

# ADSO: ADAPTIVE DATA MIXTURE & SCALE OPTIMIZATION

## A MULTI-SCALE MULTI-FIDELITY BAYESIAN OPTIMIZATION APPROACH

**Anonymous authors**

Paper under double-blind review

### ABSTRACT

LLM pre-training requires careful curation of data sources, a process that currently relies heavily on intuition or costly trial-and-error. Since existing ad hoc approaches are unlikely to transfer across domains or data types, we present a unifying framework for data mixture optimization where (mixtures, model scale, training steps) are chosen to balance cost and potential information gain. Going beyond the canonical deterministic extrapolation in scaling laws, we present a sequential decision-making framework where uncertainty in outcomes is explicitly modeled and sharpened as more measurements are gathered. In particular, we formulate a multi-scale, multi-fidelity Bayesian Optimization (BO) problem where information from smaller-scale experiments can systematically inform larger-scale training decisions. We design an adaptive algorithm that takes into account different measurement fidelities provided by model scale and training steps and empirically demonstrate it on a predictor built on 472 pre-training runs with varying data compositions. Compared to standard BO baselines, instantiating our approach with even simple kernels and acquisition functions can allow principled decisions across training models from 20M to 1B parameters and achieve **2.7x** and **6x** speedups compared to multi-fidelity BO and random search baselines in finding the best data mixture for downstream performance under fixed compute budgets. In sum, our adaptive framework underscores potential efficiency gains achievable by developing principled and transferrable data mixture optimization methods. Our code is publicly available at <https://github.com/anonWAEWA/ADSO>.

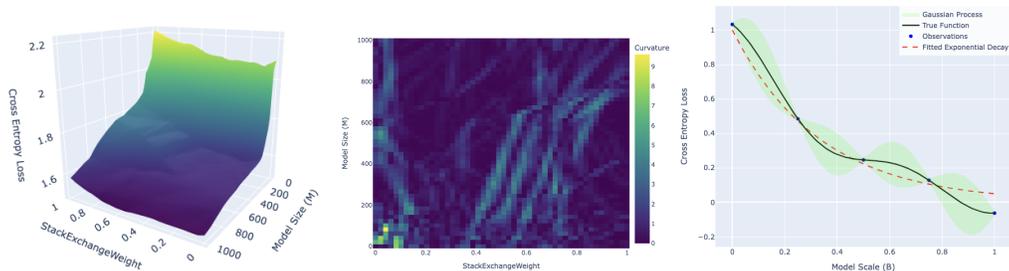


Figure 1: **Left:** The predicted loss as a function of data mixing coefficient and model sizes from a data-driven predictor on 472 runs with high  $R^2$ . Notice the highly non-smooth geometry. **Middle:** The curvature (2nd Derivative) at these points shows there are points of high irregularities. **Right:** A demonstration showing how fitted functional forms like exponential decay would demonstrate a high predictive error if fitted on smaller model points. In contrast, a Gaussian Process would capture uncertainty over the points.

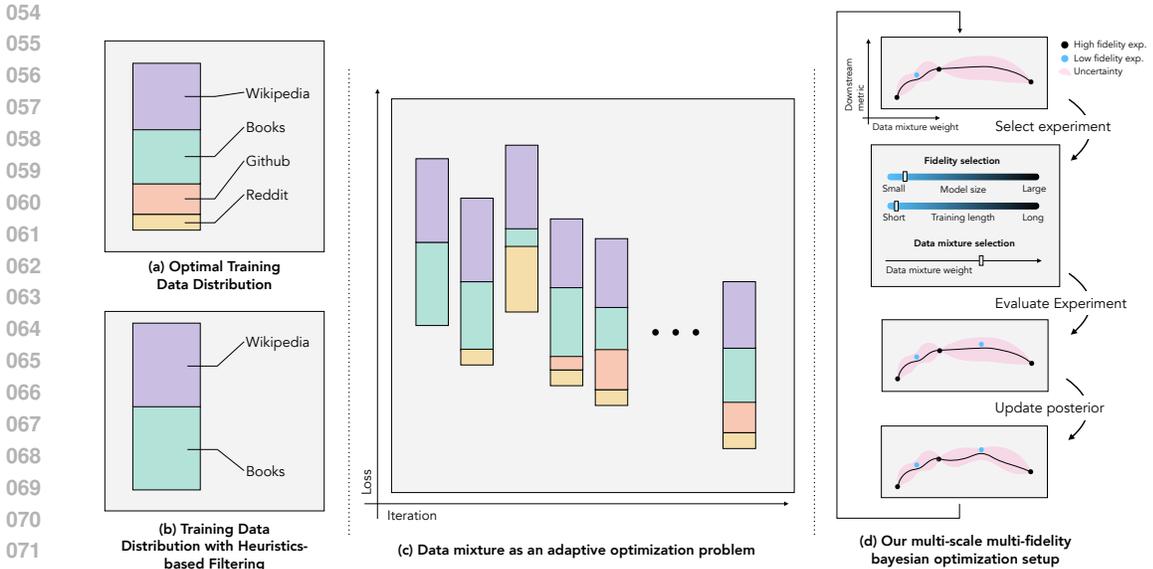


Figure 2: Our multi-scale multi-fidelity Bayesian optimization framework. (a) Given an unknown optimal training data distribution that we have to find, (b) present methods use heuristics-based filtering that guesses and checks. (c) Our algorithm treats data mixture optimization as a Bayesian Optimization problem. (d) Under cost constraint, we explore data mixtures while being cost-aware, where our costs are determined by the fidelity (model size and training steps) that we evaluate. The evaluated result then updates our probabilistic belief over the data mixture, model size, and training steps, which guides subsequent parameters.

## 1 INTRODUCTION

Data is the foundational infrastructure that all AI systems build on. Scaling data has been a key driver of progress in machine learning, particularly in language model training (Deng et al., 2009; Hoffmann et al., 2022a; Gadre et al., 2024). While this data-centric approach has yielded impressive performance gains, it incurs substantial computational and financial costs in training state-of-the-art language models (Hoffmann et al., 2022a; Luccioni et al., 2023). Beyond raw scale, the composition of training data has emerged as a critical factor: when working with heterogeneous data sources, the choice of training mixture has been shown to significantly impact model performance (Albalak et al., 2023a; Goyal et al., 2024a). This recognition has motivated extensive effort in optimizing data mixing strategies. Some institutions have developed proprietary data mixtures based on domain expertise and empirical observations (Radford et al., 2021; Jiang et al., 2023; OpenAI, 2024), while others have proposed systematic heuristics ranging from Wikipedia upsampling to perplexity-guided data selection (Thrush et al., 2024; Blakeney et al., 2024). However, these approaches are unlikely to transfer across domains and data types. For instance, when organizations in specialized sectors such as healthcare or finance seek to train custom language models on proprietary datasets, it remains unclear whether heuristics developed for public datasets are still effective. Given the substantial resources required for training high-performance language models, there is a pressing need for a principled framework to address data mixture optimization.

