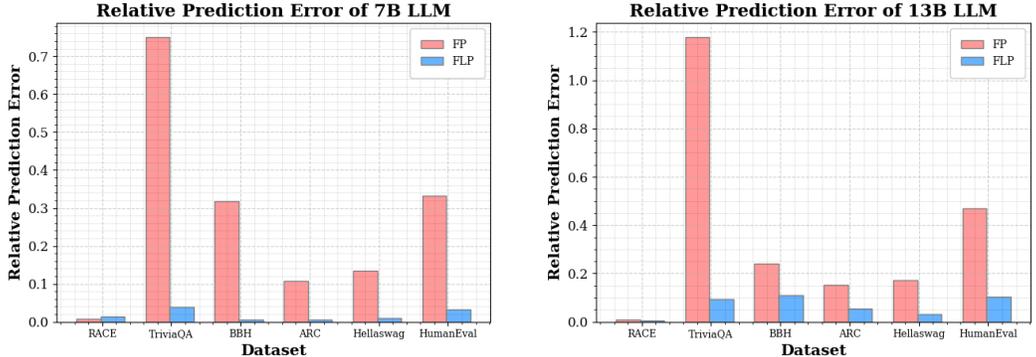
**Implementation of FLP** To fit the FLOPs-to-Loss curve, we utilize the final checkpoints from each sampling LM. In addition, during LMs training, a checkpoint is saved at every 1/30th increment

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232



233 Figure 2: The downstream performance prediction using FP and FLP fit curves. FLP can better  
234 predict the downstream performance of target 7B and 13B LLMs across all evaluation benchmarks.

235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246



247 Figure 3: The relative prediction error of 7B and 13B LLMs. FLP achieves a more accurate  
248 prediction with error margins of 5% and 10% across all benchmarks for two LLMs respectively.

249 of the total training progress. We monitor and record the pre-training loss on the training dataset,  
250 rounded to two decimal places. Only those checkpoints demonstrating an improvement in pre-training  
251 loss are retained. For these selected checkpoints, we evaluate the downstream performance and pre-  
252 training loss on the validation set. We then discard those that do not surpass the random benchmark  
253 performance by at least 5, and use the remaining data points to fit the Loss-to-Performance curve.

254 **Evaluation Metrics** In addition to presenting the fitting curves for intuitive visualization, we quantify  
255 the prediction accuracy by measuring the relative prediction error:

$$257 \text{ Relative Prediction Error} = \frac{|\text{Predictive Metric} - \text{Actual Metric}|}{\text{Actual Metric}} \quad (4)$$

258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

#### 4.4 RESULTS

The downstream performance prediction results are visualized in Fig. 2. Across all evaluation tasks,  
FLP fit curve can better predict the performance of target LLMs with 7B and 13B parameters using the  
sampling LMs up to 3B. In contrast, while FP more effectively fits the data points of sampling LMs,  
it has difficulty accounting for the “emergent phase” characterized by rapid performance shifts, due to  
the scarcity of data points from this period. As a solution, FLP utilizes pre-training loss as a more fine-  
grained indicator to monitor performance changes and effectively incorporates data from intermediate  
checkpoints, enhancing sample efficiency. The evaluation results of relative prediction error are  
shown in Fig. 3. Unlike the suboptimal predictions of FP, FLP delivers precise forecasts, maintaining  
relative error margins of 5% and 10% across all benchmarks for 7B and 13B LLMs, respectively.

Compared to FP, FLP is less effective at fitting the data points of sampling LMs, especially in HumanEval and TriviaQA. The reason is that we do not align with the “non-emergent” phase of the Loss-to-Performance curve, where LMs exhibit random performance when pre-training loss is beyond the task-specific threshold. Thus, FLP predicts higher pre-training loss for LMs with fewer FLOPs, resulting in below-random performance. This issue is not within the scope of FLP, as it is specifically designed to predict the performance of LLMs trained with significantly larger FLOPs in practice.

In addition, we discuss additional results in Appendix for the presentation purpose since adding these data points may distort the vertical axis scaling in Fig. 2. We compare FLP further with the analytical forms and approaches proposed in GPT-4 (Achiam et al., 2023) and Llama-3 (Dubey et al., 2024) technical reports. The results are shown in §B and §C respectively. We also evaluate the feasibility of employing FLP to predict the performance of a 13B LLM on MMLU (Hendrycks et al., 2020), using intermediate checkpoints from a 7B LLM (§D). Overall, the results demonstrate the general effectiveness and applicability of FLP.

## 5 FLP-M: DATA MIXING FOR DOWNSTREAM PERFORMANCE PREDICTION

Motivated by the encouraging results of FLP (§4), we propose FLP-M, a fundamental approach to meet the practical needs of integrating data from various sources (Groeneveld et al., 2024; Penedo et al., 2024). In our work, we focus on mixing general corpus with code data, considering two distinct yet overlapping data sources. This intersection offers a more realistic perspective than treating them as distinct domains (Ye et al., 2024), as real-world corpus often spans multiple domains, necessitating an analysis of the interdependence between data sources when formulating our analytical functions.

Compared to the straightforward implementation of FLP (§3), FLP-M operates on fine-grained, domain-specific pre-training loss, due to the observation that the average loss on the entire validation set fails to effectively reflect performance variations in downstream tasks in the data mixing context (§7.2). This may be due to the fact that changes in pre-training data mixtures simultaneously impact multiple capabilities of the LMs. For instance, an increase in code data loss coupled with a decrease in general data loss may leave the average validation loss unchanged, yet result in LMs with distinct capabilities and downstream performance. Note that unlike the pre-training data mixture, the validation set is deliberately curated by domain, as creating smaller, domain-specific validation sets is manageable.

### 5.1 FLOPS $\rightarrow$ DOMAIN LOSS

Given the FLOPs  $C^G$  spent on the general corpus and  $C^C$  spent on the code data, we naturally extend the power law function to the following analytical form to predict the domain-specific pre-training loss  $L^D$  on domain  $D$ :

$$L^D(C^G, C^C) = \left(\frac{C^G + C^C}{C_T}\right)^{\alpha_C} \times \left(\frac{C^G}{C_G}\right)^{\alpha_{C_1}} \times \left(\frac{C^C}{C_C}\right)^{\alpha_{C_2}} \quad (5)$$

where  $C_T$ ,  $C_G$ ,  $C_C$ ,  $\alpha_C$ ,  $\alpha_{C_1}$ , and  $\alpha_{C_2}$  are constants to be estimated. In FLP-M, we first select a sequence of total compute  $\{C_i\}_{i=1}^N$  spent on pre-training. For each selected  $C_i$ , we experiment with various ratios to mix two data sources, and decompose  $C_i$  into  $C_i^G$  and  $C_i^C$ . We measure the domain-specific pre-training loss  $L_i^D$  on a domain-specific subset  $D$  of validation data to obtain  $(C_i^G, C_i^C, L_i^D)$  data pairs. Then we can estimate the constants in Eq. 5. We also experiment with other potential analytical forms in §7.2.

### 5.2 DOMAIN LOSS $\rightarrow$ PERFORMANCE

Given the pre-training loss  $\{L^D\}_{D=1}^K$  on  $K$  domains, we train a two-layer neural network with a hidden layer size of 3 and the ReLU activation function (Agarap, 2018) to predict the downstream performance. The network is optimized using the regression loss with  $L_2$  regularization and the Adam optimizer (Diederik, 2014), employing a learning rate of 0.05 that linearly decays to 0 within 2,000 steps and a weight decay of 0.01. In FLP-M, we adopt the same strategy as in FLP to fetch the intermediate checkpoints and only retain the results that the LMs achieve above-random performance (see §3). Thus, for  $LM_i$ , we can obtain a sequence of effective data points  $(\{L_i^D\}_{D=1}^K, P_i)$ , where  $L_i^D$  is the pre-training loss on domain  $D$  and  $P_i$  is the LM’s performance. Then we can use these data points to train the neural network. We also explore other functions for fitting in §7.2.

