Dynamic Gradient Alignment For Online Data Mixing

Anonymous authors

003 004

006

008 009

010

011

012

013

014

015

016

017

018

019

021

024

025

026 027 Paper under double-blind review

ABSTRACT

The composition of training data mixtures is critical for effectively training large language models (LLMs), as it directly impacts their performance on downstream tasks. Our goal is to identify an optimal data mixture to specialize an LLM for a specific task with access to only a few examples. Traditional approaches to this problem include ad-hoc reweighting methods, importance sampling, and gradient alignment techniques. This paper focuses on gradient alignment and introduces Dynamic Gradient Alignment (DGA), a scalable online gradient alignment algorithm. DGA dynamically estimates the pre-training data mixture on which the models' gradients align as well as possible with those of the model on the specific task. DGA is the first gradient alignment approach that incurs minimal overhead compared to standard pre-training and outputs a competitive model, eliminating the need for retraining the model. Experimentally, we demonstrate significant improvements over importance sampling in two key scenarios: (i) when the pre-training set is small and importance sampling overfits due to limited data; and (ii) when there is insufficient specialized data, trapping importance sampling on narrow pockets of data. Our findings underscore the effectiveness of gradient alignment methods in optimizing training data mixtures, particularly in data-constrained environments, and offer a practical solution for enhancing LLM performance on specific tasks with limited data availability.

028 1 INTRODUCTION 029

 Large Language Models (LLMs) are typically pre-trained on extensive, generic corpora sourced from a variety of data domains (Brown et al., 2020b; Touvron et al., 2023; Zhang et al., 2022), with the composition of these corpora often depending on domain availability or heuristics (Gao et al., 2020; Together AI Team, 2023). While the diversity of natural texts allows the model to learn from various knowledge sources, not all data domains are equally beneficial according to the targeted tasks. The uncurated nature of web-crawled contents could lead to sub-optimal outcomes due to the variations in data quality (Longpre et al., 2023). Plus, some domains may contain misinformation and biases, as one potential source of hallucinations in language generation (Lin et al., 2022; Huang et al., 2023).

To better generalize to the downstream target tasks, it is critical to identify the most beneficial subset from large, generic pretraining corpora. While sample-level selection can be costly, *domain reweighting* offers an efficient group-level approach. Domain reweighting methods assume that samples from the same domain share similar features and search for optimal sampling weights across *domains* (Xie et al., 2023a; Fan et al., 2024; Liu et al., 2024; Kang et al., 2024; Grangier et al., 2024). The domains that most positively impact the target tasks should be assigned higher weights.

In this work, on top of a large, generic pretraining corpus, we assume we have access to a few examples representative of the downstream task on which we want the model to generalize, a so-called *specialized set*. For this setup, Grangier et al. (2024) recently proposed a simple and scalable *importance sampling* based method to domain reweighting, where the weight of a domain is given by the frequency of samples in the specialized set closest to the domain, where distance is measured with SentenceBert (Reimers & Gurevych, 2019) embeddings. This method determines the domain weights before any training and is model-agnostic.

Likewise, prior gradient-alignment methods determine a *static* domain weights for large-scale LM training, often relying on a small-scale proxy model (Xie et al., 2023a; Fan et al., 2024) or fitting a scaling law (Liu et al., 2024; Kang et al., 2024). While these methods show improvements over training on the natural distribution of a generic corpus, they do not dynamically update domain

weights during training to adapt to the current model state. In practical training scenarios, a large model may quickly overfit on certain domains with high weights. In such cases, an online weighting method can respond by shifting emphasis to other domains.

057 We propose Dynamic Gradient-Alignment (DGA), an online domain reweighting method that 058 estimates step-wise optimal domain weights during model training. Inspired by DOGE (Fan et al., 2024), at each reweighting step, DGA upweights the data domain whose gradient aligns more with 060 the model's gradient on the specific set. From the optimization perspective, training the model on 061 the most-aligned data domain yields the greatest reduction in the targeted loss. By incorporating 062 an exponential-moving-average (EMA) term in online domain weights updates, DGA effectively 063 mitigates overfitting and prioritizes the domains that currently benefit the target task the most. Since 064 the domain weights and model parameters are updated concurrently, inaccurate domain weights can potentially drive the model into suboptimal states, which further leads to snow-balled errors. In 065 such cases, the EMA term serves as a correction factor, guiding the model back to a more stable 066 state. As an additional contribution, we scale the domain reweighting methods into extremely fine-067 grained domains (e.g. 262k domains) by introducing a novel distribution reweighting mechanism. 068 Rather than directly reweighting 262k data domains, distribution reweighting reparameterizes the 069 high-dimensional domain weights as a convex combination of weight vectors derived from a set of distributions estimated from embedding-based importance sampling (Grangier et al., 2024). With 071 the number of distributions less than the number of training data domains, it allows DGA to scale to 072 thousands of domains and make the most of the fine-grained group-level features. 073

Our experiments demonstrate the effectiveness of DGA compared to standard pre-training and importance sampling baselines in two challenging cases: (1) the resource of training tokens in each domain is limited instead of infinite (§ 3.1), and (2) the domain granularity is extremely large, which introduces intractable computation overheads on the domain reweighting problem (§ 3.2).

077 078 079

2 DATA MIXING WITH SPECIALIZED TARGET

080 2.1 GENERIC DATASET AND SPECIFIC TASKS

082 We consider a generic training corpus $D_{gen} = \{D_1, \ldots, D_k\}$, which is partitioned into k distinct 083 data domains. We can sample from each of the k domains to train a model. Consequently, we can sample from a *mixture* of these domains and draw a batch of data following the law $x \sim \min(\alpha) \triangleq$ 084 $\sum_{i=1}^{k} \alpha_i D_i$, where $\alpha \in \mathbb{R}^k$ is the mixture *weights*, belonging to the *simplex* $\alpha \in \Delta^k \triangleq \{\alpha \in \mathcal{A}^k \in \mathcal{A}^k\}$ 085 $\mathbb{R}^k | \sum_{i=1}^k \alpha_i = 1$ and $\alpha_i \ge 0$ for all i}. Here, getting one sample from $\min(\alpha)$ means first getting a 087 random index $i \in \{k\}$ from the categorical distribution corresponding to the vector of probabilities α , and then outputting a random sample from the domain D_i . Sampling from this law is computationally efficient if we can efficiently sample from each domain. Next, we consider a model, parameterized by $\theta \in \mathbb{R}^p$, and a loss function $\ell(\theta, x)$ defined for $x \in D_{gen}$. To simplify notation, given a set of 090 samples S (which can be either a full dataset D_i , or a mini-batch), we denote the average over S of the 091 loss $\ell(\boldsymbol{\theta}, S) \triangleq \frac{1}{\#S} \sum_{\boldsymbol{x} \in S} \ell(\boldsymbol{\theta}, \boldsymbol{x})$. Since we focus on LLMs, ℓ is typically the next-token-prediction 092 loss. For a given mixture weight α , we can update θ by doing optimization steps on the generic loss 093

094 095 096

102

$$L_{\text{gen}}(\boldsymbol{\theta}, \boldsymbol{\alpha}) \triangleq \mathbb{E}_{\boldsymbol{x} \sim \min(\boldsymbol{\alpha})}[\ell(\boldsymbol{\theta}, \boldsymbol{x})] = \sum_{i=1}^{k} \alpha_i L_i(\boldsymbol{\theta}) \text{ with } L_i(\boldsymbol{\theta}) \triangleq \ell(\boldsymbol{\theta}, D_i)$$
(1)

In this paper, our goal is to use this data-mixture to train a model that performs well a *specific* task. We assume to have access to samples from this task, split into train and test sets. We call the train set the *specific dataset* D_{spe} that we use to train models. The performance on the specific set is measured with the *specific loss*

$$L_{\rm spe}(\boldsymbol{\theta}) \triangleq \ell(\boldsymbol{\theta}, D_{\rm spe}).$$
 (2)

We assume that the specific set $D_{\rm spe}$ is small; hence, optimizing $L_{\rm spe}$ directly leads to overfitting: the loss on the test data would be much higher than $L_{\rm spe}$. Instead, to get to a low specific loss, we aim to find the optimal data mixing across k data domains α at each training step to get a good model while training on the reweighted generic distribution mix(α).

107 The target specialization task can be flexible according to the application domains, ranging from reasoning, instruction following, etc., corresponding to various objective loss functions, including the

108 next-token prediction loss and preference-based losses when applied on pair-wise datasets. In this paper, we focus on next-token prediction on another dataset. Next, we introduce a general bilevel 110 formulation of the data mixing problem.