Recent works have proposed frameworks that attempt to model the relationship between data mixing coefficients and model performance (Ye et al., 2024; Ge et al., 2025a). However, these approaches make strong assumptions that warrant careful examination. The functional relationship between mixing coefficients and performance is likely context-dependent, varying across different data settings and objectives. Moreover, the assumption in these frameworks that the functional relationship is independent of model scales remains untested across different model sizes. In our empirical study, we trained a data-driven predictor on results from 472 random large language model pretraining runs at various scales, which suggests non-trivial relationships in the performance landscape (see Figure 1).

108 Instead, we propose viewing the problem of curating the optimal data mixture as an adaptive opti-  
109 mization problem where practitioners iteratively refine their mixing decisions based on empirical  
110 observations from previous experiments. This framework leverages the intuition that model perfor-  
111 mance exhibits local consistency across similar mixtures and model sizes/scales, while avoiding rigid  
112 assumptions about the global structure of the performance landscape.

113 Sequential optimization of data mixtures necessitates comprehending which data compositions suffer  
114 the highest uncertainty and sharpening beliefs on performance as more observations are gathered.  
115 In particular, good adaptive policies must distinguish between aleatoric and epistemic uncertainty:  
116 epistemic uncertainty can be reduced with more data, while aleatoric is irreducible. Measurements  
117 must be planned to maximally reduce epistemic uncertainty on future runs by balancing exploration  
118 and exploitation.

119 We formulate this sequential optimization framework as a Bayesian optimization problem: we  
120 maintain probabilistic beliefs on the performance of various data mixtures and model scales, and we  
121 use these beliefs to choose the next model scale to train, on what data mixture, and for how long.  
122 Once we fit and evaluate this new model, we use its performance to update our beliefs (Hutter et al.,  
123 2011; Falkner et al., 2018; Frazier, 2018).

124 In traditional BO, the cost of each new observation is the same, and we aim to optimize an objective  
125 while observing the smallest number of points possible. Our setting is more complicated – the cost of  
126 training a new model and observing its performance is affected by (1) the number of steps for which  
127 the model is trained and (2) the scale of the model (the number of parameters therein).

128 The number of steps for which a model is trained affects the *quality* of the observation – the more  
129 steps we use to train the model, the more accurately the results will reflect the utility of training on the  
130 data mixture in question. Previous work has handled this conundrum using so-called multi-fidelity  
131 Bayesian optimization, in which evaluations are ‘stopped early’ during the training process if it  
132 becomes clear the information revealed during additional training steps will not be worth the expense  
133 (Swersky et al., 2014; Domhan et al., 2015; Kandasamy et al., 2017; Li et al., 2018a).

134 Our setting is distinguished by the second factor above – we also want to use data gathered on smaller  
135 model scales to guide our search over parameters for larger models. Importantly, varying model  
136 scale differs fundamentally from traditional fidelity dimensions like training steps. When training for  
137  $z$  steps, we naturally obtain observations for all intermediate steps up to  $z$ . In contrast, evaluating  
138 a model of size  $m$  provides no inherent information about the performance of smaller or larger  
139 architectures. This raises interesting questions about how to appropriately treat and exploit the this  
140 structure, opening new methodological directions for investigations.

141 Luckily, in contrast to conventional hyperparameters like learning rate or momentum, where optimal  
142 configurations exhibit complex scaling behavior across model sizes (Yang et al., 2022), recent  
143 empirical evidence suggests that optimal data mixture compositions enjoy greater transferability from  
144 smaller to larger model architectures (Ye et al., 2024; Ge et al., 2025a). This transferability property  
145 enables the strategic use of smaller-scale evaluations to identify optimal data mixture configurations  
146 that remain effective at target model scales, substantially reducing the computational cost of the  
147 optimization process.

148 The main contributions of the paper are as follows:

- 149
- 150 • *We propose Multi-Fidelity Multi-Scale Bayesian Optimization settings, combining the works*  
151 *in optimizing hyperparameters and scaling laws under one intellectual framework. Our*  
152 *Bayesian Optimization approach better explores different data mixtures and model scales to*  
153 *deliver the best terminal model 2.7x faster in achieving optimal downstream task perfor-*  
154 *mance. (Section 4)*
- 155 • *We show how smaller model sizes affect the predictive utility of larger runs (e.g. how much*  
156 *does training runs below 500M help predict the losses on 1B) by ablating predictors over*  
157 *training runs of different model scales. (Section 3.3)*
- 158 • *We show how earlier training steps improve the predictive utility of others (e.g. given*  
159 *constant FLOPs for hyperparameter search, we show it’s better to have 5 full train runs*  
160 *and 10 half train runs vs 10 full train runs) by ablating predictors of different train steps.*  
161 *(Section 3.4)*

## 2 PROBLEM FORMULATION: MULTI-SCALE MULTI-FIDELITY OPTIMIZATION

We have access to a set of  $n$  datasets  $\mathcal{D} = \{d_1, d_2, \dots, d_n\}$ , and wish to train a model comprising  $m^*$  parameters for  $z^*$  training steps using  $T$  datapoints. We study the problem of finding the optimal fraction of our data budget  $T$  to draw from each of our  $n$  datasets. Let  $w_i T$  denote the number of points we draw from dataset  $i$ , with  $\mathbf{w} = \{w_1, w_2, \dots, w_n\} \in \Delta^n$ , where  $\Delta^n$  is the  $n$ -dimensional probability simplex.