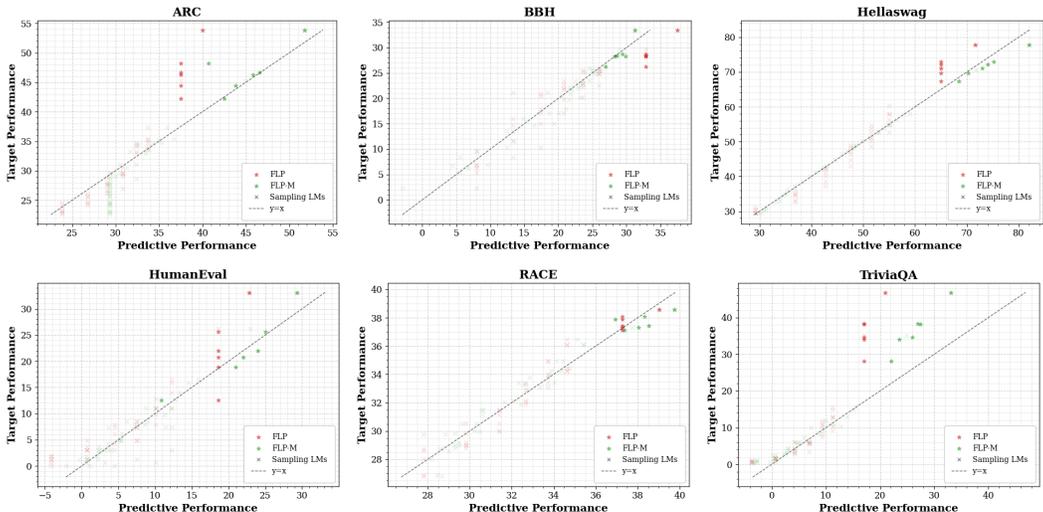


Figure 4: The downstream performance prediction using FLP and FLP-M fit curves. FLP-M can better predict the downstream performance of target LLMs across various data mixing ratios.

## 6 EXPERIMENT FOR FLP-M

### 6.1 SAMPLING AND TARGET LMS

We train a series of sampling LMs with sizes of  $\{0.12B, 0.2B, 0.32B, 0.5B, 0.72B, 1B\}$ , and the corresponding training token numbers are shown in Tab. 1. We train the LMs on the general and code data mixture with  $\{0, 0.1, 0.2, 0.3, 0.4, 0.5\}$  as the mixing ratios of code data to reflect real-world usage. We also add one sampling LM of 3B size and 0.3 mixing ratio. For evaluation, we train 3B LLMs with the other mixing ratios and a 7B LLM with 0.3 as the mixing ratio due to the limited compute budget.

### 6.2 DATA: PRE-TRAINING, VALIDATION, EVALUATION

**Pre-Training** For general corpus, we use DCLM (Li et al., 2024b), a curated high-quality pre-training corpus including heuristic cleaning, filtering, deduplication, and model-based filtering. For code data, we use The Stack v2 (Lozhkov et al., 2024), which initially contains over 3B files in 600+ programming and markup languages, created as part of the BigCode project. We mix these two data sources to create the pre-training data mixture using the ratios specified in §6.1.

**Validation** We use the same validation data mixture specified in §4.2 that includes 5 distinct domains.

**Evaluation** The evaluation benchmarks and settings are the same as those in §4.2.

### 6.3 EXPERIMENTAL SETTING

**Baseline** We implement FLP within this data mixing context as a baseline, which first predicts the average pre-training loss on the validation set and uses this to estimate downstream performance via linear regression.

**Implementation of FLP-M** We adopt the same implementation as in FLP (details in §4.3). The distinction is that we individually measure the pre-training loss on each domain of the validation mixture.

### 6.4 RESULTS

The downstream performance prediction results are visualized in Fig. 4. We update the x-axis to “predicted performance” to improve clarity, as the presence of two variables ( $C^G, C^C$ ) complicated 3D visualization. Overall, we find that FLP-M demonstrates better performance compared to FLP when considering the data mixing as an extra factor in scaling laws analysis. Using average validation loss as an indicator for assessing the performance of LMs pre-trained on mixed data sources, such

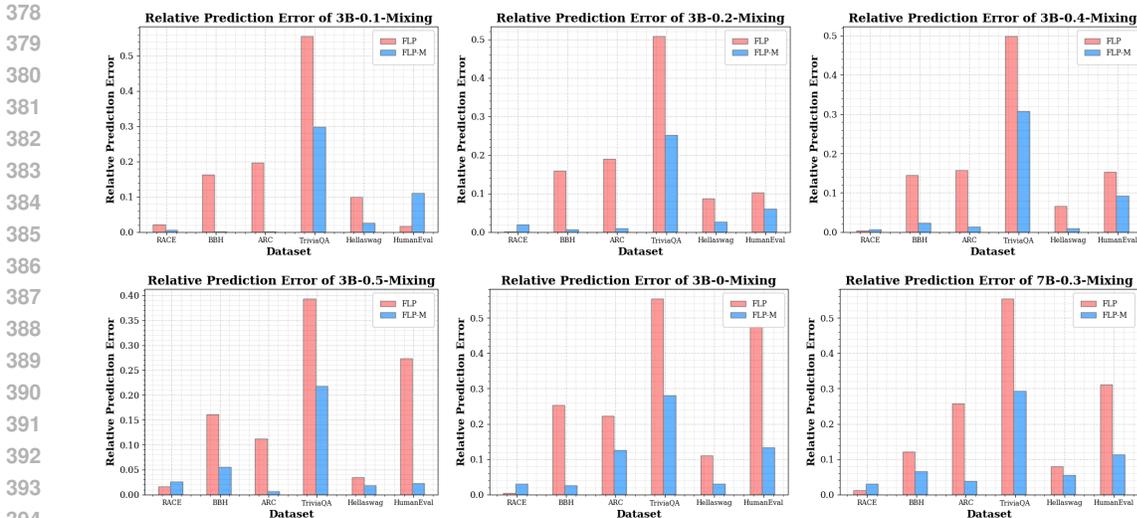


Figure 5: The relative prediction error of downstream performance prediction using FLP and FLP-M. FLP-M can better predict the performance of target LLMs across various data mixing ratios.

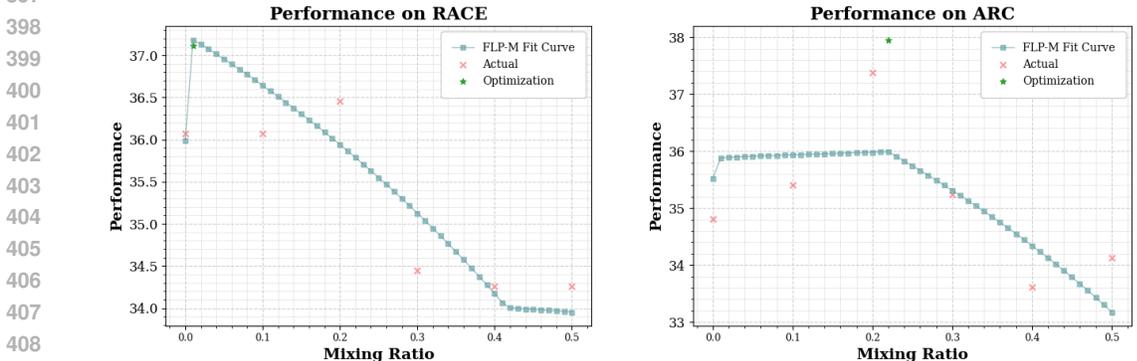


Figure 6: We use the scaling laws function derived via FLP-M to find the optimal data mixing ratio that yields the estimated best performance on the corresponding benchmarks.

as general text and code, is limited. Thus, the average loss fails to trace performance variations in downstream tasks because changes in data mixtures can affect different capabilities of the LMs. In contrast, FLP-M effectively leverages the domain-specific validation loss to capture the capabilities improvement in LMs, and thus can better predict the downstream performance. In our experiments, FLP-M accurately predicts the performance of 3B LLMs across various data mixtures and the 7B LLM with 0.3 data mixing ratio with error margins within 10% for most benchmarks.