111 2.2 **BILEVEL FORMULATION** 112

116 117

113 Since the lack of data forbids optimizing directly $L_{\rm spe}$, we look for the mixture α such that optimizing 114 the generic loss $L_{\text{gen}}(\theta, \alpha)$ yields the smallest specific loss (Grangier et al., 2023). This is formalized by the following *bilevel optimization* problem (Bracken & McGill, 1973; Dagréou et al., 2022): 115

$$\boldsymbol{\alpha}^{\star} \in \operatorname*{arg\,min}_{\boldsymbol{\alpha} \in \Delta^{k}} L_{\mathrm{spe}}(\boldsymbol{\theta}^{\star}(\boldsymbol{\alpha})), \text{ such that } \boldsymbol{\theta}^{\star}(\boldsymbol{\alpha}) \in \operatorname*{arg\,min}_{\boldsymbol{\theta}} L_{\mathrm{gen}}(\boldsymbol{\theta}, \boldsymbol{\alpha})$$
(3)

118 This bilevel formulation is intuitive: for a given weight α , the parameters obtained by minimizing the 119 generic loss L_{gen} are $\theta^*(\alpha)$, and we want those weights to yield a small specific loss. Notably, if 120 the specific loss is a mixture of generic data with an unknown weight $\tilde{\alpha}$, the bilevel formulation is 121 guaranteed to recover it. In other words: 122

Theorem 1. Assume that there exists $\tilde{\alpha}$ such that $D_{spe} = \min(\tilde{\alpha})$. Then, $\tilde{\alpha}$ is a solution to the 123 bilevel problem in Equation 3. 124

125 *Proof.* We let $\hat{\theta}$ the minimizer of L_{spe} . Then, for all α , we have by definition that $L_{\text{spe}}(\theta^*(\alpha)) \geq 1$ 126 $L_{\rm spe}(\boldsymbol{\theta})$. Furthermore, since $D_{\rm spe} = \min(\tilde{\boldsymbol{\alpha}})$, we have that $L_{\rm gen}(\boldsymbol{\theta}, \tilde{\boldsymbol{\alpha}}) = L_{\rm spe}(\boldsymbol{\theta})$ for all $\boldsymbol{\theta}$, hence 127 minimizing this yields $\theta^*(\alpha) = \hat{\theta}$. Putting these results together, we have proven that for all α , it 128 holds $L_{\rm spe}(\theta^*(\alpha)) \ge L_{\rm spe}(\theta^*(\tilde{\alpha}))$, so that $\tilde{\alpha}$ is a solution to Equation 3. 129

130 We consider two types of methods to solve Equation 3. Static methods construct a single mixture 131 weight vector $\boldsymbol{\alpha}$ and then minimize $L_{\text{gen}}(\boldsymbol{\theta}, \boldsymbol{\alpha})$; we describe in the next section how to obtain this vector α . Online methods modify the weights dynamically during model training. They produce a 132 sequence of weights $\alpha^{(t)}$ where t is the optimization iterate. In that case, at each training step, the 133 parameters $\theta^{(t)}$ are updated by doing an optimization step — with gradient descent or Adam — on 134 the function $L_{\text{gen}}(\theta, \alpha^{(t)})$. We now discuss methods to obtain a weight vector α or a sequence $\alpha^{(t)}$. 135

136 2.3 A STRONG BASELINE: IMPORTANCE SAMPLING 137

A sensible strategy is to train the model on a data mixture that most resembles the composition of the 138 targeted specialization data distribution. This is the philosophy behind importance sampling (Kloek 139 & Van Dijk, 1978). We estimate the importance sampling weights α^{IS} using the method of Grangier 140 et al. (2024). The core idea is to embed each generic domain using SentenceBert (Reimers & 141 Gurevych, 2019), and then compute the centroid of each domain $b_i = \frac{1}{\#D_i} \sum_{x \in D_i} \text{Bert}(x)$. This 142 defines a simple and cheap to compute selection function $c(x) \in \{1 \dots k\}$, assigning x to its closest 143 centroid, i.e., $c(\boldsymbol{x}) = \arg\min_{i} \|\operatorname{Bert}(\boldsymbol{x}) - \boldsymbol{b}_{i}\|$ for $\boldsymbol{x} \in D_{\operatorname{gen}} \cup D_{\operatorname{spec}}$. We use it to predict the 144 closest generic data domain for each data instance from the specific set. The importance sampling 145 weights are obtained as the ratio of examples falling in each bin: $\alpha_i^{\text{IS}} \triangleq \frac{\#\{x \in D_{\text{spe}} | c(x) = i\}}{\#D}$. One 146 $\#D_{\rm spe}$ of the main advantages of this method is its simplicity: the computation of the weights $\alpha^{\rm IS}$ is 147 decoupled from model optimization and can be performed before training. It is expected to work 148 well when the specialization set can be well approximated by the reweighted generic set, i.e., when 149 $L_{\text{spe}}(\boldsymbol{\theta}) \simeq L_{\text{gen}}(\boldsymbol{\theta}, \boldsymbol{\alpha}^{\text{IS}})$. When this is not the case, it might not lead to a good specific loss. Another 150 potential issue with this method arises when it assigns a large weight to a generic domain D_i with 151 little available data. In this case, training a model on $mix(\alpha^{IS})$ will overfit on that domain D_i , and 152 it would have been better to reduce the weight of that domain to mitigate overfitting. A last issue 153 arises when the number of specific examples, $\#D_{spe}$, is significantly smaller than the number of 154 domains k. In this situation, the importance weights become sparse, as they can have at most $\#D_{\text{spe}}$ 155 non-zero coefficients. This sparsity could be problematic, as some domains with zero weights might still be close to D_{spe} . We illustrate these shortcomings in our experiments and explain how gradient 156 alignment methods — which we introduce next — overcome them. 157

158 2.4 DGA: DYNAMIC GRADIENT ALIGNMENT

159

Algorithm. We introduce the DGA: Dynamic Gradient Alignment method for data reweighting 160 to approximately solve the bilevel problem in Equation 3. This algorithm builds upon DoGE (Fan 161 et al., 2024) and we give a precise account of their differences later. DGA keeps track of the model's

174

175

176

177

185 186 187

162 parameters θ^t and dynamic weights α^t . Once every T_r steps, we compute the gradient alignments 163 a^t , by doing 164 $(\Phi^t, \sigma^t) = \nabla \ell (\Theta^t, \sigma^t) + \nabla \ell (\Theta^t, \sigma^t)$ where $\sigma \to D$ and $\sigma \to D$. (4)

$$\boldsymbol{a}_{i}^{t} = \langle \nabla \ell(\boldsymbol{\theta}^{t}, \boldsymbol{x}_{i}), \nabla \ell(\boldsymbol{\theta}^{t}, \boldsymbol{z}) \rangle \text{ where } \boldsymbol{x}_{i} \sim D_{i} \text{ and } \boldsymbol{z} \sim D_{\text{spe}}.$$
 (4)

and update the weights by mirror descent on the simplex (Beck & Teboulle, 2003) with step $\eta > 0$:

$$\boldsymbol{\alpha}^{t+1} = \frac{\hat{\boldsymbol{\alpha}}}{\sum_{i=1}^{k} \hat{\boldsymbol{\alpha}}_{i}} \text{ where } \hat{\boldsymbol{\alpha}} = \boldsymbol{\alpha}^{t} \odot \exp(\eta \boldsymbol{a}^{t})$$
(5)

We optionally store an EMA version of the weights α^t parameterized by $\beta \in [0,1]$ to stabilize the training dynamics of the model's parameters, and define $\alpha_{\text{EMA}}^{t+1} = (1 - \beta)\alpha_{\text{EMA}}^t + \beta\alpha^{t+1}$. Finally, at each step, we update the model's parameters θ^t by doing an optimization step on $L_{\text{gen}}(\theta, \alpha_{\text{EMA}}^t)$. The full algorithm pseudo-code is given in Algorithm 1.

Rationale. This algorithm can be seen as a heuristic to solve the bilevel problem in Equation 3. Indeed, each update on θ optimizes the inner loss. The update rule on α can be seen as a mirrordescent step on $L_{\text{spe}}(\theta^*(\alpha))$ with several approximations. The first approximation consists of approximating the solution of the inner problem with one gradient descent step with step-size ρ : $\theta^*(\alpha) \simeq \theta^t - \rho \sum_{i=1}^k \alpha_i \nabla L_i(\theta^t)$. We then approximate the specific loss at θ^* by the post-update specific loss function: $L_{\text{spe}}(\theta^*(\alpha)) \simeq f(\alpha, \rho) \triangleq L_{\text{spe}}(\theta^t - \rho \sum_{i=1}^k \alpha_i \nabla L_i(\theta^t))$, that is, the drop on the specific loss after an update. When the step size ρ is small, a Taylor expansion gives

$$f(\boldsymbol{\alpha}, \rho) = L_{\rm spe}(\boldsymbol{\theta}^t) - \rho \sum_{i=1}^k \boldsymbol{\alpha}_i \langle \nabla L_i(\boldsymbol{\theta}^t), \nabla L_{\rm spe}(\boldsymbol{\theta}^t) \rangle + o(\rho)$$
(6)

Similarly, we get that the gradient of f is the gradient alignment:

$$\frac{\partial f}{\partial \boldsymbol{\alpha}_i}(\boldsymbol{\alpha}, \boldsymbol{\rho}) = -\rho \langle \nabla L_i(\boldsymbol{\theta}^t), \nabla L_{\rm spe}(\boldsymbol{\theta}^t) \rangle + o(\boldsymbol{\rho}) \tag{7}$$

We want to use this gradient of f to implement a mirror-descent method. Unfortunately, the gradients involved in the alignment are full-batch, so we approximate them with stochastic gradients obtained from mini-batches, yielding the alignments a^t from Equation 4. Overall, we get the approximation $\nabla_{\alpha} L_{\text{spe}}(\theta^*(\alpha)) \simeq -\rho a^t$; and the update rule in Equation 5 is a mirror descent step with this approximated gradient and step η/ρ .

We have explained the link between our algorithm and the bilevel problem in Equation 3. Proofs showing convergence of our method require assumptions violated in practice, e.g. most theoretical work assumes that the function $\theta \to L_{gen}(\theta, \alpha)$ is convex (Ghadimi & Wang, 2018; Arbel & Mairal, 2021; Dagréou et al., 2022). Nevertheless, successful applications of related bilevel algorithms to non-convex neural networks have been reported recently (Fan et al., 2024; Grangier et al., 2023).