Let  $\mu(\mathbf{w}, m, z)$  denote the performance of a model comprising  $m$  parameters trained for  $z$  training steps with dataset proportions  $\mathbf{w}$  on some downstream task of interest. We seek to solve the following optimization problem

$$\arg \max_{\mathbf{w}} \mu(\mathbf{w}, m^*, z^*) \tag{1}$$

We have a budget  $B$  with which we can experiment with different values of  $\mathbf{w}$ ,  $m$ , and  $z$ . Each evaluation of  $\mu(\mathbf{w}, m, z)$  incurs a cost  $c(m, z)$ .

Using  $m^*$  parameters and  $z^*$  steps every time we evaluate a new set of weights would quickly exhaust our budget. Instead, therefore, we might probe a particular set of weights on a smaller model with  $m$  parameters, or with  $z < z^*$  training steps – while the resulting observation  $\mu(\mathbf{w}, m, z)$  would be less informative than  $\mu(\mathbf{w}, m^*, z^*)$ , it would be considerably cheaper and still provide a valuable update to our posterior. This technique is called multi-fidelity Bayesian optimization.

Traditional approaches to fidelity-aware Bayesian optimization primarily address scenarios in which the model architecture  $m$  remains fixed and only the number of training steps  $z$  is varied (Swersky et al., 2014; Domhan et al., 2015; Kandasamy et al., 2017; Li et al., 2018a). We add a layer of complexity by also considering model scale. It is interesting to note that a fundamental distinction between model scale and number of training steps is that in the course of evaluating a model trained for  $z$  training steps, we must also evaluate that model for all steps  $z' < z$ . No such hierarchical relationship exists for evaluations across different model scales. This structural difference suggests promising avenues for novel methodological developments in multi-fidelity optimization theory, though such extensions lie beyond the scope of our present work.

In this paper, we propose a novel way to take advantage of this additional degree of freedom. To reduce the computational burden of evaluating our technique, we test it on a ‘predictor’ comprising a surrogate model trained on 472 pretraining runs across diverse data mixture coefficients and model scales. This predictor can accurately predict training loss trajectories for any given model scale and set of mixture coefficients. We then evaluate various Bayesian optimization methods using this surrogate model, assessing their efficacy in both optimizing the predicted utility functions and identifying optimal data mixtures among the sampled configurations at our target scale of 1B parameters.

## 3 PREDICTORS OVER DATA MIXTURE AND SCALE

In this section, we discuss the training of the predictor we will use to test our optimization methods. This predictor serves two purposes in our paper (1) as we mentioned above, it allows us to reduce the computational burden required to test our optimization method by evaluating its performance on *predicted* losses (2) it allows us to develop high-level intuition about the way runs involving smaller models or fewer training steps can inform larger runs. The training of such a predictor is *not* required to use the optimization technique we develop in this paper but serves as a justification for the overall validity of our multi-scale, multi-fidelity formulation.

### 3.1 PRETRAINING - COLLECTING PREDICTOR’S DATA

We pretrained 472 language models using the OLMo 2 package (OLMo et al., 2024) and data from SlimPajama (Shen et al., 2024), a deduplicated version of RedPajama (Weber et al., 2024). Slimpajama contains seven data categories – *Wikipedia*, *StackExchange*, *Github*, *ArXiv*, *Book*, *CommonCrawl*, and *C4*. We used only data from the first five categories to train the language models while holding out the data from *CommonCrawl* and *C4* to simulate data mixture optimization in

out-of-distribution settings. For each run, we randomly sample the data mixture proportions from a Dirichlet distribution to uniformly sample from the probability simplex and train the models for 196 training steps. Under this setup, we trained models ranging from 20M to 1B parameters. Additional details on the pretraining setup are provided in Appendix A.

### 3.2 PREDICTOR TRAINING

Our predictor is a multilayer perceptron (MLP) comprising 5,000 parameters that for a given model, predicts its cross-entropy loss over eleven different datasets: the model’s training data, each category in RedPajama listed above, and two additional datasets: *CommonCrawl* and *C4*. In addition, the model is evaluated on three downstream tasks: *hellaswag*, *piqa*, *arc\_easy*.

The model takes the following inputs (1) the model size (2) the number of steps the model is trained for (3) the proportion of each of the five dataset categories mentioned above used in training. All runs are trained with the same number of tokens. The model is trained to minimize the  $R^2$  between predicted and true values.

We now turn to various insights that can be obtained from this predictor.

### 3.3 SMALL MODELS HELP PREDICT LARGER MODELS OUTCOMES

	Train	Test
$E_1$	half of 1B runs	remaining 1B runs
$E_2$	half of 1B runs	remaining 1B runs + 700M runs
$E_3$	half of 1B runs	remaining 1B runs + all smaller runs
$E_4$	half of 700M runs	remaining 700M runs
$E_5$	half of 700M runs	remaining 700M runs + 500 runs

Table 1: Model size experiments

Dataset	E1	E2	E3	E4	E5
wikipedia	0.75	<b>0.96</b>	0.94	0.73	<b>0.88</b>
arxiv	0.68	<b>0.92</b>	0.93	0.59	<b>0.82</b>
github	0.66	0.95	<b>0.95</b>	0.62	<b>0.87</b>
book	0.83	<b>0.97</b>	<b>0.97</b>	0.79	<b>0.92</b>
stackexchange	0.73	<b>0.95</b>	<b>0.95</b>	0.68	<b>0.90</b>
commoncrawl	0.84	<b>0.98</b>	<b>0.98</b>	0.81	<b>0.94</b>
c4	0.86	<b>0.99</b>	0.98	0.82	<b>0.95</b>
arceasy	0.92	<b>0.94</b>	<b>0.94</b>	0.88	<b>0.90</b>
hellaswag	0.97	<b>0.98</b>	0.97	0.94	<b>0.96</b>
piqa	0.94	<b>0.96</b>	<b>0.96</b>	0.90	<b>0.93</b>

Table 2: Results of the experiments listed in table 1, averaged over 3 random seeds. Notice that  $E_2 > E_1$  and  $E_5 > E_4$  – our ability to predict the performance of larger models is considerably enhanced by insights from smaller models. Note also that  $E_3 \approx E_2$ ; adding information about *much* smaller models does not seem to help.

We begin by investigating the extent to which smaller model runs can inform the dynamics of larger ones. Table 1 details these experiments, and Table 2 lists the results of the experiment.

We note that – as expected – information garnered from training runs on *smaller* models seems to considerably increase the accuracy of our predictions on *larger* models, motivating our hope that a carefully crafted optimization algorithm can exploit the relationship.