However, on TriviaQA, despite significantly outperforming FLP, FLP-M shows higher relative prediction error, ranging from 20% to 30%. This discrepancy can be explained by the substantial performance improvement when scaling LLMs from under 1B to 3B parameters (increasing from below 12 to over 28). In our sampling LMs configurations (see Tab. 1), we lack sufficient data points to adequately characterize the phase of accelerated performance improvement. To better model this trend, a practical solution is to add several sampling LMs between 1B and 3B parameters.

## 7 FURTHER ANALYSIS

### 7.1 OPTIMIZING DATA MIXTURE USING FLP-M SCALING LAWS

We demonstrate how the derived scaling laws using FLP-M can be effectively applied to optimize data mixtures, enhancing downstream performance. We focus on 1B LMs in this analysis due to compute constraints. For each dataset, we use the FLP-M to estimate the function that maps expended

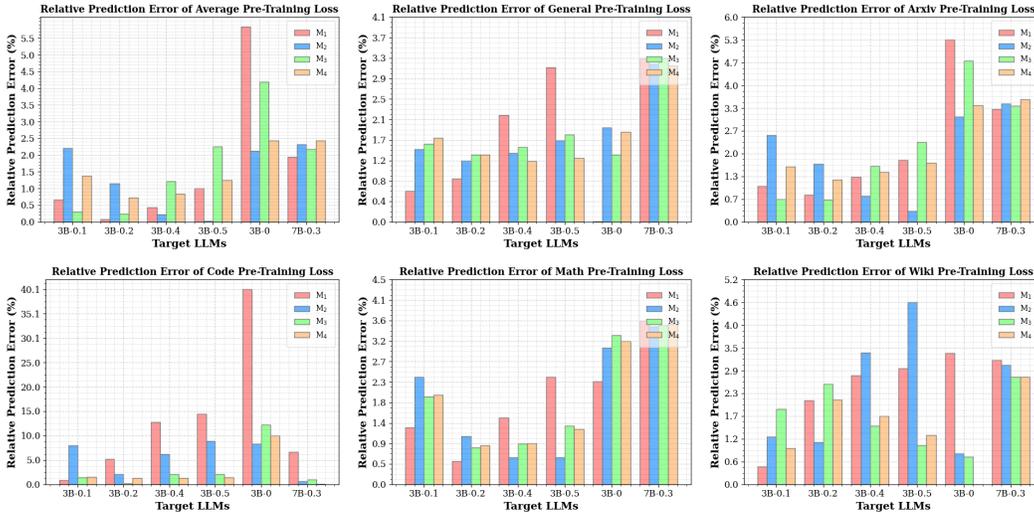


Figure 7: The relative prediction error of average and domain-specific pre-training loss.  $M_4$  provides more stable and overall more accurate predictions for domain-specific loss (within 2.5% relative prediction error across most domains).

FLOPs in each data source to the downstream performance. Then we use this function to predict performance across mixing ratios from 0 to 0.5, in intervals of 0.01.

Among all evaluation datasets, the estimated scaling laws function exhibits non-monotonic behavior on the RACE and ARC datasets, reaching its peak at mixing ratios of 0.01 and 0.22, respectively. To verify, we train 1B LMs with these two mixing ratios and measure their performance on the corresponding benchmarks. The results are shown in Fig. 6. We find that the selected optimal mixing ratio can reliably yield better performance compared to the six mixing ratios adopted for the sampling LMs, highlighting  $\text{FLP-M}$  as a practical approach for optimizing data mixtures to enhance performance on specific target tasks.

## 7.2 ABLATION STUDY

We conduct further analysis to better understand the two stages in  $\text{FLP-M}$ . Specifically, we compare various approaches to estimate the FLOPs-to-Loss and Loss-to-Performance curves in  $\text{FLP-M}$ .

**FLP-M: FLOPs  $\rightarrow$  Loss** We experiment with several candidate analytical forms listed in Tab. 3. We assess their performance in estimating the average pre-training loss across the entire validation set, as well as the domain-specific pre-training losses on corresponding subsets. We present the fit curves in Fig. 13 (§E), and the relative prediction errors for pre-training loss estimation are shown in Fig. 7. For average pre-training loss prediction, using more complex analytical models that account for the individual impact of each data source can lead to performance degradation. However, relying solely on the total compute for prediction ( $M_1$ ) can cause high prediction errors in certain domains (e.g., code) and are not stable for various mixing ratios. More complex analytical models generally perform better in predicting domain-specific loss. Among them,  $M_4$ , the adopted model in  $\text{FLP-M}$ , provides more stable (within 2.5% relative prediction error across most domains) and overall more accurate predictions (achieving the lowest average error shown in Tab. 3).

Table 3: Candidate analytical forms for fitting the FLOPs-to-Loss curve. Except for  $C^G$  and  $C^C$  representing the compute used for general and code data sources, other constants need to be estimated. The average error is computed across all domains and model types.

$L^P(C^G, C^C) =$	Analytical Form	Average Error
$M_1$	$(\frac{C^G+C^C}{C_T})^{\alpha_C}$	0.029
$M_2$	$(\frac{C^G}{C_C})^{\alpha_{C1}} \times (\frac{C^C}{C_C})^{\alpha_{C2}}$	0.026
$M_3$	$(\frac{w_0 \cdot C^G + w_1 \cdot C^C}{C_T})^{\alpha_C}$	0.017
$M_4$ (Ours)	$(\frac{C^G+C^C}{C_T})^{\alpha_C} \times (\frac{C^G}{C_G})^{\alpha_{C1}} \times (\frac{C^C}{C_C})^{\alpha_{C2}}$	<b>0.014</b>

**FLP-M: Loss  $\rightarrow$  Performance** We experiment with various approaches to estimate the function that maps the pre-training loss to the downstream performance. In this study, we utilize the actual

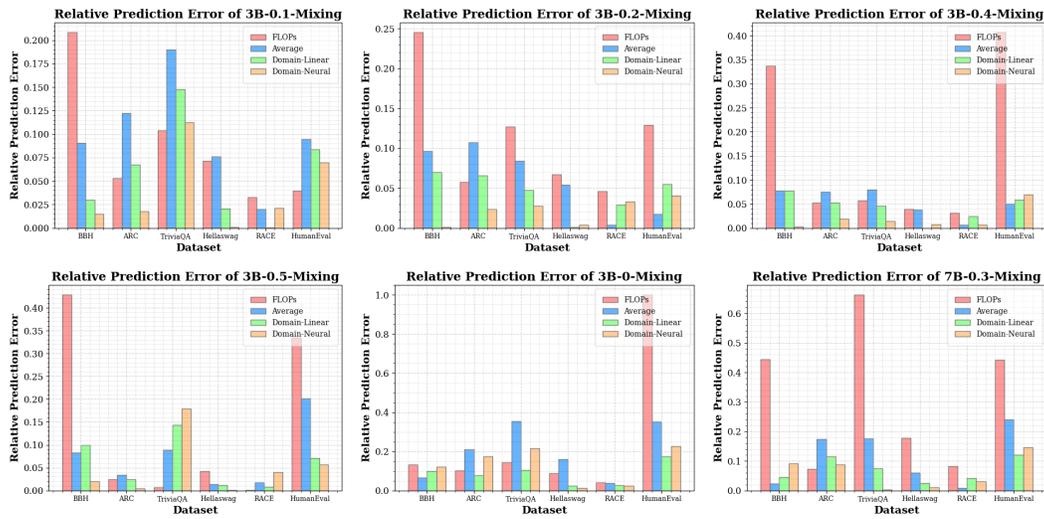


Figure 8: The relative prediction error of various approaches to estimate the Loss-to-Performance curve. Neural network estimation with domain-specific loss as input achieves the best prediction.

