000 DOMAIN2VEC: VECTORIZING DATASETS TO FIND 001 THE OPTIMAL DATA MIXTURE WITHOUT TRAINING 002 003

Anonymous authors

Paper under double-blind review

ABSTRACT

The mixture ratio of data from different source domains significantly affects the performance of language models (LM) pretraining. In this paper, we intro-012 duce DOMAIN2VEC, a novel approach that decomposes any dataset into a linear combination of several "Meta-Domains", a new concept designed to capture key underlying features of datasets. DOMAIN2VEC maintains a vocabulary of Meta-015 Domains and uses a Meta-Domain Classifier to decompose any given dataset into 016 a domain vector that corresponds to a distribution over this vocabulary. These domain vectors enable the identification of optimal data mixture ratio for LM pretraining in a training-free manner under the Distribution Alignment Assumption (DA²), which suggests that when the data distribution of the training set and the validation set is more aligned, a lower validation loss is achieved. Moreover, previous work could use DOMAIN2VEC to model the relationship between domain vectors and LM performance, greatly enhancing the scalability of previous methods without retraining as new datasets are introduced. Extensive experiments demonstrate that DOMAIN2VEC finds data mixture ratios that enhance downstream task performance with minimal computational overhead. Specifically, Do-MAIN2VEC achieves the same validation loss on Pile-CC using only 51.5% of the compute required when training on the original mixture of The Pile Dataset. Under equivalent compute budget, DOMAIN2VEC improves downstream performance by an average of 2.72%. DOMAIN2VEC serves as a strong and efficient baseline for data mixture optimization in LM pretraining, offering insights into improving data efficiency in large-scale models.

031 032

033 034

004

010 011

013

014

017

018

019

021

023

025

026

027

028

029

INTRODUCTION 1

Through training on large-scale text corpora, Large Language Models (LLMs) have demonstrated 035 strong generalization capabilities (Touvron et al., 2023; OpenAI et al., 2024; Yang et al., 2024; DeepSeek-AI et al., 2024). Training datasets for LLMs are typically divided into multiple domains 037 based on their sources. For example, a widely used dataset, The Pile (Gao et al., 2021), includes 12.07% Books3, 8.96% ArXiv, 6.12% FreeLaw, etc. Recent studies have highlighted that mixture proportions of different domains (referred to as data mixture) could significantly impact the effec-040 tiveness of language models (Hoffmann et al., 2022a; Xie et al., 2023b), with data from one domain 041 potentially influencing the outcomes of others (Guo et al., 2022). Typically, the data mixture used for 042 training large language models are determined heuristically or based on downstream performance 043 metrics, which is often unscalable and may lead to suboptimal mixtures. Thus, finding the optimal 044 data mixture in a scalable and efficient manner is a critical research question (Liu et al., 2024).

Recently, researchers have proposed various methods to predict the optimal data mixture. In this 046 paper, we categorize prior work into two lines. The first line implicitly adjusts the data mixture 047 via finding high-quality data from different domains or datasets. Lin et al. (2024) propose using 048 Selective Language Models to select useful tokens to align with the ideal data mixture. Ankner et al. (2024) and Thakkar et al. (2023) directly filter out some low-quality data at the sample level based on the perplexity or the influence score. The second line of work focuses more on modeling 051 the relationship between the data mixture and the performance of language models, which **explicitly** adjusts the data mixture of different domains or datasets. A straightforward method is to train lan-052 guage models on different data mixtures and select the one that yields the best performance, as seen in the training of Gopher (Rae et al., 2022). However, it is impossible to enumerate all possible data 054 mixtures owing to the enormous computation costs. To address this issue, Xie et al. (2023a) pro-055 pose DoReMi, which leverages a well-trained reference model to guide the training of another proxy 056 model using Group DRO (Nemirovski et al., 2009; Sagawa* et al., 2020) over different datasets. The 057 optimized data mixture derived from this process is then used to train a large model. While DoReMi 058 enhances training efficiency for identifying better data mixtures, it still relies heavily on having a well-trained reference model, and it remains difficult to determine what qualifies as a good reference model. To reduce this dependence, Fan et al. (2023) introduce DoGE, which assigns greater 060 weight to a domain based on its contribution to the learning of target domains. Inspired by scaling 061 laws Kaplan et al. (2020); Hoffmann et al. (2022b), to build a functional relationship between data 062 mixture and the performance of language models rather than providing a single data mixture (Xie 063 et al., 2023a; Fan et al., 2023), several works (Ye et al., 2024; Ge et al., 2024; Gu et al., 2024; Que 064 et al., 2024) attempt to fit nonlinear expressions through extensive experiments on smaller proxy 065 models. Gu et al. (2024) also accurately predicted that the pretrained domain loss would first rise 066 and then fall during continue pretraining, and introduced the critical mixture ratio to mitigate catas-067 trophic forgetting on the pretrained domain. Instead of using nonlinear expressions, Liu et al. (2024) 068 propose RegMix, which formulates the search for optimal data mixture as a regression task and fits a regression model to predict the performance of different data mixture. 069

While prior work has shown promising results, they have some issues as follows: 1) Higher Com-071 putational Cost: For instance, although the proxy model used in DoReMi (Xie et al., 2023a) has only 280M parameters, the estimated FLOPs of DoReMi is 3.7×10^{19} . Similarly, RHO-1 (Lin 073 et al., 2024) only calculates loss on certain tokens but still requires the entire sentence to be input 074 into the model. 2) Lack of Scalability : When building the functional relationship like Ye et al. 075 (2024) and Liu et al. (2024), the dimension of the independent variable (i.e., the number of different 076 datasets) is fixed. If we change components of training dataset (i.e., introduce some new datasets, filter some low-quality data), the previously fitted functions cannot be generalized to current datasets. 077 This necessitates resampling new data mixtures, retraining proxy models, and refitting the functions, which severely limits the scalability of these methods. 079

080 To address these issues, we introduce DOMAIN2VEC, a newly introduced concept to capture the 081 underlying features of datasets. DOMAIN2VEC maintains a vocabulary of "Meta-Domains". We 082 hypothesize that any dataset, regardless of its source, be approximated by a linear combination of 083 several Meta-Domains in certain distribution. This distribution over the vocabulary could serve as the vector representation (or domain vector) of the current dataset. To efficiently determine which 084 Meta-Domains comprise a given dataset, we propose utilizing a Meta-Domain Classifier to generate 085 the domain vector and outline a concrete pipeline to build a Meta-Domain Classifier from scratch. For finding the optimal data mixture for language model pretraining, we introduce the Distribution 087 Alignment Assumption (DA²), stating that lower validation loss can be achieved when the domain vector of training datasets aligns with domain vector of the validation datasets. Instead of modeling the relationship between data mixture and language model performance like previous work (Liu et al., 2024; Ye et al., 2024; Que et al., 2024), we focus on modeling the relationship between do-091 main vectors provided by DOMAIN2VEC and the LM performance which significantly enhances the 092 scalability of prior methods. Notably, regardless of changes to the training datasets, DOMAIN2VEC could still provide corresponding domain vector. Moreover, combining different datasets is equivalent to combining their respective domain vectors. This allows us to predict the performance of 094 various data mixtures without the need to retrain proxy models to fit these fictional relationship 095 again, further improving efficiency. 096

⁰⁹⁷ In summary, we highlight our contributions as follows:

098

099

102

103

- 1. We propose DOMAIN2VEC, a novel concept to capture the underlying features of datasets. We also propose viewing datasets as combinations of "Meta-Domains" and propose an efficient pipeline for vectorizing a dataset using a Meta-Domain Classifier.
- 2. We propose *Distribution Alignment Assumption* (DA²) for language model pretraining, a training-free method to identify the optimal data mixture. Additionally, we demonstrate how to integrate DOMAIN2VEC into prior work, which greatly enhances the scalability of prior work without retraining as training datasets changes.
- 3. We validate the effectiveness of DOMAIN2VEC from two aspects: text generation ability and downstream task performance. Experimental results show that our method could accurately predict the performance of different data mixtures without the need for training any

proxy model. Moreover, we identify data mixtures that achieve downstream performance close to DoReMi (Xie et al., 2023a), while using only 0.26% of its computational cost.