Unsurprisingly, we note that the closer in scale the smaller models are to the larger model about which we wish to make a prediction, the more useful the information is. We, therefore, expect our optimization algorithm to ‘step through’ model scales, starting with small and cheap models to

270 identify promising data mixtures, and then progressing to larger and larger models, all the while  
 271 refining the data mixtures it considers optimal.

### 273 3.4 EARLIER TRAINING STEPS HELP PREDICT LATER TRAINING STEPS

274  
 275 The second central premise of our approach is that given a fixed computer budget, it is better to  
 276 attempt many runs for fewer training steps than fewer runs for a larger number of training steps.

277 To test this hypothesis, we carry out three additional experiments. In each of these experiments,  
 278 we attempt to predict the final losses in 30% of our model runs (evenly distributed across model  
 279 sizes). The MLP for each of these experiments is trained on (1) a set of complete runs, one for each  
 280 model size (2) a set of ‘truncated’ runs, evenly distributed across model sizes. In  $E_6$ , we use 16 runs  
 281 truncated at 196 training steps, in  $E_7$ , we use 22 runs truncated at 130 training steps, and in  $E_8$ , we  
 282 use 32 runs truncated at 85 steps; thus, these experiments are trained on numbers generated with the  
 283 *same* FLOPS budget.

Dataset	$E_6$	$E_7$	$E_8$
$R^2$	0.69	0.77	<b>0.82</b>
$R^2(\log)$	0.74	0.82	<b>0.85</b>

284  
 285  
 286  
 287  
 288  
 289 Table 3: Notice the predictive power of our MLP is strongest when it is trained on many runs for  
 290 fewer steps. Results averaged over 3 runs. This validates that given a fixed compute budget, it is  
 291 better to have more runs with fewer training steps than fewer runs for a large number of training steps.

## 294 4 BAYESIAN OPTIMIZATION METHODS

295  
 296 In this section, we describe the Bayesian optimization methods we employ to solve the data mixture  
 297 problem and evaluate the efficacy of our proposed framework.

### 299 4.1 MULTI-FIDELITY MULTI-SCALE GAUSSIAN PROCESS (MFMS-GP)

---

#### 301 **Algorithm 1** Multi-fidelity Multi-scale Gaussian Process (MFMS-GP)

---

302 **Require:** Probability space  $\Delta^n$ , model-scale space  $\mathcal{M}$ , training-step space  $\mathcal{Z}$ , and cost function  
 303  $c(\cdot, \cdot)$

- 304 1: Initialize Gaussian Process (GP) surrogate model with three RBF kernels over  $\Delta^n$ ,  $\mathcal{M}$ , and  $\mathcal{Z}$   
 305 and a linear mean function
  - 306 2: Randomly sample points from  $\Delta^n$ ,  $\mathcal{M}$ , and  $\mathcal{Z}$  to initialize hyperparameters of GP.
  - 307 3: Initialize history  $\mathcal{H}$  with the randomly sampled points
  - 308 4: **for** each optimization iteration **do**
  - 309 5:   **for** each  $(m, z) \in \mathcal{M} \times \mathcal{Z}$  **do**
  - 310 6:     Optimize EI within  $(m, z)$  using gradient descent
  - 311 7:   **end for**
  - 312 8:   Select next configuration  $\lambda_{\text{next}} = (\mathbf{w}_{\text{next}}, m_{\text{next}}, z_{\text{next}})$  using Expected Improvement per Unit  
 313 (EI<sub>pu</sub>):  

$$\text{EI}_{\text{pu}}(\lambda) = \frac{\text{EI}(\lambda)}{c(m, z)} \quad \triangleright \text{EI per unit cost}$$
  - 314 9:   Evaluate  $\mu(\lambda_{\text{next}})$
  - 315 10:   Store results in  $\mathcal{H}$
  - 316 11:   Update posterior of GP with  $\mathcal{H}$
  - 317 12: **end for**
  - 318 13: **return** best configuration  $\lambda^* = \arg \max_{\lambda \in \mathcal{H}} \mu(\lambda)$
- 

320  
 321 We implement a Gaussian Process (GP) surrogate model for our multi-fidelity multi-scale setting.  
 322 The kernel of the GP is a product of three separate RBF kernels for the data proportion, the model  
 323 scale, and the training step dimensions. To enable learning the positive correlation between model  
 performance and both model scales and training steps, we use a linear mean function.

For the acquisition function, we use Expected Improvement (EI). EI aims to quantify the expected gain over the current best-observed function value,  $EI(\mathbf{x}) := \mathbb{E}[\max(y^* - f(\mathbf{x}), 0)]$ , where the expectation is taken over the posterior distribution predicted by the surrogate models, and  $y^*$  represents the current best-observed function value, given by  $y^* := f(\mathbf{x}_{\min})$  (Frazier, 2018). The EI function quantifies the expected improvement in the objective value compared to the current best, thereby encouraging the selection of points that are likely to yield better performance.

Equipped with EI, the usual Bayesian optimization approach proceeds by optimizing EI over the parameter space to find the most promising point to evaluate, using gradient-based methods such as L-BFGS-B (Zhu et al., 1997). However, motivated by the fact that the parameter space is discrete over parameter counts ( $m$ ) and training steps ( $z$ ), we optimize EI over each unique tuple  $(m, z)$ . Then, to account for the fact that evaluation for each tuple incurs varying costs ( $c(m, z)$ ), we chose to evaluate the point that has the greatest EI per unit cost (EIpu) (Lee et al., 2020).