Computational cost and memory overhead. The computation cost of DGA is compared to the cost of a regular pre-training run. For a base run iteration, the main cost is t_g , the cost of computing a gradient with a mini-batch *B*. For DGA, we need to add the cost of updating the domain weights α , which only happens every T_r iterations. This update requires computing the k + 1 gradients (one per domain, one for L_{spe}). Hence the average cost of one iteration of DGA is $(1 + (k + 1)T_r^{-1})t_g$. Therefore, DGA's compute overhead is small when T_r is large compared to the number of domains *k*.

204 During training, the memory is used by the optimizer state, the model gradients and its activations. 205 We assume the same precision for storing all vectors. The optimizer state (the model parameters and the two EMA terms for Adam) and the gradients have a storage cost of $4m_q$, where m_q denotes the 206 cost of storing the model parameters. The cost of storing the activations during backpropagation is 207 m_b . Regular pretraining with Adam therefore costs $4m_g + m_b$. DGA computes the required gradients 208 sequentially and does not require more memory to store activations. It simultaneously stores two 209 gradients instead of one (one domain gradient and one specific gradient). DGA, therefore, costs 210 $5m_q + m_b$: DGA memory overhead ranges from 0 (when $m_b \gg m_q$) to 25% (when $m_q \gg m_b$). 211

Comparison with DoGE. While our method is heavily inspired by DoGE (Fan et al., 2024), there are several key differences. First, DGA samples from the mixture: the weights θ^t are updated using samples drawn from the mixture mix(α^t), with the gradient $\nabla \ell(\theta^t, x)$ where $x \sim \min(\alpha^t)$; this is the same gradient that one would use during pre-training with weight α^t . In contrast, DoGE's weights are updated using a reweighted gradient $\sum_{i=1}^{k} \alpha_i^t \nabla \ell(\theta^t, x_i)$, where each x_i are drawn from

Alg	orithm 1 Dynamic Gradient Alignment method
1:	Input: Generic domains D_1, \ldots, D_k , specific set D_{spe} , inner optimizer state ω^0 , optimizer
	function Optimizer such as Adam or SGD, initial weights α^0 , outer learning rate η , EMA
	parameter β , weight update frequency T_r
2:	Initialize EMA weights: $\alpha_{\rm EMA}^0 = \alpha^0$
3:	for $t = 0 \dots T$ do
4:	Sample a batch from EMA generic mixture: $m{x} \sim \min(m{lpha}_{ ext{EMA}}^t)$
5:	Update the parameters $oldsymbol{ heta}^{t+1},oldsymbol{\omega}^{t+1} \leftarrow ext{Optimizer}(oldsymbol{ heta}^t,oldsymbol{\omega}^t, abla_{oldsymbol{ heta}}\ell(oldsymbol{ heta}^t,oldsymbol{x}))$
6:	if $t\%T_r = 0$ then
7:	Sample a batch from each domain: $x_i \sim D_i$ for $i = 1 \dots k$ and $y \sim D_{ m spe}$
8:	Compute gradient alignements $a_i^t \leftarrow \langle \nabla \ell(\theta^{t+1}, x_i), \nabla \ell'(\theta^{t+1}, y) \rangle$
9:	Update instantaneous weights: $\alpha^{t+1} \leftarrow \frac{\hat{\alpha}}{\sum_{i=1}^{k} \hat{\alpha}_i}$ with $\hat{\alpha} = \alpha^t \odot \exp(-\eta a^t)$
10:	Update EMA weights: $\boldsymbol{\alpha}_{\text{EMA}}^{t+1} \leftarrow \beta \boldsymbol{\alpha}_{\text{EMA}}^{t} + (1-\beta) \boldsymbol{\alpha}^{t+1}$
11:	else
12:	Do nothing: $\boldsymbol{\alpha}_{\mathrm{EMA}}^{t+1} \leftarrow \boldsymbol{\alpha}_{\mathrm{EMA}}^{t}$, and $\boldsymbol{\alpha}^{t+1} \leftarrow \boldsymbol{\alpha}^{t}$
13:	end if
14:	end for
15:	Return Optimized parameters $\theta^{(T)}$ and weights trajectory $\alpha^t, t = 0 \dots T$

the domain D_i . For a fixed number of samples available at each draw, DGA's gradient estimate has a lower variance (Seiffert et al., 2008). As explained above, DGA has a small overhead compared to regular pre-training, while DoGE updates the weights at each iteration. These two key differences mean that DGA is much closer to regular pre-training than DoGE. For instance, DGA never requires retraining a model from scratch using the mixture weights estimated from a previous run, while this is the costly strategy used for DoGE. Finally, the EMA strategy described above is novel.

Towards a convergence theory for DGA. It is hard to prove the convergence of algorithms for data reweighting with neural networks because of the non-convexity of the loss functions and the unknown link between generic and specialist datasets. We prove the convergence of DGA in a simplified setting where the generic losses are deterministic quadratic functions, and the specialist dataset is exactly a mixture of generic datasets with unknown proportions $\tilde{\alpha}$. We have:

Theorem 2. Let $\mu_1, \ldots, \mu_k \in \mathbb{R}^d$ some target vectors, and define the losses on the i^{th} generic domain as $L_i(\theta) = \frac{1}{2} \|\theta - \mu_i\|^2$. Let $\tilde{\alpha} \in \Delta_k$ a target mixture vector, and assume that the specific loss is $L_{spe}(\theta) = \sum_{i=1}^k \tilde{\alpha}_i L_i(\theta)$. Let $M = [\mu_1, \ldots, \mu_k] \in \mathbb{R}^{d \times k}$, assume that $\lambda_{\min}(M^T M) > 0$. Then, running DGA for T iterations with gradient descent as the optimizer function with step size 1, and outer learning rate $\eta = O(1/\sqrt{T})$ yields iterates α^t such that $\min_{t=1...T} \|\alpha^t - \tilde{\alpha}\|^2 = O(1/\sqrt{T})$.

This theorem demonstrates that DGA converges at the same rate as mirror descent for a simple problem and recovers the true mixture weights. To the best of our knowledge, this is the first theoretical convergence result for a gradient alignment descent.

258 3 EXPERIMENTS

254

257

Our experiments focus on two challenging cases. First, given *limited token resources* within each training domain, the model would risk overfitting with weights concentrated on a few domains.
 Second, given *large number of training domains*, applying DGA on domain reweighting could introduce intractable computation overheads linearly increasing according to the domain granularity.

Generic Datasets and Domains. For all the experiments, we use Redpajama-v2 (Together AI Team, 2023) as the generic training set D_{gen} . This is one of the largest public corpus for LLM pretraining. Redpajama-v2 contains 30 trillion filtered and deduplicated tokens from web-crawled dumps. Since this corpus does not come pre-segmented into domains, we obtain individual generic domains from D_{gen} with clustering. Specifically, we use the embedding-and-clustering pipeline from Grangier et al. (2024). We first embed all the training sequences $x \in D_{\text{gen}}$ with SentenceBert (all-MiniLM-L6-v2), yielding a 384 dimensional embedding Bert(x). We then apply k-means clustering on the sentence embeddings into k = 64 clusters yielding k domains D_1, \ldots, D_k . To get fine-grained generic domains, we apply hierarchical clustering on the top of the first level of $k_1 = 64$ clusters. Specifically, each domain is further clustered once again into 64 smaller clusters. We apply this strategy twice to get domains with granularity $k_2 = 64^2 = 4096$ and $k_3 = 64^3 = 262k$.

Model Architecture. We train small (125M), medium (350M) and large (750M) models with decoderonly transformers (Vaswani et al., 2017). We adopt most of the training settings and architectures from (Brown et al., 2020a). Their details are provided in Appendix C. For optimization, we use the AdamW optimizer (Loshchilov, 2017).

3.1 DOMAIN REWEIGHTING WITH LIMITED RESOURCES



Figure 1: Comparing data reweighting methods with *free_law* as a specific set in a low generic data regime. When there are not enough tokens, importance sampling quickly overfits, while DGA manages to explore the training distributions to avoid overfitting. We see the importance of the EMA to stabilize DGA in the low data regime. When there is no token limit, adding an EMA ($\beta = 0.1$) does not negatively affect the performance.

296

273

274

275

276

277 278

279

Previous works on domain reweighting implicitly assume infinite token resources from all training domains (Xie et al., 2023a; Fan et al., 2024; Liu et al., 2024) while it is not always applicable in real-world cases. The scenario with limited training resources is challenging for online domain reweighting. Indeed, if the weights are concentrated on a few domains, e.g. on a single domain D_i , a large model will quickly overfit when the number of tokens in D_i is small.

We expect DGA to mitigate overfitting by dynamically adjusting the domain weights. Specifically, once a model starts overfitting on D_i , the magnitude of the gradients $\nabla \ell(\theta, D_i)$ decreases as its training loss $\ell(\theta, D_i)$ is low, i.e. the domain knowledge from D_i is well-learned. Consequently, the corresponding gradient alignment score $a_i = \langle \nabla \ell(\theta, D_i), \nabla \ell(\theta, D_{spe}) \rangle$ decreases as well and DGA explores other domains with higher alignment scores. In other words, DGA down-weights domains once they are well-learned, thereby achieving a balance between *exploration* – by learning from diverse data domains – and *exploitation*, by intensively training on the most relevant domains.

However, with limited data per domain, we remark that DGA without EMA demonstrates drastic changes at each domain weight update, focusing heavily on one domain at a time. Quickly changing domain weights is problematic since we want to use the same domain weights for T_r steps in the future. This motivates the introduction of the EMA update in Algorithm 1, which regularizes the model and domain weights with the previous state when it starts to overfit.