2 DOMAIN2VEC

108

110 111

112 113

114

115

116

117

126 127 128

129 130

131

132

133

151

156

157

In this section, we introduce DOMAIN2VEC, an algorithm that decomposes a dataset into a linear combination of various "Meta-Domains". This approach allows us to represent the underlying features of datasets through a normalized vector. We also outline a pipeline for constructing the vocabulary of DOMAIN2VEC and training a Meta-Domain classifier.

118 Key Assumption DOMAIN2VEC maintains a vocabulary, a set of "Meta-Domains". Assume we 119 have *n* Meta-Domains \mathcal{D}_{j}^{*} ($0 \le j < n$), where \mathcal{D}_{j}^{*} is represented as e_{j} , a one-hot vector where the 120 *j*-th element is 1. We hypothesize that, for any given dataset \mathcal{D} , it could be represented as a domain 121 vector v, by linear combination of these Meta-Domains. Specifically,

$$\boldsymbol{v} \approx \sum_{j=0}^{n-1} v_j \cdot \boldsymbol{e}_j,\tag{1}$$

where each element v_j of v represents the projection (weight) of the dataset \mathcal{D} on \mathcal{D}_j^* . Thus, $v = [v_0, v_1, v_2, ..., v_{n-1}]^\top$ can be a representation (distribution) of the dataset \mathcal{D} over the Meta-Domains.

Construct the Vocabulary of DOMAIN2VEC First, we argue that constructed Meta-Domains, which could represent dataset from any source, requires satisfying these following three conditions:

- 1. The original data for constructing Meta-Domains should be as diverse and large as possible.
- 2. The method for constructing Meta-Domains should be computationally efficient.
- 3. There should be distinct differences between different Meta-Domains.

134 We collected data from more than 100 sources across three coarse 135 domains: English, Chinese, and Code. After deduplication, we ob-136 tained around 5.2TB of text data. First, we utilize bge-small-en-137 v1.5 and bge-small-zh-v1.5 (Xiao et al., 2023) to compute embed-138 dings for the English and Chinese data, respectively. Then, we em-139 ploy K-Means (Macqueen, 1967; Arthur & Vassilvitskii, 2006) to cluster these embeddings, resulting in 240 different Meta-Domains 140 for English and Chinese Data. We also demonstrate the relation-141 ship between the number of Meta-Domain and inertia (measuring 142 the distance between each data point and its centroid) in Figure 1. 143 As for the Code data, we directly classified these data based on 144 their programming language categories, ultimately constructing 20 145 Meta-Domains for code, covering mainstream programming lan-146 guages. Finally, we construct 260 unique Meta-Domains. 147



Figure 1: The relationship between the number of Meta-Domains and Inertia.

148 Meta-Domain Classifier In this section, we will introduce how to obtain the normalized domain 149 vector for any given dataset \mathcal{D}_i , which satisfies Equation 1. First, we trained a Meta-Domain Clas-150 sifier based on Qwen2-1.5b-base (Yang et al., 2024). For any given text $text_j \in \mathcal{D}_i$, we have

$$\boldsymbol{p}_j = [p_0, p_1, p_2, ..., p_{n-1}]^\top = \text{Classifier}(text_j)$$
(2)

where p_i represents the probability that $text_j$ belongs to the *i*-th Meta-Domain. For \mathcal{D}_i , we could sample N texts then take the average of domain vector of these samples. Thus, the domain vector v_i of dataset \mathcal{D}_i is, N-1

$$\boldsymbol{v}_i \approx \frac{1}{N} \sum_{j=0}^{N-1} \boldsymbol{p}_j \tag{3}$$

Then, we could use the vector v_i to approximately represent the feature of dataset \mathcal{D}_i from any source. Meanwhile, during the pretraining phase of large language models, we typically have training datasets $\mathcal{D}_{train} = \{\mathcal{D}_1, \mathcal{D}_2, ..., \mathcal{D}_k\}$ from multiple sources. We can convert each of these datasets into domain vectors following Equation 2 and 3. Therefore, \mathcal{D}_{train} can be approximately represented as $V_{train} = [v_1, v_2, ..., v_k]$, where $V_{train} \in \mathbb{R}^{k \times n}$. 162 Training and Evaluation The Meta-Domain Classifier is trained to determine which Meta-163 Domain an arbitrary text from the training set originally belongs to. Thus, we extract 3,000 texts 164 from each meta-domain for training and 500 documents for evaluation. We add a classifier head to Qwen2-1.5B-base (Yang et al., 2024), which has a shape of (hidden size, 260), where 260 equals the 166 number of Meta-Domains. Then, we use the Adam (Kingma & Ba, 2017) optimizer with a learning rate of 2e-5 and train the classifier for 3 epochs via cross entropy loss. After that, we evaluate 167 the performance of the meta-domain classifier on the test set, achieving a classification accuracy of 168 74.73%. Meanwhile, we also sample 1,000 examples from each sub-dataset of The Pile (Gao et al., 2021). Following Equation 3, we obtain domain vectors predicted by the Meta-Domain classifier for 170 each sub-dataset, as shown in Figure 2. It can be seen that the distribution of the Pile's sub-datasets 171 over the meta-domains is very different. This phenomena not only indicates that our classifier could 172 reasonably distinguish some base features from different datasets, but also demonstrates that the 173 various meta-domains have significant semantic differences. 174



183 Figure 2: The Domain Vector of each sub-dataset of The Pile (Gao et al., 2021), where each row corresponds to a sub-dataset and each column corresponds to a Meta-Domain. The higher the pro-185 portion of data belonging to a particular Meta-Domain, the closer the color of the corresponding cell is to blue). We only display the distribution on some English Meta-Domains for clarity. The full picture is shown in Figure 7. 187

FINDING THE OPTIMAL DATA MIXTURE USING DOMAIN2VEC 3

In this section, we will introduce how to find the optimal data mixture using DOMAIN2VEC in a 191 training free manner called "Distribution Alignment Assumption (DA²)". We will also demonstrate 192 how to incorporate our DOMAIN2VEC tools to prior work, which greatly enhance the scalability of 193 previous work as new datasets are introduced¹. 194

3.1 TASK FORMULATION

During the pretraining phase of large language models, we typically collect training datasets $\mathcal{D}_{train} = \{\mathcal{D}_1, \mathcal{D}_2, ..., \mathcal{D}_k\}$ from multiple sources (e.g., ArXiv, Wikipedia). We also pre-define a validation set \mathcal{D}_{valid} , which might be independently and identically distributed with the training dataset or might be unrelated to the training dataset (e.g., data that exceeds the training dataset cutoff date, data with quality but small quantity). Accordingly, the data mixture $\boldsymbol{r} = [r_1, r_2, ..., r_k]^{\top}, 0 \leq 1$ 202 $r_i \leq 1, \sum_{i=1}^k r_i = 1$ specifies the sampling probability distribution over different trainsets. Let the trained language model be denoted as $\hat{\theta}$, and the validation loss of the model be denoted as \mathcal{L}_{θ} . Thus, 204 the optimization objective of finding the optimal data mixture r^* is to improve the performance of language models, such as minimizing the validation loss, as shown in Equation 4. $\mathcal{L}^{\mathcal{D}_{valid}}(r)$ repre-206 sents the validation loss of the language model pretrained via data mixture r.