## 4.2 BASELINES: MULTI-FIDELITY BAYESIAN OPTIMIZATION

---

### Algorithm 2 Hyperband with Random Forest, EI

---

**Require:** Probability space  $\Delta^n$ , training-step space  $\mathcal{Z}$ , target model scale  $m^*$ , and cost function  $c(m^*, \cdot)$

- 1: Initialize random forest surrogate model RF
- 2: Set initial design with Random Sampling
- 3: Initialize history  $\mathcal{H} = \emptyset$
- 4: **for** each Hyperband iteration **do**
- 5:     Split the total computation budget into  $s$  brackets
- 6:     **for** each bracket  $s_i$  **do**
- 7:         Generate initial configurations  $\mathbf{w}_1, \dots, \mathbf{w}_n$  at lowest fidelity  $z = 1$
- 8:         **for** each fidelity  $z$  from 1 to  $z^*$  **do**
- 9:             Evaluate configurations  $\mathbf{w}_i$  at fidelity  $z$
- 10:            Store results in  $\mathcal{H}$
- 11:            Fit RF on  $\mathcal{H}$
- 12:            Select next  $\lambda_{\text{next}}$  using Expected Improvement
- 13:            Update  $\mathcal{H}$  with new evaluations
- 14:         **end for**
- 15:     **end for**
- 16: **end for**
- 17: **return** best configuration  $\lambda^* = \arg \max_{\lambda \in \mathcal{H}} \mu(\lambda, m^*, z^*)$

---

As a baseline for multi-fidelity Bayesian optimization, we use Hyperband implemented by SMAC: Sequential Model-Based Optimization for General Algorithm Configuration (Lindauer et al., 2022). This multi-fidelity Bayesian optimization uses a random forest as a surrogate model, expected improvement as the acquisition function, and uses Hyperband (Li et al., 2018b), which is an early stopping technique that focuses on efficiently evaluating multiple parameter configurations by progressively eliminating poorly performing candidates, and exploring many combinations with fewer resources. Since the multi-fidelity framework does not offer a straightforward way to incorporate the additional dimension of model scale, throughout the optimization, we fix the number of model’s parameters to the target model scale  $m^*$ .

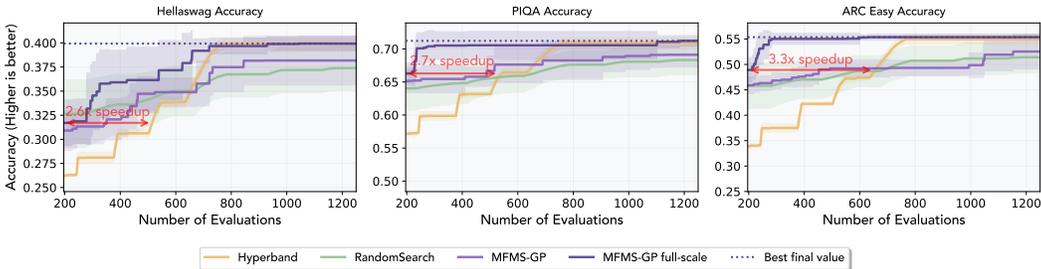
## 4.3 BASELINES: RANDOM AND GRID SEARCH

Random and Grid Search selects hyperparameters that are uniformly drawn from our data proportion space. We then run it against the largest model size and training steps.

## 5 RESULTS

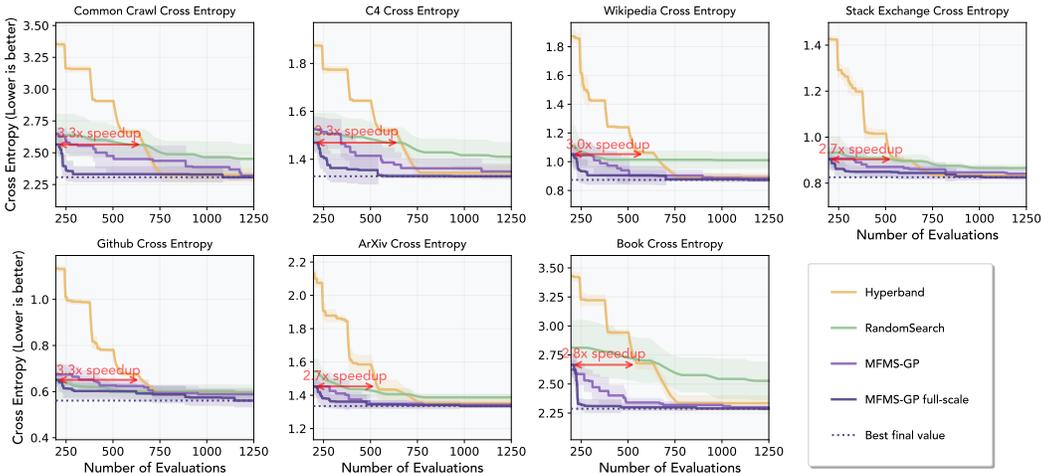
To initiate the hyperparameters of MFMS-GP, we randomly select 20 configurations up to training step  $z = 9$  to fit the kernel and mean functions’ parameters using the Adam optimizer (Kingma and Ba, 2014). The cost of evaluating these configurations are accounted for. Additionally, since it is

378  
379  
380  
381  
382  
383  
384  
385  
386  
387



388 Figure 3: On maximizing accuracy in the downstream tasks, our multi-scale multi-fidelity approach  
389 achieves more than 2.7x speedup and finds the best configuration the fastest.

391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406



407 Figure 4: On minimizing the validation cross-entropy losses, our multi-scale multi-fidelity approach  
408 achieves more than 2.7x speedup and finds the best configuration the fastest.

411 prohibitively expensive to optimize EI for each of 196 training steps, for multi-fidelity multi-scale GP,  
412 we limit the space of training steps to be  $\mathcal{Z} = \{60, 120, 197\}$ . Additional details of the experiments  
413 are available in Appendix B

414 All experiments are run over 5 seeds, and the plots show a 1 standard deviation bound. The number of  
415 evaluations (the x-axis) are in terms of training FLOPS needed to train one 1B model at 100 training  
416 steps. As an example, for Random Search, the 1200 evaluation budget would allow sampling 6 full  
417 runs.

418 Since MFMS-GP relies on GP posterior and potentially noisy EI optimizations to select model scales  
419 and training steps, it may take a while to sample points at the target scale and fidelity. Therefore, we  
420 add an additional plot, MFMS-GP full-scale, that shows the performance one would have gotten if  
421 one takes the best configuration MFMS-GP has observed, and simply set the model scale and training  
422 steps to the target  $m^*$  and  $z^*$ .

423 In Figures 3 & 4, we see that both of the plots for our MFMS-GP algorithm have a 2.6 to 3.3x speedup  
424 in finding the configuration that achieves the highest accuracy. The advantage of the MFMS-GP  
425 method speaks to the tremendous potential of considering our framework for large scale language  
426 model training.