Experiment Setup. We consider the generic domain split into k = 64 domains. We construct three scales of generic sets, either taking the full dataset or randomly sub-sampling 30M, 0.1B tokens per domain. For the targeted specific set D_{spe} , we use 5 subsets from *the Pile* (Gao et al., 2020) covering common specialized data types for LM applications: Math (*dm_mathematics*), Code (*github*, *stackexchange*), Medical (*pubmed_central*), Legal (*free_law*) and Scientific articles (*arxiv*).

We implement the importance sampling baseline described in subsection 2.3. We also compare to the *uniform baseline* with the domain weights $\alpha_{uniform}$ as the natural proportion of each data domain in the generic Redpajama-v2 dataset. For importance sampling and uniform baselines, the domain weights are fixed throughout the entire training run. For both vanilla DGA and DGA with an EMA ($\beta = 0.1$), we update domain weights α every $T_r = 100$ steps. We provide the ablation results on the step size η , frequency T_r and ema factor β in Appendix C. We use 125M models for experiments.

340

341

354 355

324 **Results.** We report the validation loss on the specialized set under various token constraints in 325 Figure 1 for *free_law* and the results on other domains in Appendix A. With 30M tokens per domain, 326 DGA with EMA effectively stabilizes the training, while vanilla DGA exhibits several loss spikes, 327 suggesting a lack of robustness. Under a 0.1B token constraint, both DGA and DGA with EMA are 328 able to dynamically adjust domain weights to mitigate overfitting. In contrast, fixed domain weights from importance sampling consistently lead to overfitting in token-limited scenarios, demonstrating the limitations of static weighting strategies in comparison to dynamic approaches like DGA. It is 330 worth noting that adding the EMA has no negative effect on the learning efficacy when there is no 331 token limit, which can be used as a robust regularization in the online domain reweighting context. 332

Domain Weights Evolution. In the experiments with a limited generic token budget (subsection 3.1),
 DGA without EMA often assigns excessive weight to one generic domain, leading to overfitting
 due to the restricted number of training tokens. This iterative over-weighting pattern on generic
 domain weights aligns with the observed loss spikes on the specific set (Figure 3a). In contrast, the
 EMA helps to regularize the weight dynamics, effectively preventing the model from overfitting by
 maintaining more balanced domain weights throughout the training process.

3.2 DISTRIBUTION REWEIGHTING: SCALING-UP DATA MIXING ON EXTREMELY FINE-GRAINED DATA DOMAINS

The computational overhead from DGA scales linearly with the number of domain k. This is intractable for datasets segmented in many fine-grained domains and, consequently, prior domain reweighting methods (Xie et al., 2023a; Fan et al., 2024; Liu et al., 2024; Kang et al., 2024) have not been applied in that setting. The fine-grained setting motivates *distribution reweighting* as an alternative to direct *domain reweighting*.

347 Distribution reweighting leverages the strength from both embedding-based (importance sampling) 348 and gradient-based (DGA) strategies. We consider a generic training set partitioned into k domains 349 with a large k (e.g. 4096, 262k). We also have a set of N auxiliary datasets $\{S_1, \ldots, S_N\}$, called 350 basis sets, each from a specific domain of interest. We compute the importance sampling histograms 351 for each basis set as $P = \{p_1, \ldots, p_N\}$, $p_i \in \Delta^k, P \in \mathbb{R}^{k \times N}$. We then use DGA to search over 352 a reparameterized space leveraging this basis. We define the domain weights $\alpha_{\text{domain}} \in \Delta^k$ as a 353 convex combination of N k-dimensional distributions derived from importance sampling,

$$\boldsymbol{\alpha}_{\text{domain}} \approx P \boldsymbol{\alpha}_{\text{dist}} = \alpha_{\text{dist},1} \cdot \boldsymbol{p}_1 + \alpha_{\text{dist},2} \cdot \boldsymbol{p}_2 + \ldots + \alpha_{\text{dist},N} \cdot \boldsymbol{p}_N \tag{8}$$

where the low-dimensional weights $\alpha_{dist} \in \Delta^N$ are learned by DGA. This allows the use of fine-356 grained domain features while eliminating intensive gradient computation on each generic domain. 357 Compared to the $(k+1)/T_r$ overheads from domain reweighting, applying distribution reweighting 358 only incurs $(N + 1)/T_r$ extra budget, where N is typically much smaller than k. Importantly, this is 359 equivalent to applying DGA with the N generic domains $\tilde{D}_1, \ldots, \tilde{D}_N$ where $\tilde{D}_i = \min(\boldsymbol{p}_i)$. Hence, 360 it does not require any modification to the base DGA algorithm; it suffices to be able to sample 361 according to each mix(p_i). We provide the pseudo-code for the distribution reweighting with DGA 362 in Appendix D. 363

Experiment Setup. We demonstrate the efficacy of distribution reweighting on the MMLU benchmark (Hendrycks et al., 2021). MMLU consists of 57 tasks from various knowledge fields, which serves as a testbed of multi-domain language modeling; by measuring the downstream accuracy, we can assess whether the improvements obtained in language modeling transfer to reasoning abilities.

We construct two specific datasets with different amounts of accessible samples: (1) MMLU_a: we take half of the examples from each task used as D_{spe} . We denote the other half of datapoints as MMLU_b, which is used for evaluation; (2) MMLU_dev: we randomly select 5 samples from each task, simulating the few-shot learning scenario. MMLU_a has 7.1k samples while MMLU_dev only has 285 samples, which yields sparse importance sampling histograms. For evaluation, we assess the language modeling performance by computing perplexity on MMLU_b. We also measure the accuracy for multiple choice question answering on MMLU_b with llm-eval (Gao et al., 2024).

We use generic domain splits with k=64,4096,262k domains. We rely on 22 auxiliary subdomains from *The Pile* (Gao et al., 2020) as our basis sets. For each auxiliary set, we take 15M tokens and compute their importance-sampling histograms as $p_1, \ldots, p_N \in \Delta^k$. To search for the optimal balance between diversity and specificity, we extend the basis

379

380

409

410

sets with the importance sampling histogram from the specific set itself (i.e. MMLU_a or MMLU_dev), yielding N = 23 distributions. For this experiment, we use 750M models.

381 DGA greatly accelerates task-adaptive language model-382 ing. We evaluate the model's ability to acquire specialized knowledge on the target task (MMLU) in Appendix (Fig-384 ure 12). DGA achieves substantial training speed-ups 385 compared to uniform sampling, accelerating by $6.5 \times$ on 386 MMLU_a and $4.3 \times$ on MMLU_dev. Moreover, with finer-387 grained clustering (k=262k), DGA outperforms impor-388 tance sampling, which suffers from sparse histograms and performance degradation, achieving $2 \times$ and $3.2 \times$ faster 389 training on MMLU_a and MMLU_dev, respectively. These 390 results highlight DGA's ability to effectively utilize fine-391 grained domain information while avoiding overfitting, 392 demonstrating its robustness in scenarios with limited spe-393 cialized samples. We include details in subsection B.1 394 with the fine-tuning results in subsection G.1. 395

Distribution reweighting makes better specialist and 396 generalist. While many task-adaptive algorithms are 397 prone to catastrophic forgetting of general knowledge, 398 we demonstrate that DGA with distribution reweighting 399 effectively balances the acquisition of specialized knowl-400 edge while preserving general knowledge. To evaluate this, 401 we report the loss on MMLU_b, representing specialized 402 knowledge, alongside the average validation loss across 403 22 domains in The Pile, which reflects general knowledge 404 drawn from broad and diverse domains (Figure 2). Compared to both importance sampling and uniform sampling 405 baselines, DGA with distribution reweighting achieves a 406 407 superior Pareto front, illustrating an improved trade-off between specialized performance and general knowledge. 408



4 DISCUSSION AND LIMITATIONS



Comparison between domain reweighting and distribution reweighting. As an efficient alterna-411 412 tive of direct domain reweighting, we assess distribution reweighting in terms of specialized task adaptation and general domain losses. As shown in Figure 13, distribution reweighting outper-413 forms domain reweighting in both specialized and general domain losses on k=64 clusters, while 414 also incurring lower compute overhead. Furthermore, distribution reweighting exhibits remarkable 415 scalability, with substantial improvements in specialized task perplexity as cluster granularity in-416 creases (k=4096,262k). In contrast, domain reweighting struggles with scalability due to its high 417 computational complexity, underscoring the efficiency and robustness of distribution reweighting in 418 fine-grained settings. 419

DGA outperforms DoGE on both specialized task adaptation and general knowledge. We
 compare DGA with DoGE (Fan et al., 2024) in the context of task-adaptive pretraining. Comparing
 to both proxy model with online tuned domain weights and the base model with fixed optimized
 domain weights, DGA greatly outperforms DoGE in both specialized loss and general loss evaluated
 on the generic set. We present more details in Appendix F.

Language modeling capability cannot be fully translated into reasoning accuracies. According to
 Table 1, both importance sampling and DGA reweighting greatly outperform the uniform baseline on
 reasoning accuracy scores. However, despite the superior language modeling ability, DGA performs
 comparably as importance sampling in terms of accuracy scores. It indicates that better language
 modeling ability may not be fully transferable to better reasoning capabilities. We report the full
 results with different model scales in Appendix B with a detailed discussion on potential reasons.