pre-training loss of target LLMs, rather than the predictive loss used in §5. We consider the following candidates with different inputs:

- (1) **FLOPs:** We adopt the analytical form used to predict the pre-training loss based on training compute (see Eq. 5), only changing the target metric to the downstream performance.
- (2) **Average Loss (Average):** We implement a linear regression model to map the average pre-training loss on the whole validation set to the downstream performance.
- (3) **Domain Loss via Linear Combination (Domain-Linear):** We apply a linear regression model to correlate pre-training loss across domains with downstream performance.
- (4) **Domain Loss via Neural Network (Ours) (Domain-Neural):** We implement a two-layer neural network to map the pre-training loss across domains to the downstream performance. The network configuration and optimization process are introduced in §5.

The fit curves are shown in Fig. 14 (§E) and the results of relative prediction error are shown in Fig. 8. Consistent with the findings in §4, directly estimating the performance based on expended compute (FLOPs) leads to highly inaccurate predictions (FLOPs vs. Loss). Pre-training loss serves as a more reliable metric for performance estimation, and decomposing it into domain-specific loss can further enhance prediction accuracy (Average vs. Domain Loss). For the predictive models, using neural network estimation can better leverage the abundant data points produced by FLP-M, resulting in better performance compared to the linear regression model (Linear vs. Neural Network).

## 8 CONCLUSION

This paper introduces a two-stage FLP solution to predict downstream performance in LLMs by leveraging pre-training loss. Encouraged by promising preliminary results, we propose FLP-M, a core solution for performance prediction that addresses the practical challenges of integrating pre-training data from diverse sources. The effectiveness of FLP-M is validated through extensive experiments.

## LIMITATIONS

Our approach FLP-M is generally applicable across various data sources, yet currently, it is demonstrated only in binary cases involving code and text data due to computational constraints. Our specific emphasis on the mixing ratio of code is deliberate, reflecting its practical significance in real-world applications. This limitation marks a key area for future expansion.

## REFERENCES

- 540  
541  
542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
543 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.  
544 arXiv preprint arXiv:2303.08774, 2023.
- 545 Abien Fred Agarap. Deep learning using rectified linear units (relu). arXiv preprint arXiv:1803.08375,  
546 2018.
- 547 Ibrahim M Alabdulmohsin, Behnam Neyshabur, and Xiaohua Zhai. Revisiting neural scaling laws  
548 in language and vision. Advances in Neural Information Processing Systems, 35:22300–22312,  
549 2022.
- 550 Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.3, knowledge capacity scaling  
551 laws. arXiv preprint arXiv:2404.05405, 2024.
- 552  
553 Yasaman Bahri, Ethan Dyer, Jared Kaplan, Jaehoon Lee, and Utkarsh Sharma. Explaining neural  
554 scaling laws. Proceedings of the National Academy of Sciences, 121(27):e2311878121, 2024.
- 555  
556 Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding,  
557 Kai Dong, Qiusi Du, Zhe Fu, et al. Deepseek llm: Scaling open-source language models with  
558 longtermism. arXiv preprint arXiv:2401.02954, 2024.
- 559  
560 Dillon Bowen, Brendan Murphy, Will Cai, David Khachaturov, Adam Gleave, and Kellin Pelrine.  
561 Scaling laws for data poisoning in llms, 2024. URL [https://arxiv.org/abs/2408.](https://arxiv.org/abs/2408.02946)  
562 02946.
- 563 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
564 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
565 few-shot learners. Advances in neural information processing systems, 33:1877–1901, 2020.
- 566  
567 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared  
568 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri,  
569 Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan,  
570 Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian,  
571 Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios  
572 Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino,  
573 Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders,  
574 Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa,  
575 Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob  
576 McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating  
large language models trained on code. 2021.
- 577  
578 Aidan Clark, Diego de Las Casas, Aurelia Guy, Arthur Mensch, Michela Paganini, Jordan Hoffmann,  
579 Bogdan Damoc, Blake Hechtman, Trevor Cai, Sebastian Borgeaud, et al. Unified scaling laws for  
580 routed language models. In International conference on machine learning, pp. 4057–4086. PMLR,  
581 2022.
- 582 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and  
583 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.  
584 arXiv preprint arXiv:1803.05457, 2018.
- 585  
586 Together Computer. Redpajama: An open source recipe to reproduce llama training dataset, 2023.  
587 URL <https://github.com/togethercomputer/RedPajama-Data>.
- 588 P Kingma Diederik. Adam: A method for stochastic optimization. 2014.
- 589  
590 Zhengxiao Du, Aohan Zeng, Yuxiao Dong, and Jie Tang. Understanding emergent abilities of  
591 language models from the loss perspective. arXiv preprint arXiv:2403.15796, 2024.
- 592  
593 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.  
arXiv preprint arXiv:2407.21783, 2024.