428  
429 **6 RELATED WORKS**

430 **Data Mixtures** Several approaches aim to move beyond heuristic methods for data mixture by  
431 leveraging algorithmic techniques. Albalak et al. (2023b) propose an online data mixing strategy

432 using a non-stochastic bandit algorithm to dynamically adjust data proportions during training,  
 433 maximizing perplexity. DoReMi Xie et al. (2023) focuses on identifying and emphasizing the  
 434 "hardest" datasets for a base model through distributionally robust language modeling to improve  
 435 training efficiency. Ge et al. (2025b) models a joint scaling behavior of domain proportions and  
 436 training steps, we push this further through modeling the model scale. Goyal et al. (2024b) delve  
 437 into the quality-quantity tradeoff in data, exploring how data filtering and repetition affect model  
 438 performance and introducing scaling laws that account for data utility decay. These works highlight  
 439 the increasing interest in principled and adaptive methods for data mixture optimization, yet often  
 440 focus on fixed model scales, contrasting with our multi-scale approach.

441 **Scaling Laws** Scaling laws provide crucial insights into the relationship between model size, training  
 442 compute, and performance in large language models (Kaplan et al., 2020). Hoffmann et al. (2022b)  
 443 established foundational scaling laws demonstrating predictable performance improvements with  
 444 increased compute, model parameters, and training data. Muennighoff et al. (2023) investigate  
 445 the impact of data repetition in data-constrained scenarios, showing diminishing returns beyond  
 446 a certain repetition threshold. Ruan et al. (2024) propose observational scaling laws based on  
 447 "principal capabilities" to explain and predict language model performance across diverse models  
 448 and benchmarks. These scaling law studies inform our framework by providing a theoretical basis for  
 449 understanding performance variations across model scales and data mixtures, allowing us to integrate  
 450 these insights into a multi-fidelity multi-scale Bayesian optimization approach.

451 **Bayesian Optimization** Data mixture optimization, like hyperparameter tuning, benefits from ef-  
 452 ficient search strategies. Approaches range from full configuration selection with methods like  
 453 Bayesian Optimization (BO) to configuration evaluation which employs early termination of un-  
 454 promising runs. Early BO methods Hutter et al. (2011) used Gaussian Processes (GPs) to model the  
 455 relationship between hyperparameters and model performance, subsequent works explored random  
 456 forests (Lindauer et al., 2022) and Parzen estimators (Bergstra et al., 2011) as surrogate models.

457 Early stopping techniques like Hyperband (Li et al., 2018b) focus on efficiently evaluating multiple  
 458 parameter configurations by progressively eliminating poorly performing candidates, and exploring  
 459 many combinations with fewer resources. More recent methods like BOHB (Falkner et al., 2018)  
 460 combine these ideas, leveraging the BO exploration of Parzen estimators with the multi-fidelity  
 461 benefits of Hyperband. Our work, ADSO, is the first to explore a multi-scale multi-fidelity approach  
 462 for data mixture optimization.

## 464 7 CONCLUSION AND FUTURE WORK

466 This work introduces a principled framework, multi-fidelity multi-scale Bayesian optimization, for  
 467 optimizing data mixture compositions in large language model training, a critical challenge in modern  
 468 AI system development. Our framework unifies recent advances in predicting optimal data mixtures  
 469 across scales with classical multi-fidelity Bayesian optimization techniques. Based on this unified  
 470 framework, we implemented the Gaussian process using the RBF kernels and expected-improvement-  
 471 per-unit acquisition function to balance the information gain and the cost of exploring new points in  
 472 the functional landscape. We find that the method achieves optimal downstream task performance 2.7  
 473 times faster than traditional multi-fidelity approaches by strategically exploring the joint space of  
 474 data mixtures and model scales.

475 In addition, we empirically demonstrate two key insights that inform future efficient optimization  
 476 of data mixtures. First, our analysis reveals that training runs on smaller models (below 500M  
 477 parameters) provide valuable predictive signals for optimizing larger architectures (1B parameters).  
 478 Second, we establish that partial training runs can effectively inform full-scale training decisions.  
 479 Specifically, our results show that a combination of full and partial training runs (e.g. 5 complete and  
 480 10 half-length runs) yields better predictive utility than an equal-compute allocation of full training  
 481 runs alone (e.g. 10 complete runs).

482 Several promising directions emerge for future research. First, extending our framework to more  
 483 settings such as language model fine-tuning, data filtering, and more diverse collections of datasets  
 484 would validate its generalizability across different data mixing scenarios. From a methodological  
 485 perspective, incorporating domain knowledge about the positive correlation between model perfor-  
 mance and both parameter count and training duration could enhance the Gaussian process kernel

design. Additionally, the fundamental differences between model scale and training steps as fidelity dimensions call for deeper methodological investigation in their appropriate treatment in the framework. Finally, exploring alternative acquisition functions, such as knowledge gradient (Poloczek et al., 2016; Wu et al., 2019), could further improve the framework’s efficiency in navigating the optimization landscape.

## IMPACT STATEMENT

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none of which we feel must be specifically highlighted here.

## REFERENCES

- Alon Albalak, Liangming Pan, Colin Raffel, and William Yang Wang. Efficient online data mixing for language model pre-training, 2023a. URL <https://arxiv.org/abs/2312.02406>.
- Alon Albalak, Liangming Pan, Colin Raffel, and William Yang Wang. Efficient online data mixing for language model pre-training, 2023b. URL <https://arxiv.org/abs/2312.02406>.
- James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyper-parameter optimization. In *Proceedings of the 25th International Conference on Neural Information Processing Systems, NIPS’11*, page 2546–2554, Red Hook, NY, USA, 2011. Curran Associates Inc. ISBN 9781618395993.
- Cody Blakeney, Mansheej Paul, Brett W. Larsen, Sean Owen, and Jonathan Frankle. Does your data spark joy? performance gains from domain upsampling at the end of training, 2024. URL <https://arxiv.org/abs/2406.03476>.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- Tobias Domhan, Jost Tobias Springenberg, and Frank Hutter. Speeding up automatic hyperparameter optimization of deep neural networks by extrapolation of learning curves. In *Proceedings of the 24th International Conference on Artificial Intelligence, IJCAI’15*, page 3460–3468. AAAI Press, 2015. ISBN 9781577357384.
- Stefan Falkner, Aaron Klein, and Frank Hutter. BOHB: robust and efficient hyperparameter optimization at scale. *CoRR*, abs/1807.01774, 2018. URL <http://arxiv.org/abs/1807.01774>.
- Peter I. Frazier. A tutorial on bayesian optimization, 2018. URL <https://arxiv.org/abs/1807.02811>.
- Samir Yitzhak Gadre, Georgios Smyrnis, Vaishaal Shankar, Suchin Gururangan, Mitchell Wortsman, Rulin Shao, Jean Mercat, Alex Fang, Jeffrey Li, Sedrick Keh, Rui Xin, Marianna Nezhurina, Igor Vasiljevic, Jenia Jitsev, Luca Soldaini, Alexandros G. Dimakis, Gabriel Ilharco, Pang Wei Koh, Shuran Song, Thomas Kollar, Yair Carmon, Achal Dave, Reinhard Heckel, Niklas Muennighoff, and Ludwig Schmidt. Language models scale reliably with over-training and on downstream tasks, 2024. URL <https://arxiv.org/abs/2403.08540>.
- Ce Ge, Zhijian Ma, Daoyuan Chen, Yaliang Li, and Bolin Ding. Bimix: A bivariate data mixing law for language model pretraining, 2025a. URL <https://arxiv.org/abs/2405.14908>.
- Ce Ge, Zhijian Ma, Daoyuan Chen, Yaliang Li, and Bolin Ding. Bimix: A bivariate data mixing law for language model pretraining, 2025b. URL <https://arxiv.org/abs/2405.14908>.
- Sachin Goyal, Pratyush Maini, Zachary C. Lipton, Aditi Raghunathan, and J. Zico Kolter. Scaling laws for data filtering—data curation cannot be compute agnostic. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 22702–22711, 2024a. doi: 10.1109/CVPR52733.2024.02142.