431 Weights Evolution on Distributions. We present the evolution of domain weights for each basis distribution from DGA in Figure 3. Comparing different levels of granularity, with k=262k, the impor-



(30M tokens per domain) 452

Figure 3: The top row presents the specific loss over time, with the two bottom rows illustrating the 454 evolution of domain (dist.) weights from DGA correspondingly, with each line representing a distinct 455 domain. Left: Weights from the limited generic token experiment (subsection 3.1). Middle and 456 **Right**: Weights from the distribution reweighting experiment (subsection 3.2). The thick black line highlights the dynamic weights assigned by DGA on the MMLU importance sampling distribution, 458 which serves as a fixed training distribution for the importance sampling runs. 459

457

453

461 tance of the MMLU distribution is more emphasized than with k=4096, with the help of fine-grained 462 domain features. Additionally, with sufficient samples from the specific domain (MMLU_a, Figure 3b), 463 the MMLU distribution is consistently up-weighted across 262k generic domains. In contrast, on 464 MMLU_dev, while the distribution on MMLU is initially up-weighted, it declines gradually in the late stage of training. Owing to the number of accessible samples from the specific set, the importance 465 sampling distribution on MMLU_dev across 262k generic domains is very sparse. During the training, 466 the learnability of the few activated generic domains diminishes, making other distributions more 467 beneficial to the model. 468

469 In addition to the importance sampling distribution from the specific sets (MMLU_a and MMLU_dev), 470 DGA effectively identifies other relevant distributions from The Pile that contribute to the learning on MMLU. These influential distributions, which include philpapers, freelaw, and 471 dm_mathematics, are all considered to contain high-quality, academic-related contents. We 472 present detailed curves with domain labels in subsection B.4. This ability to adaptively select benefi-473 cial distributions enhances the model's generalization and helps mitigate overfitting by leveraging a 474 broader yet pertinent set of data sources during pretraining. 475

476 Impact of generic domain granularity. In Figure 4, we present the validation loss on the specific 477 domain according to the number of clusters within the generic dataset. From k = 64 to 4096, both DGA and importance sampling demonstrate significant improvement in language modeling in 478 terms of validation loss (i.e., log of perplexity). However, when the number of clusters exceeds the 479 scale of the accessible specific set, the importance sampling method overfits the limited number of 480 activated generic domains, failing to capture broader domain knowledge. In contrast, DGA effectively 481 leverages extremely fine-grained domain information across 262k generic domains with only 7k482 samples from MMLU_a. In the few-shot context (MMLU_dev), DGA mitigates a large performance 483 degradation by utilizing diverse domain knowledge from other relevant distributions. 484

Scaling performance of DGA on model scales. To examine the scaling performance of the 485 DGA algorithm, we train three models of varying scales (125M, 350M, and 750M) using both uniform sampling and DGA on *k*=64 generic clusters. Model pretrained on DGA consistently outperforms uniform sampling baseline on specialized loss while exhibits degradation in general loss, i.e. validation loss on RedPajama-v2. We provide more details in subsection G.2.

5 RELATED WORK

486

487

488

489

490

513

527

491 Task-adaptive Data Selection for Domain-Specific 492 LLMs. Many works have shown that one can effectively 493 improve the LLM's performance on a specific downstream 494 task with data selection according to the relevance of 495 generic data for the targeted data domain. Gururangan et al. (2020) show that continued pretraining on data with 496 high vocabulary overlap can boost its performance on the 497 specific end-task. On machine translation task, Aharoni & 498 Goldberg (2020) identify task-relevant pretraining datasets 499 from a generic corpus using nearest neighbor of a small 500 specialist dataset based on SentenceBert sentence repre-501 sentation. Wang et al. (2020); Grangier et al. (2023) train a 502 small proxy model to give an importance weight per sam-503 ple. Xie et al. (2023b) proposed DSIR as a lexical-based 504 importance sampling method using n-gram features. 505

Other than feature-based importance sampling (Grangier et al., 2024), influence function-based method select data points which leads to the greatest loss drop on the target from the optimization perspective (Koh & Liang, 2020; Kwon et al., 2024; Agarwal et al., 2017). However, these methods often introduce intensive computational overheads from the second-order gradient computations, which is not applicable on large generic pretraining corpus.



Data Resampling through Domain Reweighting.
Given the large scale of the generic pretraining corpus, sample-level selection strategies are hard to implement for LLM pretraining. Alternatively, domain reweighting methods (Xie et al., 2023a; Fan et al., 2024; Liu et al., 2024;



Kang et al., 2024) apply group-level selection by adjusting data sampling weights across different 519 domains to reflect their importance. Based on the weak-to-strong generalization strategy (Burns et al., 520 2023), existing domain re-weighting methods typically estimate the optimal domain weights for a 521 larger model based on the preferences of a small-scale proxy model. Xie et al. (2023a) apply group 522 distributed robust optimization to optimize the worst-case loss gap between two small-scale proxies. 523 Fan et al. (2024) use gradient alignment to dynamically adjust domain weights during proxy model 524 training. Specifically, it identifies the most beneficial domains by aligning the gradients of the training 525 data with the target task. However, it trains the proxy model on reweighted domain gradients to simulate the resampling scenario, which introduces more variance in the domain weights estimation. 526

528 6 CONCLUSION

529 To tackle two key challenges of online domain reweighting, we introduce Dynamic Gradient Align-530 ment (DGA) as a stable and scalable data mixing method for LLM pretraining. Given a target task, DGA is an online algorithm that adjusts the training data distribution according to the current 531 model status. This adaptation relies on an estimate of the progress on the target task from gradient 532 alignments. We show that under limited tokens within generic domains, DGA with EMA can notably 533 mitigate overfitting and yields superior performance on the end-task by balancing exploitation and 534 exploration. We also propose a novel distribution reweighting strategy, which enables DGA to scale 535 up to extremely fine-grained data domains without incurring intensive computations. Our experiments 536 on MMLU show that applying distribution reweighting with DGA effectively leverages fine-grained 537 domain knowledge to balance specialty and diversity during training. Our work demonstrates the 538 scalability of gradient-alignment-based data reweighting methods, as well as their efficiency in data-constrained settings.

540 REFERENCES

549

552

553

554

567

583

584

Naman Agarwal, Brian Bullins, and Elad Hazan. Second-order stochastic optimization for machine
 learning in linear time, 2017.

- Roee Aharoni and Yoav Goldberg. Unsupervised domain clusters in pretrained language models. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 7747–7763, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.692. URL https://aclanthology.org/2020.acl-main.692.
- Michael Arbel and Julien Mairal. Amortized implicit differentiation for stochastic bilevel optimization.
 arXiv preprint arXiv:2111.14580, 2021.
 - Amir Beck and Marc Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31(3):167–175, 2003.
- Jerome Bracken and James T McGill. Mathematical programs with optimization problems in the constraints. *Operations research*, 21(1):37–44, 1973.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhari-558 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-559 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz 561 Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In 563 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), Advances in Neu-564 ral Information Processing Systems, volume 33, pp. 1877-1901. Curran Associates, Inc., 565 2020a. URL https://proceedings.neurips.cc/paper_files/paper/2020/ 566 file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020b. URL https: //arxiv.org/abs/2005.14165.
- Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner, Bowen Baker, Leo Gao, Leopold Aschenbrenner,
 Yining Chen, Adrien Ecoffet, Manas Joglekar, Jan Leike, Ilya Sutskever, and Jeff Wu. Weak-to strong generalization: Eliciting strong capabilities with weak supervision, 2023. URL https:
 //arxiv.org/abs/2312.09390.
- Mathieu Dagréou, Pierre Ablin, Samuel Vaiter, and Thomas Moreau. A framework for bilevel optimization that enables stochastic and global variance reduction algorithms. *Advances in Neural Information Processing Systems*, 35:26698–26710, 2022.
 - Simin Fan, Matteo Pagliardini, and Martin Jaggi. Doge: Domain reweighting with generalization estimation, 2024. URL https://arxiv.org/abs/2310.15393.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. The pile: An 800gb dataset of diverse text for language modeling, 2020. URL https://arxiv.org/abs/2101. 00027.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster,
 Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff,
 Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika,
 Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot
 language model evaluation, 07 2024. URL https://zenodo.org/records/12608602.

- 594 Saeed Ghadimi and Mengdi Wang. Approximation methods for bilevel programming. arXiv preprint 595 arXiv:1802.02246, 2018. 596 David Grangier, Pierre Ablin, and Awni Hannun. Adaptive training distributions with scalable online 597 bilevel optimization. arXiv preprint arXiv:2311.11973, 2023. 598 David Grangier, Angelos Katharopoulos, Pierre Ablin, and Awni Hannun. Specialized language 600 models with cheap inference from limited domain data. arXiv preprint arXiv:2402.01093, 2024. 601 602 Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. Don't stop pretraining: Adapt language models to domains and tasks, 2020. URL 603 https://arxiv.org/abs/2004.10964. 604 605 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob 606 Steinhardt. Measuring massive multitask language understanding, 2021. URL https://arxiv. 607 org/abs/2009.03300. 608 Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong 609 Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. A survey on hallucination 610 in large language models: Principles, taxonomy, challenges, and open questions, 2023. URL 611 https://arxiv.org/abs/2311.05232. 612 613 Feiyang Kang, Yifan Sun, Bingbing Wen, Si Chen, Dawn Song, Rafid Mahmood, and Ruoxi Jia. 614 Autoscale: Automatic prediction of compute-optimal data composition for training llms, 2024. 615 URL https://arxiv.org/abs/2407.20177. 616 Teun Kloek and Herman K Van Dijk. Bayesian estimates of equation system parameters: an 617 application of integration by monte carlo. Econometrica: Journal of the Econometric Society, pp. 618 1-19, 1978. 619 620 Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions, 2020. 621 Yongchan Kwon, Eric Wu, Kevin Wu, and James Zou. Datainf: Efficiently estimating data influence 622 in lora-tuned llms and diffusion models, 2024. 623 624 Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human 625 falsehoods, 2022. URL https://arxiv.org/abs/2109.07958. 626 Qian Liu, Xiaosen Zheng, Niklas Muennighoff, Guangtao Zeng, Longxu Dou, Tianyu Pang, Jing 627 Jiang, and Min Lin. Regmix: Data mixture as regression for language model pre-training, 2024. 628 URL https://arxiv.org/abs/2407.01492. 629 630 Shayne Longpre, Gregory Yauney, Emily Reif, Katherine Lee, Adam Roberts, Barret Zoph, Denny 631 Zhou, Jason Wei, Kevin Robinson, David Mimno, and Daphne Ippolito. A pretrainer's guide to 632 training data: Measuring the effects of data age, domain coverage, quality, toxicity, 2023. URL 633 https://arxiv.org/abs/2305.13169. 634 I Loshchilov. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017. 635 636 Nils Reimers and Iryna Gurevych. Sentence-BERT: Sentence embeddings using Siamese BERT-637 networks. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), Proceedings of the 638 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pp. 3982–3992, Hong Kong, 639 China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1410. 640 URL https://aclanthology.org/D19-1410. 641 642 Chris Seiffert, Taghi M Khoshgoftaar, Jason Van Hulse, and Amri Napolitano. Resampling or 643 reweighting: A comparison of boosting implementations. In 2008 20th IEEE international 644 conference on tools with artificial intelligence, volume 1, pp. 445–451. IEEE, 2008. 645
- Together AI Team. Redpajama-data-v2: An open dataset with 30 trillion tokens for training large language models, October 2023. URL https://www.together.ai/blog/ redpajama-data-v2.