- 594 Elias Frantar, Carlos Riquelme, Neil Houlsby, Dan Alistarh, and Utku Evci. Scaling laws for  
595 sparsely-connected foundation models. [arXiv preprint arXiv:2309.08520](#), 2023.
- 596
- 597 Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. In  
598 [International Conference on Machine Learning](#), pp. 10835–10866. PMLR, 2023a.
- 599 Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster,  
600 Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff,  
601 Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika,  
602 Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot  
603 language model evaluation, 12 2023b. URL <https://zenodo.org/records/10256836>.
- 604
- 605 Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya  
606 Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. [arXiv preprint](#)  
607 [arXiv:2406.04093](#), 2024.
- 608 Sachin Goyal, Pratyush Maini, Zachary C Lipton, Aditi Raghunathan, and J Zico Kolter. Scaling  
609 laws for data filtering—data curation cannot be compute agnostic. In [Proceedings of the IEEE/CVF](#)  
610 [Conference on Computer Vision and Pattern Recognition](#), pp. 22702–22711, 2024.
- 611 Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord,  
612 Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, et al. Olmo: Accelerat-  
613 ing the science of language models. [arXiv preprint arXiv:2402.00838](#), 2024.
- 614
- 615 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and  
616 Jacob Steinhardt. Measuring massive multitask language understanding. [arXiv preprint](#)  
617 [arXiv:2009.03300](#), 2020.
- 618 Tom Henighan, Jared Kaplan, Mor Katz, Mark Chen, Christopher Hesse, Jacob Jackson, Heewoo  
619 Jun, Tom B Brown, Prafulla Dhariwal, Scott Gray, et al. Scaling laws for autoregressive generative  
620 modeling. [arXiv preprint arXiv:2010.14701](#), 2020.
- 621
- 622 Danny Hernandez, Jared Kaplan, Tom Henighan, and Sam McCandlish. Scaling laws for transfer.  
623 [arXiv preprint arXiv:2102.01293](#), 2021.
- 624
- 625 Danny Hernandez, Tom Brown, Tom Conerly, Nova DasSarma, Dawn Drain, Sheer El-Showk, Nelson  
626 Elhage, Zac Hatfield-Dodds, Tom Henighan, Tristan Hume, et al. Scaling laws and interpretability  
627 of learning from repeated data. [arXiv preprint arXiv:2205.10487](#), 2022.
- 628
- 629 Joel Hestness, Sharan Narang, Newsha Ardalani, Gregory Diamos, Heewoo Jun, Hassan Kianinejad,  
630 Md Mostofa Ali Patwary, Yang Yang, and Yanqi Zhou. Deep learning scaling is predictable,  
631 empirically. [arXiv preprint arXiv:1712.00409](#), 2017.
- 632
- 633 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza  
634 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al.  
635 Training compute-optimal large language models. [arXiv preprint arXiv:2203.15556](#), 2022.
- 636
- 637 Shengding Hu, Xin Liu, Xu Han, Xinrong Zhang, Chaoqun He, Weilin Zhao, Yankai Lin, Ning Ding,  
638 Zebin Ou, Guoyang Zeng, et al. Predicting emergent abilities with infinite resolution evaluation.  
639 In [The Twelfth International Conference on Learning Representations](#), 2023.
- 640
- 641 Yuzhen Huang, Jinghan Zhang, Zifei Shan, and Junxian He. Compression represents intelligence  
642 linearly. [arXiv preprint arXiv:2404.09937](#), 2024.
- 643
- 644 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,  
645 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.  
646 Mistral 7b. [arXiv preprint arXiv:2310.06825](#), 2023.
- 647
- 648 Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly  
649 supervised challenge dataset for reading comprehension. [arXiv preprint arXiv:1705.03551](#), 2017.
- 650
- 651 Feiyang Kang, Yifan Sun, Bingbing Wen, Si Chen, Dawn Song, Rafid Mahmood, and Ruoxi Jia.  
652 Autoscale: Automatic prediction of compute-optimal data composition for training llms. [arXiv](#)  
653 [preprint arXiv:2407.20177](#), 2024.

- 648 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott  
649 Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models.  
650 [arXiv preprint arXiv:2001.08361](https://arxiv.org/abs/2001.08361), 2020.
- 651
- 652 Jakub Krajewski, Jan Ludziejewski, Kamil Adamczewski, Maciej Pióro, Michał Krutul, Szymon  
653 Antoniak, Kamil Ciebiera, Krystian Król, Tomasz Odrzygóźdź, Piotr Sankowski, et al. Scaling  
654 laws for fine-grained mixture of experts. [arXiv preprint arXiv:2402.07871](https://arxiv.org/abs/2402.07871), 2024.
- 655
- 656 Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. Race: Large-scale reading  
657 comprehension dataset from examinations. [arXiv preprint arXiv:1704.04683](https://arxiv.org/abs/1704.04683), 2017.
- 658
- 659 Bozhou Li, Hao Liang, Zimo Meng, and Wentao Zhang. Are bigger encoders always better in vision  
660 large models?, 2024a. URL <https://arxiv.org/abs/2408.00620>.
- 661
- 662 Jeffrey Li, Alex Fang, Georgios Smyrnis, Maor Ivgi, Matt Jordan, Samir Gadre, Hritik Bansal, Etash  
663 Guha, Sedrick Keh, Kushal Arora, Saurabh Garg, Rui Xin, Niklas Muennighoff, Reinhard Heckel,  
664 Jean Mercat, Mayee Chen, Suchin Gururangan, Mitchell Wortsman, Alon Albalak, Yonatan Bitton,  
665 Marianna Nezhurina, Amro Abbas, Cheng-Yu Hsieh, Dhruva Ghosh, Josh Gardner, Maciej Kilian,  
666 Hanlin Zhang, Rulin Shao, Sarah Pratt, Sunny Sanyal, Gabriel Ilharco, Giannis Daras, Kalyani  
667 Marathe, Aaron Gokaslan, Jieyu Zhang, Khyathi Chandu, Thao Nguyen, Igor Vasiljevic, Sham  
668 Kakade, Shuran Song, Sujay Sanghavi, Fartash Faghri, Sewoong Oh, Luke Zettlemoyer, Kyle Lo,  
669 Alaaeldin El-Nouby, Hadi Pouransari, Alexander Toshev, Stephanie Wang, Dirk Groeneveld, Luca  
670 Soldaini, Pang Wei Koh, Jenia Jitsev, Thomas Kollar, Alexandros G. Dimakis, Yair Carmon, Achal  
671 Dave, Ludwig Schmidt, and Vaishaal Shankar. Datacomp-lm: In search of the next generation of  
672 training sets for language models, 2024b.
- 673
- 674 Haowei Lin, Baizhou Huang, Haotian Ye, Qinyu Chen, Zihao Wang, Sujian Li, Jianzhu Ma, Xiaojun  
675 Wan, James Zou, and Yitao Liang. Selecting large language model to fine-tune via rectified scaling  
676 law. [arXiv preprint arXiv:2402.02314](https://arxiv.org/abs/2402.02314), 2024.
- 677
- 678 Lucas Lingle. A large-scale exploration of  $\mu$ -transfer. [arXiv preprint arXiv:2404.05728](https://arxiv.org/abs/2404.05728), 2024.
- 679
- 680 Qian Liu, Xiaosen Zheng, Niklas Muennighoff, Guangtao Zeng, Longxu Dou, Tianyu Pang, Jing  
681 Jiang, and Min Lin. Regmix: Data mixture as regression for language model pre-training. [arXiv  
682 preprint arXiv:2407.01492](https://arxiv.org/abs/2407.01492), 2024.
- 683
- 684 Zhengzhong Liu, Aurick Qiao, Willie Neiswanger, Hongyi Wang, Bowen Tan, Tianhua Tao, Junbo  
685 Li, Yuqi Wang, Suqi Sun, Omkar Pangarkar, et al. Llm360: Towards fully transparent open-source  
686 llms. [arXiv preprint arXiv:2312.06550](https://arxiv.org/abs/2312.06550), 2023.
- 687
- 688 Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane  
689 Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, Tianyang Liu, Max Tian, Denis Kocetkov,  
690 Arthur Zucker, Younes Belkada, Zijian Wang, Qian Liu, Dmitry Abulkhanov, Indraneil Paul,  
691 Zhuang Li, Wen-Ding Li, Megan Risdal, Jia Li, Jian Zhu, Terry Yue Zhuo, Evgenii Zheltonozhskii,  
692 Nii Osa Osa Dade, Wenhao Yu, Lucas Krauß, Naman Jain, Yixuan Su, Xuanli He, Manan  
693 Dey, Edoardo Abati, Yekun Chai, Niklas Muennighoff, Xiangru Tang, Muhtasham Oblokulov,  
694 Christopher Akiki, Marc Marone, Chenghao Mou, Mayank Mishra, Alex Gu, Binyuan Hui, Tri  
695 Dao, Armel Zebaze, Olivier Dehaene, Nicolas Patry, Canwen Xu, Julian McAuley, Han Hu, Torsten  
696 Scholak, Sebastien Paquet, Jennifer Robinson, Carolyn Jane Anderson, Nicolas Chapados, Mostofa  
697 Patwary, Nima Tajbakhsh, Yacine Jernite, Carlos Muñoz Ferrandis, Lingming Zhang, Sean Hughes,  
698 Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. Starcoder 2 and the stack v2:  
699 The next generation, 2024.
- 700
- 701 Qian Ma, Haitao Mao, Jingzhe Liu, Zhehua Zhang, Chunlin Feng, Yu Song, Yihan Shao, Tianfan  
Fu, and Yao Ma. Do neural scaling laws exist on graph self-supervised learning?, 2024. URL  
<https://arxiv.org/abs/2408.11243>.
- 702
- 703 Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier  
Amatriain, and Jianfeng Gao. Large language models: A survey. [arXiv preprint arXiv:2402.06196](https://arxiv.org/abs/2402.06196),  
2024.