- 540 Sachin Goyal, Pratyush Maini, Zachary C. Lipton, Aditi Raghunathan, and J. Zico Kolter. Scaling  
541 laws for data filtering – data curation cannot be compute agnostic, 2024b. URL [https://](https://arxiv.org/abs/2404.07177)  
542 [arxiv.org/abs/2404.07177](https://arxiv.org/abs/2404.07177).  
543
- 544 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza  
545 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom  
546 Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy,  
547 Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre.  
548 Training compute-optimal large language models, 2022a. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2203.15556)  
549 [2203.15556](https://arxiv.org/abs/2203.15556).
- 550 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza  
551 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom  
552 Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy,  
553 Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre.  
554 Training compute-optimal large language models, 2022b. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2203.15556)  
555 [2203.15556](https://arxiv.org/abs/2203.15556).
- 556 Frank Hutter, Holger H. Hoos, and Kevin Leyton-Brown. Sequential model-based optimization for  
557 general algorithm configuration. In *Proceedings of the 5th International Conference on Learning*  
558 *and Intelligent Optimization, LION'05*, page 507–523, Berlin, Heidelberg, 2011. Springer-Verlag.  
559 ISBN 9783642255656. doi: 10.1007/978-3-642-25566-3\_40. URL [https://doi.org/10.](https://doi.org/10.1007/978-3-642-25566-3_40)  
560 [1007/978-3-642-25566-3\\_40](https://doi.org/10.1007/978-3-642-25566-3_40).
- 561 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,  
562 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier,  
563 L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas  
564 Wang, Timoth  e Lacroix, and William El Sayed. Mistral 7b, 2023. URL [https://arxiv.](https://arxiv.org/abs/2310.06825)  
565 [org/abs/2310.06825](https://arxiv.org/abs/2310.06825).
- 566 Kirthevasan Kandasamy, Gautam Dasarathy, Jeff Schneider, and Barnab  s P  czos. Multi-fidelity  
567 Bayesian optimisation with continuous approximations. In Doina Precup and Yee Whye Teh,  
568 editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of  
569 *Proceedings of Machine Learning Research*, pages 1799–1808. PMLR, 06–11 Aug 2017. URL  
570 <https://proceedings.mlr.press/v70/kandasamy17a.html>.  
571
- 572 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child,  
573 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models,  
574 2020. URL <https://arxiv.org/abs/2001.08361>.
- 575 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization.  
576 *CoRR*, abs/1412.6980, 2014. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:6628106)  
577 [6628106](https://api.semanticscholar.org/CorpusID:6628106).
- 578 Eric Hans Lee, Valerio Perrone, C  dric Archambeau, and Matthias W. Seeger. Cost-aware bayesian op-  
579 timization. *CoRR*, abs/2003.10870, 2020. URL <https://arxiv.org/abs/2003.10870>.  
580
- 581 Lisha Li, Kevin Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. Hyperband:  
582 A novel bandit-based approach to hyperparameter optimization. *Journal of Machine Learning*  
583 *Research*, 18(185):1–52, 2018a. URL [http://jmlr.org/papers/v18/16-558.html](http://jmlr.org/papers/v18/li18a.html).  
584
- 585 Lisha Li, Kevin Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. Hyperband:  
586 A novel bandit-based approach to hyperparameter optimization, 2018b. URL [https://arxiv.](https://arxiv.org/abs/1603.06560)  
587 [org/abs/1603.06560](https://arxiv.org/abs/1603.06560).
- 588 Marius Lindauer, Katharina Eggenberger, Matthias Feurer, Andr   Biedenkapp, Difan Deng, Carolin  
589 Benjamins, Tim Ruhkopf, Ren   Sass, and Frank Hutter. Smac3: A versatile bayesian optimization  
590 package for hyperparameter optimization. *Journal of Machine Learning Research*, 23(54):1–9,  
591 2022. URL [http://jmlr.org/papers/v23/21-0888.html](http://jmlr.org/papers/v23/lindauer22.html).
- 592 Alexandra Sasha Luccioni, Sylvain Viguier, and Anne-Laure Ligozat. Estimating the carbon footprint  
593 of bloom, a 176b parameter language model. *Journal of Machine Learning Research*, 24(253):  
1–15, 2023. URL [http://jmlr.org/papers/v24/23-0069.html](http://jmlr.org/papers/v24/luccioni23.html).