$$\boldsymbol{r}^* = \arg\min(\min_{\boldsymbol{\rho}} \mathcal{L}^{\mathcal{D}_{valid}}_{\boldsymbol{\theta}}(\boldsymbol{r})) \triangleq \arg\min_{\boldsymbol{\rho}} \mathcal{L}^{\mathcal{D}_{valid}}(\boldsymbol{r}) \tag{4}$$

3.2 **DISTRIBUTION ALIGNMENT ASSUMPTION (DA**²)

Empirically, when the data distribution of the training set \mathcal{D}_{train} and the validation set \mathcal{D}_{valid} is consistent, we could achieve a lower validation loss $\mathcal{L}^{\mathcal{D}_{valid}}$ on the validation set². The most essential

175

176

177

179

181

182

188

189 190

195 196

197

199

200

201

203

205

207 208 209

212

213

214

215

¹The pseudo code of DOMAIN2VEC+DA² and DOMAIN2VEC + RegMix are shown in Appendix A.2. ²We also have provided the detailed description in the Appendix A.1.

question is "How do we model the data distribution of various datasets?". Fortunately, according to Section 2, for the training dataset \mathcal{D}_{train} , we can obtain its vector representation $V_{train} \in \mathbb{R}^{k \times n}$, which models some base features of \mathcal{D}_{train} . Correspondingly, for the validation set \mathcal{D}_{valid} , we also have its vector representation $q_{valid} = [q_0, q_1, q_2, ..., q_{n-1}]^{\top}$. After mixing \mathcal{D}_{train} with data mixture r, the final distribution over Meta-Domains of \mathcal{D}_{train} is given by $V_{train} \cdot r$. Therefore, based on the distribution alignment assumption, Equation 4 can be equivalently written as:

$$\boldsymbol{r}^* = \operatorname*{arg\,min}_{\boldsymbol{r}} \operatorname{Dist}(\boldsymbol{V}_{train} \cdot \boldsymbol{r}, \boldsymbol{q}_{valid}) \tag{5}$$

where $Dist(\cdot, \cdot)$ is a distance function used to measure the similarity between two vectors. In this paper, we use Huber Loss (Huber, 1964; Hastie et al., 2009) to measure the similarity.

3.3 APPLYING DOMAIN2VEC TO PRIOR WORK

As mentioned before, we could directly combine DOMAIN2VEC with prior work, which could address the issue of needing to refit the relationship $\mathcal{L}^{\mathcal{D}_{valid}}(\mathbf{r})$ between the data mixture \mathbf{r} and the validation loss \mathcal{L} as new datasets are introduced. In this section, we will introduce how to integrate DOMAIN2VEC with RegMix (Liu et al., 2024) to search for the optimal data mixture. Following Liu et al. (2024), we train a Linear Regression Model ($\hat{y} = \boldsymbol{\omega}^{\top} \cdot \boldsymbol{x}$) like Equation 6 to fit $\mathcal{L}(\boldsymbol{p})$ on \mathcal{D}_i^* (notated as $\mathcal{L}^{D_i^*}(\boldsymbol{p})$) to build the relationship between the validation loss on \mathcal{D}_i^* and domain vector $\boldsymbol{p} = [p_0, p_1, p_2, ..., p_{n-1}]^{\top}, 0 \le p_i \le 1, \sum_{i=0}^{n-1} p_i = 1$. Formally,

$$\boldsymbol{\omega}_{i}^{*} = \operatorname*{arg\,min}_{\boldsymbol{\omega}} \| \mathcal{L}^{D_{i}^{*}}(\boldsymbol{p}) - \boldsymbol{\omega}^{\top} \cdot \boldsymbol{p} \|$$
(6)

Because any validation set \mathcal{D}_{valid} can also be viewed as a linear combination of multiple Meta-Domains, i.e., $\mathcal{D}_{valid} \approx \sum_{i=0}^{n-1} q_i \cdot \mathcal{D}_i^*$. Meanwhile, the validation loss over different Meta-Domains are additive. Thus, the validation loss on the \mathcal{D}_{valid} of the data mixture p is,

$$\mathcal{L}^{\mathcal{D}_{valid}}(\boldsymbol{p}) = \sum_{i=0}^{n-1} q_i \cdot \mathcal{L}^{D_i^*}(\boldsymbol{p}) = \sum_{i=0}^{n-1} q_i \cdot (\boldsymbol{\omega}_i^*)^\top \cdot \boldsymbol{p}$$
(7)

245After mixing \mathcal{D}_{train} according to the data mixture r, the final dis-246tribution over the Meta-Domain is given by $V_{train} \cdot r$. Therefore,247we can conclude that $p = V_{train} \cdot r$. Substituting $p = V_{train} \cdot r$,248Equation 4 can be equivalently written as follows,

$$\boldsymbol{r}^* = \arg\min_{\boldsymbol{r}} \sum_{i=0}^{n-1} q_i \cdot \mathcal{L}^{\mathcal{D}_i^*}(\boldsymbol{V}_{train} \cdot \boldsymbol{r})$$
(8)

To fit Equation 6 for each Meta-Domain, we sampled 10,500 diverse data mixture from a Dirichlet distribution based on the token distribution of Meta-Domains. Then we used these data mixtures to train different models with 85M parameters on 1B tokens. We used LightGBM to fit Equation 6 for each Meta-Domain. We also reserved data mixtures that were not trained by LightGBM to eval-

Figure 3: The relationship between the number of trained data mixture and the Spearman Correlation Coefficient.

reserved data mixtures that were not trained by LightGBM to evaluate whether fitted equations can accurately predict the validation loss for unseen data mixture.
The Spearman Correlation Coefficient between the actual loss and the predicted loss by LightGBM
is shown in Figure 3. As the trained data mixture increases, the predictions made by LightGBM
become more accurate.

262

264 265

222 223 224

225

226 227

228

236 237

242 243 244

249 250

251

4 DOMAIN2VEC HELPS FIND THE OPTIMAL DATA MIXTURE WITH LESS COMPUTATION, EVEN WITHOUT TRAINING

The motivation for finding the optimal data mixture is to "Enhance the performance of large language models". The performance of large language models can be evaluated from two perspectives:
1) Text generation ability, which refers to the language modeling loss or perplexity on the hold-out validation dataset. 2) Downstream task performance, such as MMLU (Hendrycks et al., 2021) and GSM8K (Cobbe et al., 2021). Therefore, for the text generation ability, we should find the optimal



Figure 4: The validation loss on the EuroParl (The Pile) and Stackexchange (RedPajama) of models trained using data mixture in Table 1. The loss on other validation sets are shown in Appendix A.4.

data mixture to minimize the validation loss. For downstream task performance, we should find the optimal data mixture to maximize the downstream task performance. By deploying DOMAIN2VEC, we could accurately predict the validation loss of any training dataset with different mixture ratios on any validation dataset, even without the need for training some proxy models. Moreover, we used only 0.26% of the computational costs required by DoReMi (Xie et al., 2023a) to find a data mixture with performance comparable to baselines like DoReMi.

289 4.1 MINIMIZE THE VALIDATION LOSS

291 4.1.1 PILOT STUDY

279

280

281 282

283

284

285

286

287 288

290

Using the validation loss of the large language
model as a metric to evaluate its generation capabilities is very straightforward. *However, is there a data mixture that can simultaneously achieve the lowest loss across all validation sets?* Can the *optimal data mixture generalize across models of different model size?* These questions are essential for the study of data mixture of large language

Table 1: The data mixture we used to mix C4 (Raffel et al., 2020) and Knowledge Pile (Fei et al., 2024).