648	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
649	Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand
651	Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language
650	models, 2023. UKL https://arxiv.org/abs/2302.139/1.
652	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz
654	Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio,
655	H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), Advances in Neural Information
656	Processing Systems, volume 30. Curran Associates, Inc., 2017.
657	Xinvi Wang Hieu Pham Paul Michel Antonios Anastasonoulos Jaime Carbonell and Graham
658	Neubig. Optimizing data usage via differentiable rewards. In <i>International Conference on Machine</i>
659	Learning, pp. 9983–9995. PMLR, 2020.
660	
661	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maartan Daama, Danny Zhay, Danald Matalan Ed H. Chi, Tatayangi Hashimata, Orial Vinyala
662	Dergy Liang, Leff Deen, and William Fedux. Emergent abilities of large language models, 2022
663	IIRL https://arxiv.org/abs/2206_07682
664	GRE heeps., / arxiv.org/abs/2200.07002.
665	Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy Liang, Quoc V.
666	Le, Tengyu Ma, and Adams Wei Yu. Doremi: Optimizing data mixtures speeds up language model
667	pretraining, 2023a. URL https://arxiv.org/abs/2305.10429.
668	Sang Michael Xie, Shibani Santurkar, Tengyu Ma, and Percy Liang. Data selection for language
669	models via importance resampling, 2023b. URL https://arxiv.org/abs/2302.03169.
670	
671	Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher
672	Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Iodor Mihaylov, Myle Ott, Sam Shleiter, Kurt
673	Opt: Open pre-trained transformer language models 2022 UPL https://orwiw.org/obc/
674	2205_01068
675	2203.01000.
676	
677	
678	
679	
680	
681	
682	
683	
684	
685	
686	
687	
688	
689	
601	
602	
603	
694	
695	
696	
697	
698	
699	
700	

A TRAINING WITH LIMITED GENERIC TOKENS

A.1 VALIDATION LOSS ON THE TARGETED END-TASK

We present the complete results on all six target domains (arxiv, free_law, dm_mathematics, pubmed_central, github, stackexchange) as follows. Across all six target domains, DGA with EMA ($\beta = 0.1$) consistently stablize the training and yields better language modelling performance under token-limited contexts.



Figure 5: Results on all the domains for the low data experiment (subsection 3.1). The specific domain is free_law.



Figure 6: Results on all the domains for the low data experiment (subsection 3.1). The specific domain is arxiv



Figure 7: Results on all the domains for the low data experiment (subsection 3.1). The specific domain is dm-mathematics



Figure 8: Results on all the domains for the low data experiment (subsection 3.1). The specific domain is github



Figure 9: Results on all the domains for the low data experiment (subsection 3.1). The specific domain is pubmed-central



Figure 10: Results on all the domains for the low data experiment (subsection 3.1). The specific domain is stackexchange

A.2 DOMAIN WEIGHTS EVOLUTION

812 We present the domain weights evolution on 64 generic domains from DGA with and w.o. EMA
813 regularization. With both stackexchange and free_law as the specific set, EMA effectively
814 smoothes the spiky domain weights, which therefore stablize the training process.



Figure 11: Comparing data reweighting methods with stackexchange (resp. free_law) as the specific set, in a low generic data regime. When there are not enough tokens, importance sampling quickly overfits, while DGA manages to explore the training distributions to avoid overfitting. We see the importance of the EMA to stabilize DGA in the low data regime.

B DISTRIBUTION REWEIGHTING

B.1 DGA WITH DISTRIBUTION REWEIGHTING ACCELERATE TASK-ADAPTIVE TRAINING

We evaluate the perplexity on the target task (MMLU) in Figure 12, which reflects the model's ability to acquire specialized knowledge. Notably, DGA achieves significant acceleration compared to uniform sampling, with training speed-ups of $6.5 \times$ and $4.3 \times$ when targeting MMLU_a and MMLU_dev, respectively. Comparing to the importance sampling baseline, DGA achieves comparable performance when trained with MMLU_a on k = 4096 generic clusters. However, with specialized dataset as MMLU_dev, where samples from the specialized task are limited, DGA demonstrates a $2 \times$ speed-up on k = 4096 clusters. Additionally, with a finer-grained generic clustering (k = 262k), DGA accelerates training by $2 \times$ with MMLU_a and $3.2 \times$ with MMLU_dev. The remarkable improvements highlight DGA's capability to effectively leverage fine-grained domain information and mitigate overfitting issues often seen with importance sampling with insufficient samples from specialized task.



Figure 12: Perplexity on the Specific Task (MMLU). (a,c) present the specialized perplexity over time with target at MMLU_a; (b,d) present the results with target at MMLU_dev.

B.2 COMPARISON BETWEEN DOMAIN REWEIGHTING AND DISTRIBUTION REWEIGHTING

In Figure 13, we compare distribution reweighting and domain reweighting according to the taskadaptive capability and the general knowledge preserved during task-adaptive training process. With distribution reweighting on N = 23 distributions and k = 64 domains, both specialized and generic perplexity are improved above domain reweighting on k = 64 domains. Meanwhile, the computational overheads are reduced from $k/T_r = 64\%$ to $N/T_r = 23\%$. The perplexity on the specialized domain can be significantly improved by increasing the domain granularity $(k = 64 \rightarrow 4096 \rightarrow 262k)$, while the domain reweighting algorithm cannot scale up to large number of domains because of its high complexity.



(c) Generic Perplexity w. target at MMLU_dev

(d) Specialized Perplexity w. target at MMLU_dev

Figure 13: Comparison between Domain Reweighting and Distribution Reweighting. (a,c) present the perplexity across 22 domains in *The Pile*, which defined as the *generic perplexity*; (b,d) present the *specialized perplexity* on MMLU.

972 B.3 EVALUATION RESULTS ON MMLU

976 B.3.1 REASONING ACCURACY

Despite the superior performance on language modeling, DGA can hardly outperform importance sampling in terms of accuracy scores on MMLU benchmark (Table 1), which indicates the improvement on language modeling ability may not be able to fully translated to reasoning capacities. For the limited improvement on accuracy, we provide two potential explanations: (1) since the target objective in the outerloop (L_{spe} in Equ. 3) is the next token prediction loss for language modeling, the language modeling performance is expected to be improved. However, it does not necessarily translate to the improved reasoning accuracy (2) MMLU is considered to be a challenging reasoning task, where the accuracies can hardly be improved when the model capacity is below some specific threshold. This phenomenon is well-illustrated by Figure 11. (G) in (Wei et al., 2022). How to derive an optimization objective which directly benefits the reasoning accuracy would be a compelling future direction.

Table 1: MMLU accuracies with domain reweighting methods. Both importance sampling and
 DGA reweighting greatly improve the accuracy above uniform baseline, while DGA does not show significant improvement above importance sampling.

Method		MMLU_a	MMLU_dev
Uniform		26.1 %	26.1 %
Importance S.	$\begin{array}{c} k = 4096 \\ k = 262k \end{array}$	27.7 % 28.4 %	27.7 % 27.0 %
DGA dist. reweighting	$\begin{array}{c} k = 4096 \\ k = 262k \end{array}$	26.8 % 28.0 %	27.4 % 27.0 %

1026 B.3.2 Full Evaluation Results across Different Model Scales

1028We present the complete evaluation results on MMLU benchmark on small- (125M) and large-1029(750M) scale models. k denotes the number of generic domains, N denotes the number of reweighted1030importance sampling distributions from basis sets. N=22 indicates we only reweight 22 distributions1031from 22 The Pile subsets, while N=23 includes the importance sampling histgram from the specific1032set (MMLU). Since the 125M model shows marginal difference in accuracy because of limited1033capacity, we only scored 750M model on MMLU reasoning accuracies.

Table 2: Results on the domain reweighting experiment, with half MMLU as train set. The best results is **Bolded** and the second best is <u>Underlined</u>.