- 702 Niklas Muennighoff, Alexander Rush, Boaz Barak, Teven Le Scao, Nouamane Tazi, Aleksandra  
703 Piktus, Sampo Pyysalo, Thomas Wolf, and Colin A Raffel. Scaling data-constrained language  
704 models. *Advances in Neural Information Processing Systems*, 36, 2024.
- 705  
706 Guilherme Penedo, Hynek Kydlíček, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell, Colin  
707 Raffel, Leandro Von Werra, and Thomas Wolf. The fineweb datasets: Decanting the web for the  
708 finest text data at scale, 2024.
- 709 Haoran Que, Jiaheng Liu, Ge Zhang, Chenchen Zhang, Xingwei Qu, Yinghao Ma, Feiyu Duan,  
710 Zhiqi Bai, Jiakai Wang, Yuanxing Zhang, et al. D-cpt law: Domain-specific continual pre-training  
711 scaling law for large language models. *arXiv preprint arXiv:2406.01375*, 2024.
- 712  
713 Rafael Rafailov, Yaswanth Chittooru, Ryan Park, Harshit Sikchi, Joey Hejna, Bradley Knox, Chelsea  
714 Finn, and Scott Niekum. Scaling laws for reward model overoptimization in direct alignment  
715 algorithms. *arXiv preprint arXiv:2406.02900*, 2024.
- 716 Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. Are emergent abilities of large language  
717 models a mirage? *Advances in Neural Information Processing Systems*, 36, 2024.
- 718  
719 Rulin Shao, Jacqueline He, Akari Asai, Weijia Shi, Tim Dettmers, Sewon Min, Luke Zettlemoyer,  
720 and Pang Wei Koh. Scaling retrieval-based language models with a trillion-token datastore. *arXiv*  
721 *preprint arXiv:2407.12854*, 2024.
- 722  
723 Zhiqiang Shen, Tianhua Tao, Liqun Ma, Willie Neiswanger, Joel Hestness, Natalia Vassilieva, Daria  
724 Soboleva, and Eric Xing. Slimpajama-dc: Understanding data combinations for llm training. *arXiv*  
*preprint arXiv:2309.10818*, 2023.
- 725  
726 Hui Su, Zhi Tian, Xiaoyu Shen, and Xunliang Cai. Unraveling the mystery of scaling laws: Part i.  
727 *arXiv preprint arXiv:2403.06563*, 2024.
- 728  
729 Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung,  
730 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.
- 731  
732 Chaofan Tao, Qian Liu, Longxu Dou, Niklas Muennighoff, Zhongwei Wan, Ping Luo, Min Lin, and  
733 Ngai Wong. Scaling laws with vocabulary: Larger models deserve larger vocabularies. *arXiv*  
*preprint arXiv:2407.13623*, 2024.
- 734  
735 Yi Tay, Mostafa Dehghani, Jinfeng Rao, William Fedus, Samira Abnar, Hyung Won Chung, Sharan  
736 Narang, Dani Yogatama, Ashish Vaswani, and Donald Metzler. Scale efficiently: Insights from  
737 pre-training and fine-tuning transformers. *arXiv preprint arXiv:2109.10686*, 2021.
- 738  
739 Howe Tissue, Venus Wang, and Lu Wang. Scaling law with learning rate annealing, 2024. URL  
<https://arxiv.org/abs/2408.11029>.
- 740  
741 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay  
742 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation  
743 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 744  
745 Zhongwei Wan, Xin Wang, Che Liu, Samiul Alam, Yu Zheng, Zhongnan Qu, Shen Yan, Yi Zhu,  
746 Quanlu Zhang, Mosharaf Chowdhury, et al. Efficient large language models: A survey. *arXiv*  
*preprint arXiv:2312.03863*, 1, 2023.
- 747  
748 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama,  
749 Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models.  
*arXiv preprint arXiv:2206.07682*, 2022.
- 750  
751 Tianwen Wei, Liang Zhao, Lichang Zhang, Bo Zhu, Lijie Wang, Haihua Yang, Biye Li, Cheng Cheng,  
752 Weiwei Lü, Rui Hu, et al. Skywork: A more open bilingual foundation model. *arXiv preprint*  
753 *arXiv:2310.19341*, 2023.
- 754  
755 Greg Yang, Edward J Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick Ryder,  
Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tensor programs v: Tuning large neural networks  
via zero-shot hyperparameter transfer. *arXiv preprint arXiv:2203.03466*, 2022.

756 Yiqun Yao, Xiusheng Huang, Xuezhi Fang, Xiang Li, Ziyi Ni, Xin Jiang, Xuying Meng, Peng Han,  
757 Shuo Shang, Kang Liu, et al. nanolm: an affordable llm pre-training benchmark via accurate loss  
758 prediction across scales. [arXiv preprint arXiv:2304.06875](#), 2023.

759  
760 Jiasheng Ye, Peiju Liu, Tianxiang Sun, Yunhua Zhou, Jun Zhan, and Xipeng Qiu. Data mixing  
761 laws: Optimizing data mixtures by predicting language modeling performance. [arXiv preprint](#)  
762 [arXiv:2403.16952](#), 2024.

763 Longfei Yun, Yonghao Zhuang, Yao Fu, Eric P Xing, and Hao Zhang. Toward inference-optimal  
764 mixture-of-expert large language models. [arXiv preprint arXiv:2404.02852](#), 2024.

765  
766 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine  
767 really finish your sentence? [arXiv preprint arXiv:1905.07830](#), 2019.

768 Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. In  
769 [Proceedings of the IEEE/CVF conference on computer vision and pattern recognition](#), pp. 12104–  
770 12113, 2022.

771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

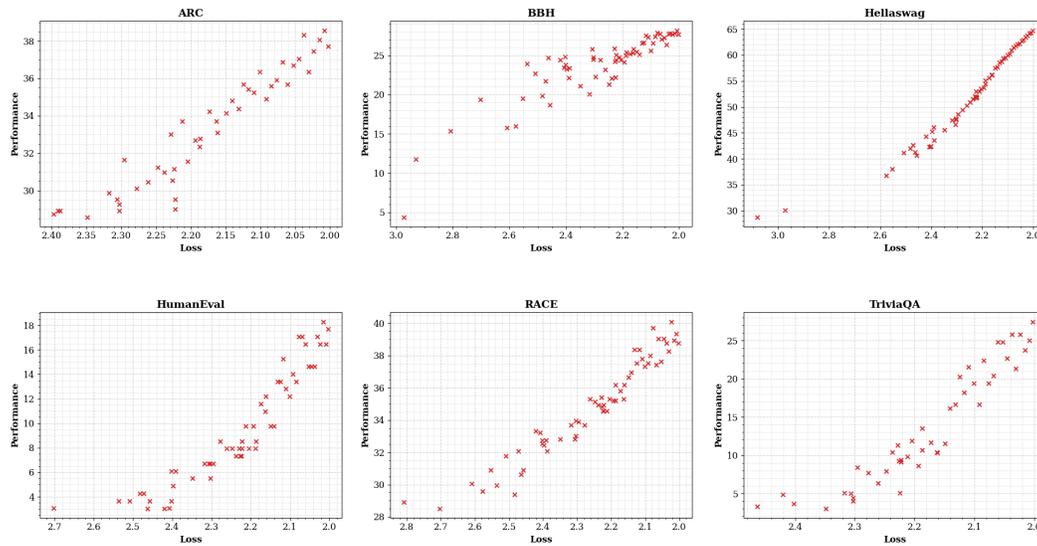


Figure 9: We visualize the relation between pre-training loss and task performance for all LMs that surpass random baseline performance on the target benchmark, observing a generally linear trend.