Dataset	Data Mixture					
C4 0	0.2	0.4	0.6	0.8	1.0	
Knowledge Pile 1.0	0.8	0.6	0.4	0.2	0.0	

models. To answer these questions, we first mix C4 (Raffel et al., 2020) and Knowledge Pile (Fei 300 et al., 2024) with different data mixtures as the training set as shown in Tabel 1. We pretrain two 301 Transformer (Vaswani et al., 2017) Decoder-only models with 83M and 1.6B parameters from 302 scratch using a next-token prediction loss. During pretraining, we evaluate the validation loss of 303 models trained with different mixture ratios on 20 subsets of The Pile (Gao et al., 2021) and RedPa-304 jama (Computer, 2023), as shown in Figure 4. We find that, for different validation sets, the ranking 305 of mixture ratios varies significantly. For each validation dataset, we also rank all the data mixtures 306 based on their validation loss and calculate the Spearman and Pearson correlation coefficients of the 307 data mixture ranking between the 83M model and the 1.6B model on various validation sets. The Spearman correlation coefficient is 0.9743, and the Pearson correlation coefficient is 0.9947. Thus, 308 for the same validation set, the data mixture ranking of validation loss on identical validation 309 dataset does not change with the variation in model parameters. This phenomena indicates that 310 we could find the optimal data mixture without the need to train a large model. Based on these find-311 ings, we will demonstrate how DOMAIN2VEC could predict the ranking of different mixture ratios 312 even without training some small proxy models. 313

3143154.1.2 EXPERIMENTAL SETUP

Dataset & Data Mixture Colossal Clean Crawled Corpus (C4) (Raffel et al., 2020) is a colossal and cleaned version of Common Crawl's web crawl corpus. Knowledge Pile (Fei et al., 2024) is a high-quality 735 GB dataset which could significantly improves the performance of large language models in knowledge-related and mathematical reasoning tasks. We mix C4 and Knowledge Pile with different data mixtures as the training set as shown in Tabel 1.

321

Training Setup We pretrained some Transformer (Vaswani et al., 2017) Decoder-only models
 with 83M and 1.6B parameters from scratch using a next-token prediction loss. All the models have a batch size of 1.5M tokens, and the maximum sequence length is 4096. We use the Adam (Kingma



Figure 5: The validation loss on the Pile-CC subset. DOMAIN2VEC achieves the comparable validation loss of Human (The model using original data mixture from The Pile), which only uses almost 51.5% training computational costs of Human. Using the same training cost, DOMAIN2VEC can reduce the validation loss by approximately 4.72% compared to Human.

& Ba, 2017) optimizer with gradient clip of 1.0. The learning rate linearly warms up to a maximum learning rate of 2e-4 over the first 100 steps, then decreases to 2e-5 using a cosine learning rate scheduler with 10,000 steps. The detailed parameters of models we used are shown in the Table 6.

Evaluation Because the optimal mixture ratio varies for different validation datasets, it is impossible to find a data mixture that is optimal for all validation sets. Therefore, we turn to predict the ranking of loss on 20 validation datasets from The Pile (Gao et al., 2021) and RedPajama (Computer, 2023) for the six different mixture ratios shown in Table 1. Then, we evaluate our proposed method using the Spearman correlation coefficient and the Pearson correlation coefficient between the predicted ranking and the actual ranking.

347 348 4.1.3 EXPERIMENTAL RESULTS

First, we present the validation loss
curves for various data mixtures in
Figure 4 and the Appendix A.1. It
can be observed that, on most validation sets, incorporating a certain
amount of Knowledge Pile significantly reduces the model's valida-

Table 2: The results of deploying the DOMAIN2VEC to predict the ranking of different Validation sets.

Metrics	Random	DOMAIN2VEC+DA ²	DOMAIN2VEC+RegMix
Pearson	0.0300	0.5833	0.3881
Spearman	0.0497	0.6657	0.4629

tion loss, even on the C4 validation set from RedPajama. This indicates the high quality of the 356 training data in the Knowledge Pile. Then, we sample 10,000 samples from C4 and Knowledge Pile 357 respectively, and 1,000 samples from each validation set. After that, we apply DOMAIN2VEC to 358 rank the data mixture, as shown in Table 1. As demonstrated in Table 2, the ranking predicted by 359 DOMAIN2VEC exhibits a strong positive correlation with the actual ranking, significantly outper-360 forming random guessing. Interestingly, we find that DOMAIN2VEC + RegMix even predicted that 361 a mixture of 20% Knowledge Pile and 80% C4 could achieve the lowest validation loss on C4 validation set from RedPajama. We hypothesize that this is due to the higher data quality of Knowledge 362 Pile compared to C4, as well as the overlap between these two datasets in certain Meta-Domains. As 363 a result, incorporating a portion of Knowledge Pile into the mixture likely enhances the training of 364 C4. It is also important to note that our method is a *training-free approach*, unlike prior works that rely on training small proxy models to rank data mixtures. Despite this more challenging setup, our 366 method accurately predicts the rankings of different data mixtures. We believe these experimental 367 results could offer valuable insights for the community.

368 369 370

333

334

335

336 337

338

339

340

4.2 MAXIMIZE THE DOWNSTREAM TASK PERFORMANCE.

In this section, we demonstrate how to use DOMAIN2VEC to identify the optimal data mixture for maximizing downstream task performance. A key question is how to model the relationship between data mixture and downstream task performance. Fortunately, Liu et al. (2024) finds that the validation loss on Pile-CC has the highest correlation with the downstream performance across their evaluations. To make a comparison with previous work, we use the same evaluation datasets as Liu et al. (2024). Thus, our task is to find a data mixture that minimizes the validation loss on Pile-CC. Experimental results reveal that DOMAIN2VEC predicts a data mixture with performance comparable to DoReMi (Xie et al., 2023a), while using only 0.26% computational cost.

Table 3: Downstream Task Performance of different models pretrained on different data mixture. 378 Similiar to Liu et al. (2024), Human refers the original data mixture from The Pile. Pile-CC Only 379 refers only training on the Pile-CC subset. The data mixture and estimated flops of DoReMi and 380 RegMix are from Liu et al. (2024). All the data mixture we used are shown in Table 4 and Table 5. 381 The results of 106M models pretrained on 2B tokens are showin in Table 7 owing to the page limit. 382