1037	125M model	MMLU loss	MMLU acc.	avg. Pile loss
1030	Uniform	3.56	_	3.27
1040	Importance S. (k=4096)	3.32	-	3.16
1040	Importance S. (k=262k)	3.22	-	3.19
1041	DGA domain reweighting (k=64)	3.31	-	<u>3.10</u>
1042	DGA dist. reweighting (N=22, k=4096)	3.34	-	3.13
1043	DGA dist. reweighting (N=22, k=262k)	3.34	-	3.05
1044	DGA dist. reweighting (N=23, k=4096)	3.33	-	3.13
1045	DGA dist. reweighting (N=23, k=262k)	<u>3.25</u>	-	<u>3.10</u>
1010				
1046	750M model	MMLU loss	MMLU acc.	avg. Pile loss
1046 1047 1048	750M model Uniform	MMLU loss 3.03	MMLU acc. 26.1 %	avg. Pile loss 2.82
1046 1047 1048 1049	750M model Uniform Importance S. k=4096	MMLU loss 3.03 2.82	MMLU acc. 26.1 % 27.7 %	avg. Pile loss 2.82 2.75
1046 1047 1048 1049	750M model Uniform Importance S. k=4096 Importance S. k=262k	MMLU loss 3.03 2.82 2.82	MMLU acc. 26.1 % 27.7 % 28.4 %	avg. Pile loss 2.82 2.75 2.81
1046 1047 1048 1049 1050	750M model Uniform Importance S. k=4096 Importance S. k=262k DGA domain reweighting (k=64)	MMLU loss 3.03 2.82 2.82 2.97	MMLU acc. 26.1 % 27.7 % 28.4 % 27.1 %	avg. Pile loss 2.82 2.75 2.81 2.77
1046 1047 1048 1049 1050 1051	750M model Uniform Importance S. k=4096 Importance S. k=262k DGA domain reweighting (k=64) DGA dist. reweighting (N=22, k=4096)	MMLU loss 3.03 2.82 2.82 2.97 2.86	MMLU acc. 26.1 % 27.7 % 28.4 % 27.1 % 27.2 %	avg. Pile loss 2.82 2.75 2.81 2.77 2.73
1046 1047 1048 1049 1050 1051 1052	750M model Uniform Importance S. k=4096 Importance S. k=262k DGA domain reweighting (k=64) DGA dist. reweighting (N=22, k=4096) DGA dist. reweighting (N=22, k=262k)	MMLU loss 3.03 2.82 2.82 2.97 2.86 2.84	MMLU acc. 26.1 % 27.7 % 28.4 % 27.1 % 27.2 % 27.0 %	avg. Pile loss 2.82 2.75 2.81 2.77 <u>2.73</u> 2.66
1046 1047 1048 1049 1050 1051 1052 1053	750M model Uniform Importance S. k=4096 Importance S. k=262k DGA domain reweighting (k=64) DGA dist. reweighting (N=22, k=4096) DGA dist. reweighting (N=22, k=262k) DGA dist. reweighting (N=23, k=4096)	MMLU loss 3.03 2.82 2.82 2.97 2.86 2.84 2.83	MMLU acc. 26.1 % 27.7 % 28.4 % 27.1 % 27.2 % 27.0 % 26.8 %	avg. Pile loss 2.82 2.75 2.81 2.77 <u>2.73</u> 2.66 2.73
1046 1047 1048 1049 1050 1051 1052 1053 1054	750M model Uniform Importance S. k=4096 Importance S. k=262k DGA domain reweighting (k=64) DGA dist. reweighting (N=22, k=4096) DGA dist. reweighting (N=22, k=262k) DGA dist. reweighting (N=23, k=4096) DGA dist. reweighting (N=23, k=262k)	MMLU loss 3.03 2.82 2.82 2.97 2.86 2.84 2.83 2.74	MMLU acc. 26.1 % 27.7 % 28.4 % 27.1 % 27.2 % 27.0 % 26.8 % <u>28.0</u> %	avg. Pile loss 2.82 2.75 2.81 2.77 <u>2.73</u> 2.66 <u>2.73</u> 2.76

Table 3: Results on the domain reweighting experiment, with 5 examples per task of MMLU as train
 set. We score only the 750M models.

125M model	MMLU loss	MMLU acc.	avg. Pile loss
Uniform	3.56	-	3.27
Importance S. (k=4096)	3.40	-	3.16
Importance S. (k=262k)	3.41	-	3.29
DGA domain reweighting (k=64)	3.46	-	3.19
DGA dist. reweighting (N=22, k=4096)	3.37	-	3.12
DGA dist. reweighting (N=22, k=262k)	<u>3.36</u>	-	3.03
DGA dist. reweighting (N=23, k=4096)	3.37	-	3.14
DGA dist. reweighting (N=23, k=262k)	3.35	-	<u>3.08</u>
750M model	MMLU loss	MMLU acc.	avg. Pile loss
750M model Uniform	MMLU loss 3.03	MMLU acc. 26.1 %	avg. Pile loss 2.82
750M model Uniform Importance S. k=4096	MMLU loss 3.03 2.96	MMLU acc. 26.1 % 27.7 %	avg. Pile loss 2.82 2.75
750M model Uniform Importance S. k=4096 Importance S. k=262k	MMLU loss 3.03 2.96 3.13	MMLU acc. 26.1 % 27.7 % 27.0 %	avg. Pile loss 2.82 2.75 2.99
750M model Uniform Importance S. k=4096 Importance S. k=262k DGA domain reweighting (k=64)	MMLU loss 3.03 2.96 3.13 3.01	MMLU acc. 26.1 % 27.7 % 27.0 % 27.0 %	avg. Pile loss 2.82 2.75 2.99 2.77
750M model Uniform Importance S. k=4096 Importance S. k=262k DGA domain reweighting (k=64) DGA dist. reweighting (N=22, k=4096)	MMLU loss 3.03 2.96 3.13 3.01 2.89	MMLU acc. 26.1 % 27.7 % 27.0 % 27.0 % 26.8 %	avg. Pile loss 2.82 2.75 2.99 2.77 2.76
750M model Uniform Importance S. k=4096 Importance S. k=262k DGA domain reweighting (k=64) DGA dist. reweighting (N=22, k=4096) DGA dist. reweighting (N=22, k=262k)	MMLU loss 3.03 2.96 3.13 3.01 <u>2.89</u> 2.93	MMLU acc. 26.1 % 27.7 % 27.0 % 27.0 % 26.8 % 27.0 %	avg. Pile loss 2.82 2.75 2.99 2.77 2.76 2.68
750M model Uniform Importance S. k=4096 Importance S. k=262k DGA domain reweighting (k=64) DGA dist. reweighting (N=22, k=4096) DGA dist. reweighting (N=22, k=262k) DGA dist. reweighting (N=23, k=4096)	MMLU loss 3.03 2.96 3.13 3.01 <u>2.89</u> 2.93 2.88	MMLU acc. 26.1 % 27.7 % 27.0 % 27.0 % 26.8 % 27.0 % 27.0 % 27.4 %	avg. Pile loss 2.82 2.75 2.99 2.77 2.76 2.68 2.75

1080 B.4 WEIGHTS EVOLUTION ON DISTRIBUTIONS

1082 We present the weights assigned to each distribution over time corresponding to the loss on the 1083 specialized task in Figure 14. The top-10 up-weighted distributions include one from the specific 1084 dataset D_{spe} and other academic related domains (e.g. pubmed_central, phil_papers, 1085 dm_mathematics), which are greatly relevant to the sub-topics in MMLU benchmark.



- 1129
- 1130
- 1131 1132
- 1133

1134 C HYPERPARAMETERS

Table 4 provides the model architectures and hyperparameters used in this paper.

Table 4: Architecture hyperparameters for various model scales used in the paper. All models are vanilla Transformer decoder-only models.

	Layers	Attention heads	Embed dim	Hidden dim	Context limit	learning rate
125M	12	12	768	3072	1024	1×10^{-4}
350M	24	16	1024	4096	1024	1×10^{-4}
750M	36	20	1280	5120	1024	1×10^{-4}

D DGA FOR DISTRIBUTION REWEIGHTING

Algorithm 2 explains the distribution reweighting with DGA. The implementation can be easily adapted from domain reweighting DGA with minor modifications.

1153 Algorithm 2 Distribution Reweighting w. DGA. (Difference from domain reweighting are marked in blue)

1.	Input: Constinue domains D , D , IS distributions $A \stackrel{\Delta}{=} [n, n]$ specific set D .
1:	input: Generic domains D_1, \ldots, D_k , i.s. distributions $\mathcal{A}_{dist} = [p_1, \ldots, p_N]$, specific set D_{spe} , input optimizer state ω^0 optimizer function Optimizer such as Adam or SGD initial weights.
	α^0 , outer learning rate η , weight update frequency T_r
2:	Initialize distribution weights: $\alpha_{1,i}^{0} = \alpha_{1,i}^{0}$, i.e. init. domain weights: $\alpha_{1,i}^{0} = \alpha_{1,i}^{0} \otimes \mathcal{A}_{dist}$.
3:	for $t = 0 \dots T$ do
4:	Sample batch from generic mixture: $\boldsymbol{x} \sim \min(\boldsymbol{\alpha}_{domain}^t)$
5:	Update the parameters $\theta^{t+1}, \omega^{t+1} \leftarrow \texttt{Optimizer}(\theta^t, \omega^t, \nabla_{\theta} \ell(\theta^t, x))$
6:	if $t\%T_r = 0$ then
7:	Sample a batch from each <i>distribution</i> : $x_i \sim \min(p_i)$ for $i = 1 \dots N$ and $y \sim D_{spe}$
8:	Compute gradient alignements $m{a}_i^t \leftarrow \langle abla \ell(m{ heta}^{t+1}, m{x}_i), abla \ell'(m{ heta}^{t+1}, m{y}) angle$
9:	Update <i>distribution weights</i> : $\alpha_{\text{dist}}^{t+1} \leftarrow \frac{\hat{\alpha}}{\sum_{i=1}^{k} \hat{\alpha}_i}$ with $\hat{\alpha} = \alpha_{\text{dist}}^t \odot \exp(-\eta a^t)$,
10:	Updated <i>domain weights</i> : $\alpha_{\text{domain}}^{t+1} = \alpha_{\text{dist}}^{t+1} \otimes \mathcal{A}_{dist}$.
11:	else
12:	Do nothing: $oldsymbol{lpha}_{ ext{dist}}^{t+1} \leftarrow oldsymbol{lpha}_{ ext{dist}}^t$
13:	end if
14:	end for
15:	Return Optimized parameters $\theta^{(T)}$ and weights trajectory $\alpha^t, t = 0 \dots T$