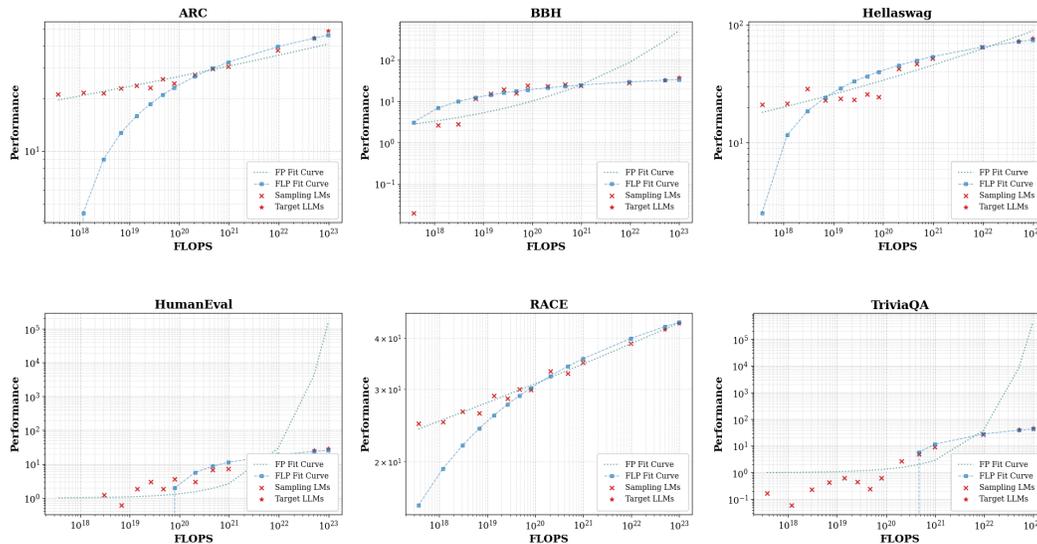


Figure 10: The downstream performance prediction using FP (Achiam et al., 2023) and FLP fit curves. FLP can better predict the downstream performance of target 7B and 13B LLMs across all evaluation benchmarks, while FP’s predictions are very unstable (e.g., HumanEval, TriviaQA).

## APPENDIX

### A LINEAR RELATION BETWEEN LOSS AND PERFORMANCE

We gather data points from intermediate checkpoints of all sampling LMs and visualize the relationship between pre-training loss and corresponding task performance in Fig. 9. We observe a generally linear trend across all benchmarks, which motivates our selection of linear analytical form to characterize the mapping from pre-training loss to downstream performance.

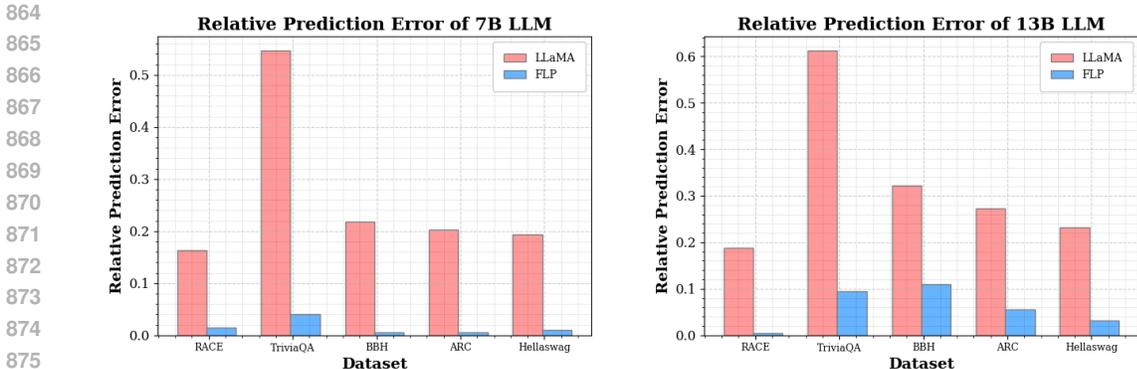


Figure 11: The comparison to the downstream task prediction approach in Llama-3 development (Dubey et al., 2024). We find that initially estimating the negative log-likelihood of the target answer does not effectively predict performance based on our data points.

### B ANALYTICAL FORM TO FIT FLOPS-TO-PERFORMANCE CURVE

We also experiment with the analytical form proposed in Achiam et al. (2023) to estimate the FLOPs-to-Performance curve:

$$\log P(C) = \left(\frac{C}{C_M}\right)^{\alpha_M}, \tag{6}$$

where  $C_M$  and  $\alpha_M$  are constant terms to be estimated. The fit curves are shown in Fig. 10. We observe that FLP still consistently outperforms FP across all evaluation benchmarks. In addition, FP can yield very unstable predictions on certain datasets, like HumanEval and TriviaQA, due to a lack of sufficient data for accurate modeling.

### C COMPARE WITH LLAMA-3 APPROACH

We compare with the Llama-3 approach for downstream task prediction (Dubey et al., 2024). They suggest initially estimating the negative log-likelihood (NLL) of the target answer based on the computational cost in FLOPs, followed by using this NLL to model the task performance through a sigmoid function. The comparison results are shown in Fig. 11. We find that the two-stage approach proposed in Dubey et al. (2024) fails to effectively estimate the performance based on our data points, compared to FLP.

### D MMLU EXPERIMENT

Our sampling LMs, up to 3B, exhibit random performance (*i.e.*, 25%) on the MMLU benchmark (Hendrycks et al., 2020). Consequently, these models do not provide effective data points for estimation. Accordingly, we utilize intermediate checkpoints from 7B LLMs to estimate the performance of 13B LLMs on MMLU using FLP. The results are shown in Fig. 12, and the relative prediction error is 3.54%. FLP can also effectively predict the performance on MMLU by leveraging intermediate LMs checkpoints that emerge on this task.

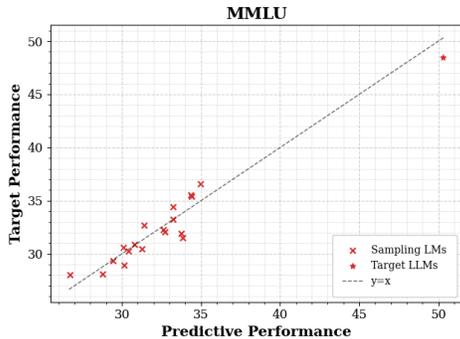


Figure 12: The performance prediction on MMLU using FLP.