383	Benchmark	Human	DoReMi	Pile-CC Only	RegMix	$\mathbf{DOMAIN2VEC} + \mathbf{DA}^2$	DOMAIN2VEC + RegMix
384				290M Model Pre	trained on 6B	Tokens	
385	Social IQA	0.364	0.373	0.374	0.371	0.371	0.368
386	HellaSwag	0.295	0.312	0.317	0.315	0.307	0.312
300	PiQA	0.605	0.631	0.639	0.642	0.624	0.633
387	UpenBookQA Lombodo	0.201	0.271	0.271	0.262	0.268	0.200
388	Lambada	0.175	0.208	0.200	0.210	0.182	0.208
500	ABC Facy	0.711	0.082	0.005	0.074	0.670	0.097
389	COPA	0.393	0.410	0.419	0.417	0.420	0.412
200	RACE	0.052	0.000	0.082	0.037	0.027	0.042
390	LogiOA	0.283	0.200	0.200	0.276	0.235	0.292
391	WinoGrande	0.511	0.506	0.509	0.524	0.498	0.504
202	MultiRC	0.507	0.555	0.509	0.545	0.521	0.517
392	Average Performance	0.417	0.432	0.431	0.431	0.421	0.428
393				595M Model Pre	trained on 6R	Tokens	
394	Social IOA	0.378	0.387	0 300	0 304	0.383	0.388
305	HellaSwan	0.378	0.387	0.390	0.394	0.355	0.366
333	PiOA	0.538	0.577	0.580	0.585	0.555	0.500
396	OpenBookOA	0.273	0.050	0.005	0.007	0.288	0.039
207	Lambada	0.255	0.294	0.332	0.310	0.269	0.292
391	SciO	0.233	0.757	0.770	0.510	0.209	0.769
398	ARC Easy	0.439	0.453	0.478	0.481	0.453	0.460
200	COPA	0.642	0.680	0.672	0.663	0.668	0.667
299	RACE	0.289	0.309	0.311	0.311	0.288	0.303
400	LogiOA	0.263	0.268	0.252	0.267	0.263	0.267
404	WinoGrande	0.509	0.515	0.506	0.509	0.512	0.503
401	MultiRC	0.516	0.533	0.522	0.507	0.506	0.527
402	Average Performance	0.442	0.459	0.464	0.465	0.450	0.456
403				1B Model Pretra	ained on 20B	Tokens	
	Social IOA	0.297	0.411	0.406	0.406	0.304	0.401
404	HelloSwog	0.337	0.427	0.400	0.400	0.394	0.410
405	PiOA	0.575	0.427	0.491	0.450	0.684	0.680
	OpenBookOA	0.038	0.004	0.300	0.304	0.299	0.302
406	Lambada	0.301	0.359	0.348	0.353	0.334	0.339
407	SciO	0.802	0.822	0.809	0.828	0.821	0.818
	ARC Easy	0.482	0.508	0.512	0.518	0.500	0.499
408	COPA	0.683	0.692	0.713	0.708	0.678	0.698
409	RACE	0.306	0.319	0.313	0.314	0.305	0.300
	LogiQA	0.259	0.258	0.269	0.272	0.268	0.267
410	WinoGrande	0.513	0.527	0.541	0.512	0.535	0.533
411	MultiRC	0.523	0.504	0.510	0.530	0.529	0.548
	Average Performance	0.464	0.484	0.487	0.489	0.480	0.483
412		 	2.7×10^{19}		2.5×10^{18}	0.66 × 10 ¹⁶	0.66×10^{16}
413	Estimated FLOPs	0	3.7×10^{-1}	0	3.3×10^{-1}	9.00×10 (0.26%)	9.00×10 (0.26%)
110			(10070)		(9.4070)	(0.2070)	(0.2070)

4.2.1 EXPERIMENTAL SETUP

Dataset & Baseline The Pile dataset (Gao et al., 2021) is an 825 GB English text corpus for the pretraining of large language models. Following RegMix (Liu et al., 2024), we also just use the 17 components of The Pile that do not have copyright issues. And we should find the data mixture to achieve lower validation loss on Pile-CC for better downstream task performance. We also compare our approach with various baselines, such as Human (Based on the data size), DoReMi (Xie et al., 2023a), and RegMix (Liu et al., 2024). Pile-CC Only (Just train the model on the Pile-CC sub dataset) is designed for verifing that there is a strong correlation between Pile-CC's validation loss and downstream performance. The data mixture of different baselines are shown in Table 4.

423 424

414

415 416

417

418

419

420

421

422

425 **Training Setup** We pretrained various sizes of Transformer (Vaswani et al., 2017) Decoder-only 426 models from scratch using a next-token prediction loss. The model parameters range from 106M 427 to 1B. Following (Hoffmann et al., 2022b), the computed token number for the different models is 428 20 times the parameter number of current model. All the models have a batch size of 1M tokens, 429 and the maximum sequence length is 4096. We use the Adam Kingma & Ba (2017) optimizer with gradient clip of 1.0. The learning rate linearly warms up to a maximum learning rate of 6e-4 over 430 the first 1,000 steps, then decreases to 0 using a cosine learning rate scheduler at the end of training 431 stage. The detailed parameters of models we used are shown in the Table 6.

432 Evaluation First, we observed the performance on Pile-CC's validation loss on different model 433 sizes as shown in Figure 5. Then we evaluated the performance of different data mixture using 434 following benchmarks: Social IQA (Sap et al., 2019), HellaSwag (Zellers et al., 2019), PiQA (Bisk 435 et al., 2019), OpenBookQA (Mihaylov et al., 2018), Lambada (Paperno et al., 2016), SciQ (Welbl 436 et al., 2017), ARC Easy (Clark et al., 2018), COPA (Gordon et al., 2012), RACE (Lai et al., 2017), LogiQA (Liu et al., 2021), WinoGrande (Sakaguchi et al., 2021), and MultiRC (Khashabi et al., 437 2018). We utilize LM Evaluation Harness (Gao et al., 2024) to evaluate these models and report the 438 average score across 0-shot to 5-shot settings in Table 3. 439

440 441

4.2.2 EXPERIMENTAL RESULTS

442 First, we sample 1,000 samples from each component from The Pile and Pile-CC validation set 443 and use the Meta-Domain Classifier to calculate the domain vector of each dataset. We generate 444 100,000 different data mixture from a Dirichlet distribution based on the token distribution. Using 445 these mixtures, we predict the optimal data mixture by applying Equation 5 (DOMAIN2VEC+DA²) 446 and Equation 7 (DOMAIN2VEC+RegMix). To avoid the over-fitting of language models, each subset 447 of The Pile is trained for at most one epoch. We also apply rejection sampling to eliminate all the 448 unreasonable data mixtures. As a result, the optimal data mixture predicted by DOMAIN2VEC may 449 vary depending on the size of language models.

450 As illustrated in Figure 5, our proposed DOMAIN2VEC + 451 DA² and DOMAIN2VEC + REGMIX could significantly im-452 prove the training efficiency on Pile-CC compared to Hu-453 man (Using the original data mixture from The Pile). Specif-454 ically, DOMAIN2VEC + DA^2 and DOMAIN2VEC + REGMIX 455 require only about 55.38% and 51.50% of the training steps, respectively, to reach the same validation loss as Human. Fur-456 thermore, under equivalent compute budget, DOMAIN2VEC 457 + DA^2 and DOMAIN2VEC + REGMIX reduce the validation 458 loss by approximately 4.04% and 4.64%, respectively, com-459 pared to Human. In Table 3, we report the performance of 460 models trained on data mixtures derived from various base-461 lines across a wide range of downstream tasks. It can be ob-462 served that the Pile-CC Only shows an average accuracy im-463 provement of 4.27% over Human, indicating that training on 464 more tokens from Pile-CC does enhance the downstream task 465 performance of language models. More importantly, our proposed DOMAIN2VEC + DA^2 and DOMAIN2VEC + REG-466 467 MIX, utilizing only about 0.26% of the FLOPs required by DoReMi, could identify data mixtures that achieve perfor-468 mance comparable to DoReMi, RegMix and Pile-CC Only. 469 This demonstrates both the effectiveness and computational 470



Figure 6: We use t-SNE to visualize domain vectors from different sub-datasets from The Pile. This figure indicates that different subdatasets may contain data from the same Meta-Domain, which explains why different datasets can mutually benefit the training of each other.

efficiency of DOMAIN2VEC. At the same time, it is important to note that while achieving compa-471 rable performance, our method does not allocate an excessively high proportion on Pile-CC training 472 dataset as DoReMi and RegMix do. This suggests that different datasets might mutually benefit the 473 training of other datasets. To investigate the cause of this phenomenon, we used t-SNE (Van der 474 Maaten & Hinton, 2008) to visualize the domain vector of each component of The Pile, as shown 475 in Figure 6. This figure reveals that different datasets can contain data belonging to same meta-476 domain, and datasets like Pile-CC, Wikipedia, PhilPapers encompass data from many different 477 meta-domains. This overlap between datasets suggests that the domain vector effectively captures the underlying features of different datasets, explai 478

5 RELATED WORK

479

480

Recently, there has been a amount of research focusing on finding the optimal data mixture, which
could be broadly categorized into two lines. The first line **implicitly adjusts** the data mixture by
down-sampling data from various datasets via finding high-quality data. For instance, Lin et al.
(2024) propose RHO-1, which leverages Selective Language Models to select useful tokens to align
the data mixture with the ideal ratio. Rather than selecting high-quality data at the token level,
Ankner et al. (2024) utilize the perplexity of small reference models to filter out low-quality samples.