¹¹⁸⁸ E PROOF OF THEOREM 2

To ease notations, we define $m(\alpha) = \sum_{i=1}^{k} \alpha_i \mu_i$. Let α^t be the current weight estimate of DGA and θ^t be the estimate of the parameters. The generalist gradient is

$$\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^t, \boldsymbol{\alpha}^t) = \sum_{i=1}^k \alpha_i (\boldsymbol{\theta}^t - \boldsymbol{\mu}_i)$$

1196 Hence, doing a gradient descent step on θ^t with step 1 yields

$$\boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^t - \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^t, \boldsymbol{\alpha}^t)$$
(9)

$$=\sum_{i=1}^{k} \alpha_i \boldsymbol{\mu}_i = m(\boldsymbol{\alpha}) \tag{10}$$

1202 Then, the gradient alignment is

$$a_i = \langle \nabla_{\boldsymbol{\theta}} L_i(\boldsymbol{\theta}^{t+1}), \nabla_{\boldsymbol{\theta}} L_{\text{spe}}(\boldsymbol{\theta}^{t+1})$$
(11)

1204
1205
$$a_{i} = \langle \nabla_{\theta} L_{i}(\boldsymbol{\sigma}^{+}), \nabla_{\theta} L_{\text{spe}}(\boldsymbol{\sigma}^{+}) \rangle$$

$$= \langle m(\boldsymbol{\alpha}^{t}) - \boldsymbol{\mu}_{i}, m(\boldsymbol{\alpha}) - m(\tilde{\boldsymbol{\alpha}}) \rangle]$$
(11)
(12)

1206

$$= \langle m(\boldsymbol{\alpha}^{t}), m(\boldsymbol{\alpha}) - m(\tilde{\boldsymbol{\alpha}}) \rangle - \Sigma_{j=1}^{k} \langle \boldsymbol{\mu}_{i}, \boldsymbol{\mu}_{j} \rangle (\boldsymbol{\alpha}^{t} - \tilde{\boldsymbol{\alpha}})$$
(13)

$$= \langle m(\boldsymbol{\alpha}^t), m(\boldsymbol{\alpha}) - m(\tilde{\boldsymbol{\alpha}}) \rangle - \left[M M^T (\boldsymbol{\alpha}^t - \tilde{\boldsymbol{\alpha}}) \right]$$
(14)

1209 The first part does not depend on i, so it will have no effect in the mirror descent step. Hence, the 1210 mirror descent step is equivalent to a mirror descent step to minimize the function

$$f(oldsymbol{lpha}) = rac{1}{2} \|M(oldsymbol{lpha} - ilde{oldsymbol{lpha}})\|^2$$

The convergence theory of mirror descent yields the bound (Beck & Teboulle, 2003, Theorem 4.2)

$$\min_{t=1...T} f(\boldsymbol{\alpha}^t) = O(\frac{1}{\sqrt{T}})$$

which in turn, thanks to the strong-convexity of f, gives

$$\min_{t=1...T} \|\boldsymbol{\alpha}^t - \tilde{\boldsymbol{\alpha}}\|^2 \le \frac{1}{\lambda_{min}(MM^T)} \min_{t=1...T} f(\boldsymbol{\alpha}^t) = O(\frac{1}{\sqrt{T}})$$

¹²⁴² F COMPARISON BETWEEN DGA AND DOGE

We compare DGA with DoGE (Fan et al., 2024) in the context of task-adaptive training. Specifically, we use the RedPajama-v2 dataset clustered into k=64 domains as the generic dataset, and Arxiv from *The Pile* as the specialized domain.

Following the methodology outlined by Fan et al. (2024), we implement a two-stage pipeline for domain-reweighted pretraining. First, we optimize the domain weights α on a proxy model. Next, we compute the average domain weights over all time steps and use them as fixed sampling weights to train the base model. For this setup, we train proxy models at two scales (31M and 125M parameters) and a base model at 125M parameters. For DGA, we follow Algorithm 1 to train the 125M base model. Unlike DoGE, DGA updates the domain weights dynamically in an online manner, eliminating the need for a separate proxy model or precomputed fixed sampling weights.

1254 1255 We report both the specialized loss on Arxiv over time (measured in GPT hours) and the validation 1256 loss on RedPajama-v2, which reflects the retention of general knowledge from the generic dataset. 1257 To evaluate the performance of DoGE as both an online and offline domain reweighting method, we 1258 present results from both the proxy model and the base model. All the models are trained with 1 \times 1258 Nvidia A100 GPU.

According to Figure 15 (a), DGA greatly outperforms online DoGE (125M proxy model) and achieves comparable specialized loss as the offline DoGE (125M model trained w. 31M/125M Doge weights). However, as an online algorithm, DGA does not require any proxies, which significantly outperform DoGE in terms of efficiency. In addition, as shown in Figure 15 (b), DGA reaches a lower loss on RedPajama than both online and offline DoGE. It indicates that DoGE is more vulnerable from catastrophic forgetting compared to DGA.





1290 1291

¹²⁹⁶ G SUPPLEMENT RESULTS ON LANGUAGE MODELING

1298 G.1 TASK-ADAPTIVE PRETRAINING WITH FINE-TUNING 1299

1300 We further assess the model's performance after fine-tuning on a dataset drawn from the downstream 1301 task. Specifically, we fine-tune 125M pretrained model checkpoints trained using uniform sampling 1302 and DGA reweighting with k=64, on tokens from the Stackexchange subset of *The Pile*. We 1303 chose different numbers of tokens available for fine-tuning. For each number of tokens, we train the 1304 model on those tokens only using a small learning rate (10^{-5}) , and report the best validation loss 1305 across the runs.

As shown in Figure 16, the model pretrained with DGA demonstrates superior performance on the specialized domain (Stackexchange) after fine-tuning, consistently outperforming the uniform sampling and importance sampling baselines across various scales of available fine-tuning tokens.



1350 G.2 SCALING PERFORMANCE OF DGA ACROSS VARIOUS MODEL SCALES

To examine the scaling performance of the DGA algorithm, we train three models of varying scales (125M, 350M, and 750M) using both uniform sampling and DGA with k=64 clusters. As illustrated in Figure 19 (a), DGA consistently outperforms uniform sampling by a significant margin across all three model scales. However, as shown in Figure 19 (b), DGA exhibits some degradation in general knowledge compared to uniform sampling, which presents the highest level of sample diversity. Notably, this performance gap in general knowledge narrows on the largest model (750M), highlighting the potential scalability benefits of DGA in larger model regimes.





1404 G.3 ABLATION ON HYPER-PARAMETERS

Ablations on step size η **and update frequency** T_r . We conducted the ablation experiments on both step size η and frequency T_r used for updating domain weights α , as described in Algorithm 1. According to Figure 18, using a step size that is too small will result in slow updates to the domain weights. On the other hand, applying a large step size (η) can accelerate the learning of domain weights but may also lead to training instability due to overly up-weighted domains. Since EMA can effectively stabilize learning, we recommend that practitioners choose η values between 0.1 and 0.5.

1412Regarding the reweight frequency, using a smaller T_r generally improves the final specialized loss but1413increases computation costs. To balance cost and performance, we recommend setting T_r between141430 and 100. However, these values should be tailored to specific use cases and levels of domain1415granularity.



Figure 18: Grid search on step size η and update frequency T_r

Ablations on EMA factor β . We perform an ablation on the hyper-parameter β applied in the exponential moving average update in Algorithm 1. We choose two extreme values: (1) β =0.1, which is close to 0 and (2) β =0.9, which is close to 1.0. We present the specialized loss on various target domains over time with different values of β . With a strict token limit (30M per domain), both β =0.1, 0.9 can effectively smooth the loss curve while efficiently acquiring specialized knowledge. However, in the cases with 0.1B token per domain or unlimited resources, setting β =0.9 severely slow down the learning of domain weights, which hurts training efficiency.





H CATASTROPHIC FORGETTING ON GENERAL KNOWLEDGE

As a prevalent issue in task-adaptive training, catastrophic forgetting often arises when models lose general knowledge while adapting to specialized tasks. Our experimental results demonstrate that DGA can preserve more general knowledge during task-adaptive pretraining, while other baselines, including importance sampling (Figure 2) and DoGE (Figure 15) are more vulnerable to catastrophic forgetting.

1512	In addition, compared with direct domain rewaighting, our proposed distribution rewaighting along
1513	rithm can rate in the general domain knowledge better without secrificing specialized loss (Figure 12)
1514	runni can retain the general dollarin knowledge better without sacrificing specialized loss (Figure 15).
1515	
1516	
1517	
1512	
1510	
1519	
1520	
1521	
1522	
1523	
1524	
1525	
1526	
1527	
1528	
1529	
1530	
1531	
1532	
1533	
1534	
1535	
1536	
1537	
1538	
1539	
1540	
1541	
1542	
1543	
1544	
1545	
1546	
1547	
1548	
1549	
1550	
1551	
1552	
1553	
1554	
1555	
1556	
1557	
1558	
1559	
1560	
1561	
1562	
1563	
1564	
1565	
1000	