### E FLP-M: FIT CURVE FOR ABLATION STUDY

The FLOPs-to-Loss fit curves are in Fig. 13 and the Loss-to-Performance fit curves are in Fig. 14. We observe that  $M_4$  in Tab. 3 offers more stable and accurate predictions for domain-specific loss,

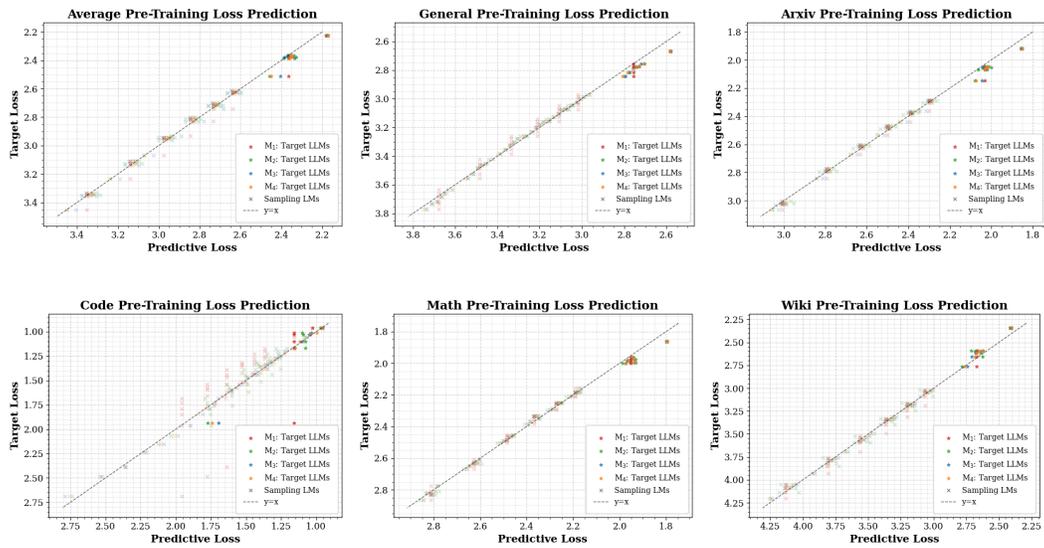


Figure 13: The pre-training loss prediction using various analytical forms.  $M_4$  provides more stable and overall more accurate predictions for domain-specific loss.

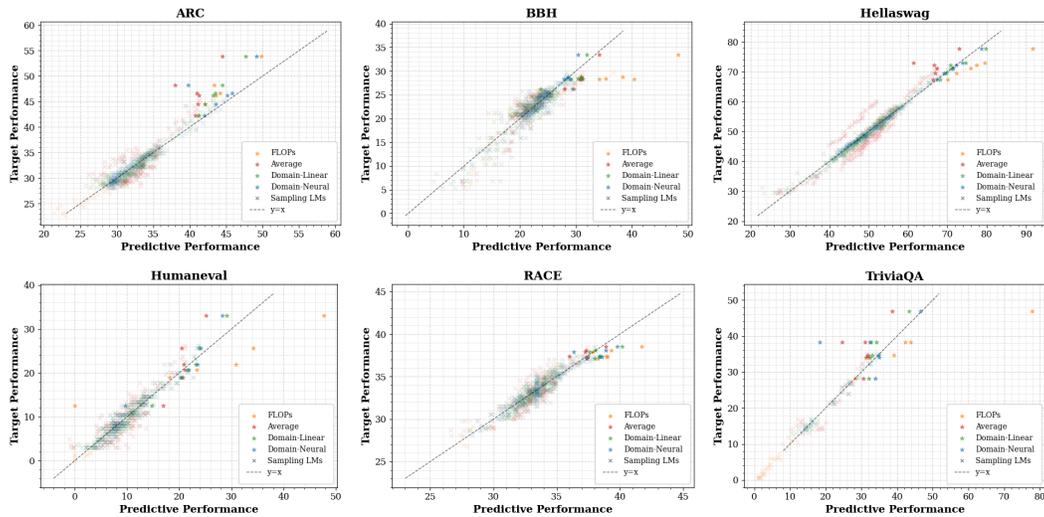


Figure 14: The downstream performance prediction using various approaches. The domain loss coupled with neural network estimation demonstrates the best prediction performance.

with the combined approach of domain loss and neural network estimation delivering the best overall downstream performance prediction.

## F USING DOMAIN LOSS IN FLP

We explore the application of  $FLP-M$  during pre-training on a consistent distribution (the experimental setting described in §4), and compare it with  $FLP$ . The fitting curves are shown in Fig. 15 and the results of relative prediction error are shown in Fig. 16. We show that  $FLP-M$  fails to effectively predict the performance of target LLMs when sampling LLMs are pre-trained on a fixed distribution. This ineffectiveness is attributed to the closely related domain-specific validation losses among the sampling LLMs within the same training distribution, which suggests that decomposing the pre-training validation loss yields no additional information in this pre-training setting. Thus, estimating five domain-

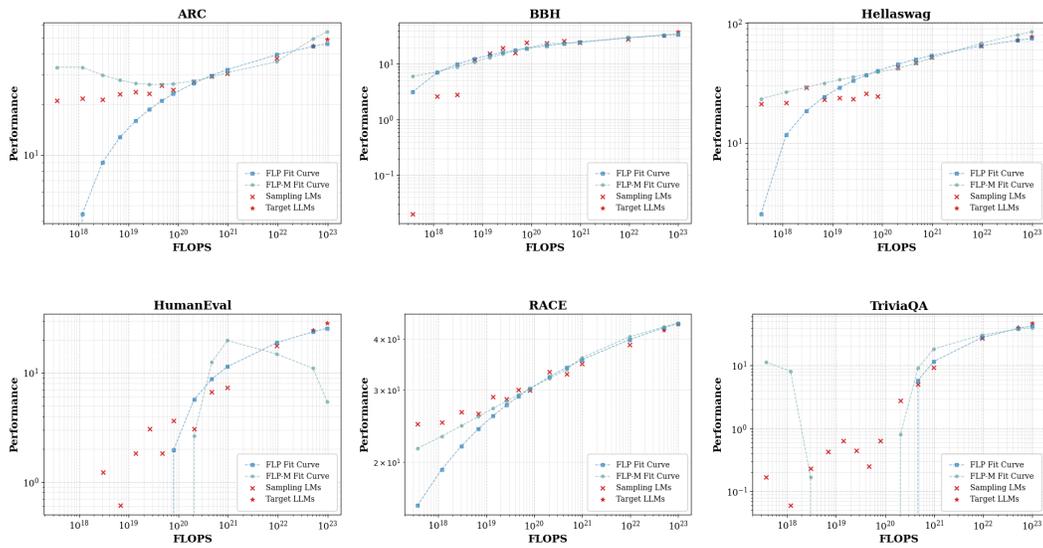


Figure 15: The downstream performance prediction using FLP and FLP-M fit curves. FLP can better predict the downstream performance of target LLMs with 7B and 13B parameters.

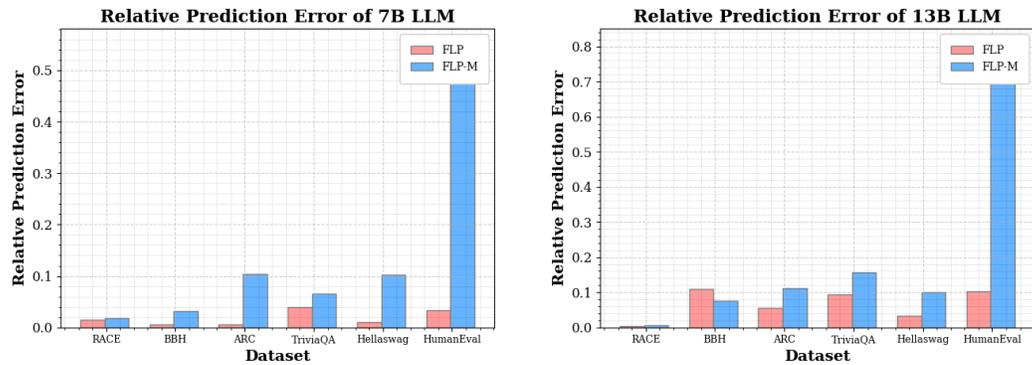


Figure 16: The relative prediction error of 7B and 13B LLMs using FLP and FLP-M. FLP achieves significantly better performance.

specific loss, rather than a single average validation loss, can further increase the risk of error propagation. Moreover, using highly correlated features as neural network inputs may lead to overfitting.