486 Additionally, Thakkar et al. (2023) demonstrate that the Influence Score could guide the process of 487 data re-weighting. After that, Thakkar et al. (2023) propose an online data selection method that 488 eliminates the need of any reference model. The second line of research emphasizes modeling the 489 relationship between data mixture and the performance of language models, which explicitly ad-490 justs the data mixture across different datasets. The most straightforward approach is to observe the performance of various data mixtures and then select the optimal one, as demonstrated during the 491 training of Gopher (Rae et al., 2022). However, this approach comes with high training costs, making 492 it challenging to scale for larger models. To address this issue, Xie et al. (2023a) propose DoReMi, 493 which utilizes a small proxy model to re-weight data from different domains, improving the training 494 efficiency of larger models to some extent. However, DoReMi still requires a well-trained reference 495 model beforehand, which introduce additional computational costs, and it is challenging to define 496 what constitutes an ideal reference model. In response, Fan et al. (2023) introduce DoGE, a method 497 that uses a min-max optimization to train a proxy model for obtaining better domain weights. This 498 approach assigns larger weights to domains that either contribute to learning in other domains or are 499 themselves more challenging to learn. Chen et al. (2023) also propose a skills-based framework to 500 dynamically adjust data mixtures during model training. While the aforementioned methods con-501 sider the relationship between data mixture and the performance of language models, they typically provide a single data mixture rather than modeling a functional relationship. Inspired by the scal-502 ing law (Kaplan et al., 2020; Hoffmann et al., 2022b), Ye et al. (2024) propose Data Mixing Laws, 503 which describes this relationship using an exponential form. Similarly, Ge et al. (2024) introduce 504 BiMix, a scaling law that accounts for both compute consumption and the data mixture. Both Que 505 et al. (2024) and Gu et al. (2024) develop scaling laws for continued pretrain, considering the data 506 mixtures between pretrained and continued pretrained datasets. Notably, Gu et al. (2024) accurately 507 predict that the pretrained domain loss would first increase and then decrease during continued pre-508 training, and introduce critical mixture ratios to mitigate catastrophic forgetting in the pretrained 509 domain. More recently, Liu et al. (2024) propose using a Linear Regression Model to fit the valida-510 tion loss of different data mixtures, demonstrating a strong correlation.

511 While prior works have shown promising results, they have some issues as follows: 1) Computa-512 tional Efficiency: For instance, the estimated FLOPs of DoReMi and RegMix is 3.7×10^{19} and 513 3.5×10^{18} . 2) Lack of Scalability: When the components of the training dataset change (i.e., add 514 some new datasets), the previously fitted functions like Ye et al. (2024) and Liu et al. (2024) can-515 not be directly applied to the updated scenario. This limitation arises because the dimension of the 516 independent variable (i.e., the number of different datasets) in these fitted relationships is fixed. As 517 a result, we need to resample different data mixtures, then retrain some proxy models, and perform 518 the fitting again. In this paper, we propose a novel concept DOMAIN2VEC, which decomposes any dataset into a linear combination of several Meta-Domains to capture underlying features of datasets. 519 DOMAIN2VEC shares some ideas with some prior works in the field of Meta-Learning, such as Jo-520 maa et al. (2021) and Chen et al. (2024). These works have explored dataset representation in latent 521 spaces. While sharing the concept of latent space representation for datasets, DOMAIN2VEC dif-522 fers in both purpose and implementation and we focus on language model pretraining data mixture. 523 Then we propose **D**ISTRIBUTION ALIGNMENT ASSUMPTION, a training-free manner to identify 524 the optimal data mixture for language model pretraining. Importantly, using DOMAIN2VEC tools 525 we provided, all fitting experiments are conducted in the dimension of Meta-Domains. When train-526 ing datasets change, we can still map them as linear combinations of several Meta-Domains, which 527 greatly enhance the scalability of prior works (Xie et al., 2023a; Ye et al., 2024; Liu et al., 2024).

528 529 530

6 CONCLUSIONS

531 In this work, we introduce DOMAIN2VEC, a novel concept to capture the underlying features 532 of datasets by decomposing datasets into a linear combination of several "Meta-Domains". We 533 also propose an efficient method to acquire vectorized representation (domain vector) for any 534 given dataset. Based on the domain vector, we introduce a training-free approach by *Distribution* Alignment Assumption (DA^2) for language models pretraining to find the optimal data mixture. 536 By leveraging DOMAIN2VEC, we greatly enhance the scalability of previous methods without re-537 training as training datasets change. Experimental results show that DOMAIN2VEC could use less computation costs to find the data mixture with better text generation ability and downstream task 538 performance. DOMAIN2VEC could serve as a strong and efficient baseline, and we hope that this work will provide some insights into the data mixture optimization for language models pretraining.

540 REFERENCES

Zachary Ankner, Cody Blakeney, Kartik Sreenivasan, Max Marion, Matthew L. Leavitt, and Man sheej Paul. Perplexed by perplexity: Perplexity-based data pruning with small reference models,
 2024. URL https://arxiv.org/abs/2405.20541.

- David Arthur and Sergei Vassilvitskii. k-means++: The advantages of careful seeding. Technical report, Stanford, 2006.
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about
 physical commonsense in natural language. In AAAI Conference on Artificial Intelligence, 2019.
 URL https://api.semanticscholar.org/CorpusID:208290939.
- Jintai Chen, Zhen Lin, Qiyuan Chen, and Jimeng Sun. Cross-table pretraining towards a universal function space for heterogeneous tabular data. *arXiv preprint arXiv:2406.00281*, 2024.
- Mayee F Chen, Nicholas Roberts, Kush Bhatia, Jue WANG, Ce Zhang, Frederic Sala, and Christopher Re. Skill-it! a data-driven skills framework for understanding and training language models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=IoizwOINLf.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. ArXiv, abs/1803.05457, 2018. URL https://api.semanticscholar.org/ CorpusID:3922816.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems, 2021. URL https://arxiv. org/abs/2110.14168.
- Together Computer. Redpajama: An open source recipe to reproduce llama training dataset, 04
 2023. URL https://github.com/togethercomputer/RedPajama-Data.
- 570 DeepSeek-AI, Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Dengr, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, 571 Fangyun Lin, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Hanwei Xu, Hao 572 Yang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian 573 Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jin Chen, Jingyang Yuan, Junjie Qiu, Junxiao Song, Kai 574 Dong, Kaige Gao, Kang Guan, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Liyue 575 Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming 576 Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. 577 Chen, R. L. Jin, Ruiqi Ge, Ruizhe Pan, Runxin Xu, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan 578 Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, 579 Shuiping Yu, Shunfeng Zhou, Size Zheng, T. Wang, Tian Pei, Tian Yuan, Tianyu Sun, W. L. 580 Xiao, Wangding Zeng, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wentao Zhang, X. Q. 581 Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Liu, Xin Xie, Xingkai 582 Yu, Xinnan Song, Xinyi Zhou, Xinyu Yang, Xuan Lu, Xuecheng Su, Y. Wu, Y. K. Li, Y. X. Wei, 583 Y. X. Zhu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui 584 Wang, Yi Zheng, Yichao Zhang, Yiliang Xiong, Yilong Zhao, Ying He, Ying Tang, Yishi Piao, 585 Yixin Dong, Yixuan Tan, Yiyuan Liu, Yongji Wang, Yongqiang Guo, Yuchen Zhu, Yuduan Wang, 586 Yuheng Zou, Yukun Zha, Yunxian Ma, Yuting Yan, Yuxiang You, Yuxuan Liu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhewen Hao, Zhihong Shao, 588 Zhiniu Wen, Zhipeng Xu, Zhongyu Zhang, Zhuoshu Li, Zihan Wang, Zihui Gu, Zilin Li, and 589 Ziwei Xie. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model, 590 2024. URL https://arxiv.org/abs/2405.04434.
- Simin Fan, Matteo Pagliardini, and Martin Jaggi. DOGE: Domain reweighting with generalization
 estimation. In Second Agent Learning in Open-Endedness Workshop, 2023. URL https://openreview.net/forum?id=qiKqsqwYXm.