- 594 Niklas Muennighoff, Alexander M. Rush, Boaz Barak, Teven Le Scao, Aleksandra Piktus, Nouamane  
595 Tazi, Sampo Pyysalo, Thomas Wolf, and Colin Raffel. Scaling data-constrained language models,  
596 2023. URL <https://arxiv.org/abs/2305.16264>.
- 597 Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia,  
598 Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Tafjord, Taira  
599 Anderson, David Atkinson, Faeze Brahman, Christopher Clark, Pradeep Dasigi, Nouha Dziri,  
600 Michal Guerquin, Hamish Ivison, Pang Wei Koh, Jiacheng Liu, Saumya Malik, William Merrill,  
601 Lester James V. Miranda, Jacob Morrison, Tyler Murray, Crystal Nam, Valentina Pyatkin, Aman  
602 Rangapur, Michael Schmitz, Sam Skjonsberg, David Wadden, Christopher Wilhelm, Michael  
603 Wilson, Luke Zettlemoyer, Ali Farhadi, Noah A. Smith, and Hannaneh Hajishirzi. 2 olmo 2 furious,  
604 2024. URL <https://arxiv.org/abs/2501.00656>.
- 605 OpenAI. Gpt-4 technical report, 2024. URL <https://arxiv.org/abs/2303.08774>.
- 606 Matthias Poloczek, Jialei Wang, and Peter I. Frazier. Multi-information source optimization, 2016.  
607 URL <https://arxiv.org/abs/1603.00389>.
- 609 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
610 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever.  
611 Learning transferable visual models from natural language supervision, 2021. URL <https://arxiv.org/abs/2103.00020>.
- 612 Yangjun Ruan, Chris J. Maddison, and Tatsunori Hashimoto. Observational scaling laws and the  
613 predictability of language model performance, 2024. URL <https://arxiv.org/abs/2405.10938>.
- 614 Zhiqiang Shen, Tianhua Tao, Liqun Ma, Willie Neiswanger, Zhengzhong Liu, Hongyi Wang, Bowen  
615 Tan, Joel Hestness, Natalia Vassilieva, Daria Soboleva, and Eric Xing. Slimpajama-dc: Under-  
616 standing data combinations for llm training, 2024. URL <https://arxiv.org/abs/2309.10818>.
- 617 Kevin Swersky, Jasper Snoek, and Ryan Prescott Adams. Freeze-thaw bayesian optimization, 2014.  
618 URL <https://arxiv.org/abs/1406.3896>.
- 619 Tristan Thrush, Christopher Potts, and Tatsunori Hashimoto. Improving pretraining data using  
620 perplexity correlations, 2024. URL <https://arxiv.org/abs/2409.05816>.
- 621 Maurice Weber, Daniel Fu, Quentin Anthony, Yonatan Oren, Shane Adams, Anton Alexandrov,  
622 Xiaozhong Lyu, Huu Nguyen, Xiaozhe Yao, Virginia Adams, Ben Athiwaratkun, Rahul Chalamala,  
623 Kezhen Chen, Max Ryabinin, Tri Dao, Percy Liang, Christopher Ré, Irina Rish, and Ce Zhang.  
624 Redpajama: an open dataset for training large language models, 2024. URL <https://arxiv.org/abs/2411.12372>.
- 625 Jian Wu, Saul Toscano-Palmerin, Peter I. Frazier, and Andrew Gordon Wilson. Practical multi-  
626 fidelity bayesian optimization for hyperparameter tuning. *CoRR*, abs/1903.04703, 2019. URL  
627 <http://arxiv.org/abs/1903.04703>.
- 628 Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy Liang, Quoc V.  
629 Le, Tengyu Ma, and Adams Wei Yu. Doremi: Optimizing data mixtures speeds up language model  
630 pretraining, 2023. URL <https://arxiv.org/abs/2305.10429>.
- 631 Greg Yang, Edward J. Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick  
632 Ryder, Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tensor programs v: Tuning large neural  
633 networks via zero-shot hyperparameter transfer, 2022. URL <https://arxiv.org/abs/2203.03466>.
- 634 Jiasheng Ye, Peiju Liu, Tianxiang Sun, Yunhua Zhou, Jun Zhan, and Xipeng Qiu. Data mixing  
635 laws: Optimizing data mixtures by predicting language modeling performance, 2024. URL  
636 <https://arxiv.org/abs/2403.16952>.
- 637 Ciyou Zhu, Richard H. Byrd, Peihuang Lu, and Jorge Nocedal. Algorithm 778: L-bfgs-b: Fortran  
638 subroutines for large-scale bound-constrained optimization. *ACM Trans. Math. Softw.*, 23(4):  
639 550–560, December 1997. ISSN 0098-3500. doi: 10.1145/279232.279236. URL <https://doi.org/10.1145/279232.279236>.
- 640  
641  
642  
643  
644  
645  
646  
647

## A PRETRAINING RUNS DETAILS

We use the OLMo 2 OLMo et al. (2024) package for training our language models. The model configurations are

Group	d_model	n_heads	n_layers	Runs
20M	256	8	8	115
60M	512	8	8	71
150M	768	12	12	53
300M	1024	16	16	74
500M	1280	16	16	39
700M	1536	16	16	52
1B	2048	16	16	68

Table 4: Model Architecture Details by Group with Number of Runs

The training configurations (learning rate, momentum etc.) are directly taken from OLMo’s configuration files (e.g. 700M) We study the compute optimal regime Hoffmann et al. (2022a): for each 1B model runs, we used 20B tokens in total for training. In the interest of collecting more runs, all other model scales are trained on 10B tokens.

## B BAYESIAN OPTIMIZATION DETAILS

For MFMS-GP, the cost of evaluating a run at a particular model scale is taken from the number of FLOPS the corresponding model scale costs during the pretraining runs. The costs are scaled appropriately such that a unit of cost corresponds FLOPS required to train 1B model for 1 training step.

The GP hyperparameters are trained using the Adam optimizer with 0.1 learning rate for 50 iterations.

To search for optimal EI within each  $(m, z)$  tuple, we initiate 5 random probability weights and perform a gradient search over the probability simplex.

Occasionally, the GP would be too certain of its posterior prediction such that the optimized EI are all small in magnitude. Therefore, when the optimal EI is below a certain threshold, we lower the length scales of the RBF kernels to encourage more exploration. The threshold is set to be  $1e^{-4}$ , and the length scales would be lowered to 95% of their original values.

As a measure to encourage selecting higher cost evaluations later in the optimization cycle, instead of using  $EI_{pu}(\lambda) = \frac{EI(\lambda)}{c(m,z)}$ , we introduce an additional parameter  $\alpha$  that controls the importance of cost, and pick the configuration that maximizes  $\frac{EI(\lambda)}{c(m,z)^\alpha}$ . Initially  $\alpha = 1$ , and it decays by 1% for every step of the Bayesian optimization. You may include other additional sections here.