608

619

627

634

594 Zhaoye Fei, Yunfan Shao, Linyang Li, Zhiyuan Zeng, Conghui He, Hang Yan, Dahua Lin, and 595 Xipeng Qiu. Query of cc: Unearthing large scale domain-specific knowledge from public corpora, 596 2024. URL https://arxiv.org/abs/2401.14624. 597

- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason 598 Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. The pile: An 800gb dataset of diverse text for language modeling. CoRR, abs/2101.00027, 2021. URL 600 https://arxiv.org/abs/2101.00027.
- 602 Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Fos-603 ter, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muen-604 nighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lin-605 tang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework 606 for few-shot language model evaluation, 07 2024. URL https://zenodo.org/records/ 607 12608602.
- Ce Ge, Zhijian Ma, Daoyuan Chen, Yaliang Li, and Bolin Ding. Data mixing made efficient: A 609 bivariate scaling law for language model pretraining, 2024. URL https://arxiv.org/ 610 abs/2405.14908. 611
- 612 Andrew Gordon, Zornitsa Kozareva, and Melissa Roemmele. SemEval-2012 task 7: Choice of plau-613 sible alternatives: An evaluation of commonsense causal reasoning. In Eneko Agirre, Johan Bos, 614 Mona Diab, Suresh Manandhar, Yuval Marton, and Deniz Yuret (eds.), *SEM 2012: The First 615 Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main 616 conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop 617 on Semantic Evaluation (SemEval 2012), pp. 394–398, Montréal, Canada, 7-8 June 2012. Association for Computational Linguistics. URL https://aclanthology.org/S12-1052. 618
- Jiawei Gu, Zacc Yang, Chuanghao Ding, Rui Zhao, and Fei Tan. Cmr scaling law: Predicting 620 critical mixture ratios for continual pre-training of language models, 2024. URL https:// 621 arxiv.org/abs/2407.17467. 622
- 623 Shangmin Guo, Yi Ren, Stefano V Albrecht, and Kenny Smith. Sample relationships through the 624 lens of learning dynamics with label information. In First Workshop on Interpolation Regular-625 izers and Beyond at NeurIPS 2022, 2022. URL https://openreview.net/forum?id= mIl1mMA7Uz. 626
- Trevor Hastie, Robert Tibshirani, Jerome H Friedman, and Jerome H Friedman. The elements of 628 statistical learning: data mining, inference, and prediction, volume 2. Springer, 2009. 629
- 630 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Ja-631 cob Steinhardt. Measuring massive multitask language understanding. In International Confer-632 ence on Learning Representations, 2021. URL https://openreview.net/forum?id= 633 d7KBjmI3GmQ.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza 635 Rutherford, Diego de las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hen-636 nigan, Eric Noland, Katherine Millican, George van den Driessche, Bogdan Damoc, Aurelia 637 Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Oriol Vinyals, Jack William Rae, and Lau-638 rent Sifre. An empirical analysis of compute-optimal large language model training. In Al-639 ice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), Advances in Neural 640 Information Processing Systems, 2022a. URL https://openreview.net/forum?id= 641 iBBCRUlOAPR. 642
- 643 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza 644 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, 645 Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. 646 Training compute-optimal large language models, 2022b. URL https://arxiv.org/abs/ 647 2203.15556.

649

650

651

652

653

661

672

684

685

- Peter J. Huber. Robust Estimation of a Location Parameter. The Annals of Mathematical Statistics, 35(1):73-101, 1964. doi: 10.1214/aoms/1177703732. URL https://doi.org/10.1214/ aoms/1177703732.
- Hadi S Jomaa, Lars Schmidt-Thieme, and Josif Grabocka. Dataset2vec: Learning dataset metafeatures. Data Mining and Knowledge Discovery, 35(3):964–985, 2021.
- 654 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, 655 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language 656 models, 2020. URL https://arxiv.org/abs/2001.08361. 657
- Daniel Khashabi, Snigdha Chaturvedi, Michael Roth, Shyam Upadhyay, and Dan Roth. Look-658 ing beyond the surface: A challenge set for reading comprehension over multiple sentences. 659 In Marilyn Walker, Heng Ji, and Amanda Stent (eds.), Proceedings of the 2018 Conference of 660 the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pp. 252-262, New Orleans, Louisiana, June 662 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1023. URL https: 663 //aclanthology.org/N18-1023. 664
- 665 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. URL https://arxiv.org/abs/1412.6980. 666
- 667 Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. RACE: Large-scale ReAding 668 comprehension dataset from examinations. In Martha Palmer, Rebecca Hwa, and Sebastian Riedel 669 (eds.), Proceedings of the 2017 Conference on Empirical Methods in Natural Language Process-670 ing, pp. 785–794, Copenhagen, Denmark, September 2017. Association for Computational Lin-671 guistics. doi: 10.18653/v1/D17-1082. URL https://aclanthology.org/D17-1082.
- Zhenghao Lin, Zhibin Gou, Yeyun Gong, Xiao Liu, Yelong Shen, Ruochen Xu, Chen Lin, Yujiu 673 Yang, Jian Jiao, Nan Duan, and Weizhu Chen. Rho-1: Not all tokens are what you need, 2024. 674 URL https://arxiv.org/abs/2404.07965. 675
- 676 Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. Logiqa: a chal-677 lenge dataset for machine reading comprehension with logical reasoning. In Proceedings of the 678 Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI'20, 2021. ISBN 9780999241165. 679
- 680 Qian Liu, Xiaosen Zheng, Niklas Muennighoff, Guangtao Zeng, Longxu Dou, Tianyu Pang, Jing 681 Jiang, and Min Lin. Regmix: Data mixture as regression for language model pre-training, 2024. 682 URL https://arxiv.org/abs/2407.01492. 683
 - J Macqueen. Some methods for classification and analysis of multivariate observations. In Proceedings of 5-th Berkeley Symposium on Mathematical Statistics and Probability/University of California Press, 1967.
- 687 Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct 688 electricity? a new dataset for open book question answering. In Conference on Empirical Methods 689 in Natural Language Processing, 2018. URL https://api.semanticscholar.org/ 690 CorpusID: 52183757. 691
- A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to 692 stochastic programming. SIAM Journal on Optimization, 19(4):1574–1609, 2009. doi: 10.1137/ 693 070704277. URL https://doi.org/10.1137/070704277. 694
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Floren-696 cia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red 697 Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brock-699 man, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, 700 Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, 701 Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey

702 Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila 704 Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, 705 Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gib-706 son, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan 708 Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun 710 Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-711 mali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook 712 Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel 713 Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen 714 Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel 715 Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, 716 Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, 717 Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, 718 Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel 719 Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Ra-720 jeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, 721 Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel 722 Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe 723 de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, 724 Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, 725 Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra 726 Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, 727 Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Sel-728 sam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, 729 Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, 730 Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Pre-731 ston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vi-732 jayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan 733 Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, 734 Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Work-735 man, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming 736 Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao 737 Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. URL https://arxiv.org/abs/2303.08774. 739

- Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Ngoc Quan Pham, Raffaella Bernardi,
 Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. The LAMBADA dataset:
 Word prediction requiring a broad discourse context. In Katrin Erk and Noah A. Smith (eds.), *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1525–1534, Berlin, Germany, August 2016. Association for Computational Linguistics. doi: 10.18653/v1/P16-1144. URL https://aclanthology.org/
 P16–1144.
- Haoran Que, Jiaheng Liu, Ge Zhang, Chenchen Zhang, Xingwei Qu, Yinghao Ma, Feiyu Duan,
 Zhiqi Bai, Jiakai Wang, Yuanxing Zhang, Xu Tan, Jie Fu, Wenbo Su, Jiamang Wang, Lin Qu,
 and Bo Zheng. D-cpt law: Domain-specific continual pre-training scaling law for large language
 models, 2024. URL https://arxiv.org/abs/2406.01375.
- Alec Radford. Improving language understanding by generative pre-training. 2018.
- Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John
 Aslanides, Sarah Henderson, Roman Ring, Susannah Young, Eliza Rutherford, Tom Hennigan,
 Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche, Lisa Anne Hendricks,
 Maribeth Rauh, Po-Sen Huang, Amelia Glaese, Johannes Welbl, Sumanth Dathathri, Saffron

756 757 758 759 760 761 762 763 764 765 766 766	Huang, Jonathan Uesato, John Mellor, Irina Higgins, Antonia Creswell, Nat McAleese, Amy Wu, Erich Elsen, Siddhant Jayakumar, Elena Buchatskaya, David Budden, Esme Sutherland, Karen Simonyan, Michela Paganini, Laurent Sifre, Lena Martens, Xiang Lorraine Li, Adhiguna Kun- coro, Aida Nematzadeh, Elena Gribovskaya, Domenic Donato, Angeliki Lazaridou, Arthur Men- sch, Jean-Baptiste Lespiau, Maria Tsimpoukelli, Nikolai Grigorev, Doug Fritz, Thibault Sottiaux, Mantas Pajarskas, Toby Pohlen, Zhitao Gong, Daniel Toyama, Cyprien de Masson d'Autume, Yu- jia Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin, Aidan Clark, Diego de Las Casas, Au- relia Guy, Chris Jones, James Bradbury, Matthew Johnson, Blake Hechtman, Laura Weidinger, Iason Gabriel, William Isaac, Ed Lockhart, Simon Osindero, Laura Rimell, Chris Dyer, Oriol Vinyals, Kareem Ayoub, Jeff Stanway, Lorrayne Bennett, Demis Hassabis, Koray Kavukcuoglu, and Geoffrey Irving. Scaling language models: Methods, analysis & insights from training go- pher, 2022. URL https://arxiv.org/abs/2112.11446.
768 769 770	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>J. Mach. Learn. Res.</i> , 21(1), jan 2020. ISSN 1532-4435.
771 772 773 774	Shiori Sagawa*, Pang Wei Koh*, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust neural networks. In <i>International Conference on Learning Representations</i> , 2020. URL https://openreview.net/forum?id=ryxGuJrFvS.
775 776 777	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: an adver- sarial winograd schema challenge at scale. <i>Commun. ACM</i> , 64(9):99–106, August 2021. ISSN 0001-0782. doi: 10.1145/3474381. URL https://doi.org/10.1145/3474381.
778 779 780 781 782 783 784	Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. Social IQa: Commonsense reasoning about social interactions. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pp. 4463–4473, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1454. URL https://aclanthology.org/D19-1454.
785 786 787 788 788	Megh Thakkar, Tolga Bolukbasi, Sriram Ganapathy, Shikhar Vashishth, Sarath Chandar, and Partha Talukdar. Self-influence guided data reweighting for language model pre-training. In <i>The 2023 Conference on Empirical Methods in Natural Language Processing</i> , 2023. URL https://openreview.net/forum?id=rXn9W04M2p.
789 790 791 792 793 794 795 796 797 798 799 800 801 802	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko- lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023. URL https://arxiv.org/abs/2307.09288.
803 804 805	 Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. Journal of machine learning research, 9(11), 2008. Ashish Vaswani. Noam Shazeer. Niki Parmar, Jakob Hazkoroit, Llion Jones, Aiden N. Correct.
500	Asinshi yaswani, nuani shazeet, miki rannai, jakoo Uszkoren, Liion jones, Aldan N Gomez,

Ashish Vaswahi, Noam Shazeer, Niki Parmar, Jakob Uszkorett, Lhon Jones, Atdah N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf.

- Johannes Welbl, Nelson F. Liu, and Matt Gardner. Crowdsourcing multiple choice science questions. In Leon Derczynski, Wei Xu, Alan Ritter, and Tim Baldwin (eds.), *Proceedings of the 3rd Workshop on Noisy User-generated Text*, pp. 94–106, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi: 10.18653/v1/W17-4413. URL https://aclanthology.org/W17-4413.
- Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighoff. C-pack: Packaged resources to advance general chinese embedding, 2023.
- Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy Liang,
 Quoc V Le, Tengyu Ma, and Adams Wei Yu. Doremi: Optimizing data mixtures speeds up
 language model pretraining. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023a. URL https://openreview.net/forum?id=lXuByUeHhd.
 - Sang Michael Xie, Shibani Santurkar, Tengyu Ma, and Percy Liang. Data selection for language models via importance resampling. In *Thirty-seventh Conference on Neural Information Process-ing Systems*, 2023b. URL https://openreview.net/forum?id=uPSQv0leAu.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Daviheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jin-gren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wen-bin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. Qwen2 technical report, 2024. URL https://arxiv.org/abs/2407.10671.
- Jiasheng Ye, Peiju Liu, Tianxiang Sun, Yunhua Zhou, Jun Zhan, and Xipeng Qiu. Data mixing laws:
 Optimizing data mixtures by predicting language modeling performance, 2024. URL https: //arxiv.org/abs/2403.16952.
 - Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? In *Annual Meeting of the Association for Computational Linguistics*, 2019. URL https://api.semanticscholar.org/CorpusID:159041722.

A Appendix

A.1 DETAILED DESCRIPTION OF THE DISTRIBUTION ALIGNMENT ASSUMPTION

In this section, we will introduce the detailed description of the Distribution Alignment Assumptionfor language model pretaining.

⁸⁷⁰ In the scenario of finding the optimal data mixture for language model pretraining, the validation set ⁸⁷¹ \mathcal{D}_{valid} is fixed, and we should adjust the data mixture to construct the training set \mathcal{D}_{train} to achieve ⁸⁷² lower validation loss calculated by Equation 9, where $\hat{\theta}$ is parameters of a pretrained language ⁸⁷³ model. ⁸⁷⁴

$$\mathbb{E}_{X \sim \mathcal{D}_{valid}} - \log P(X|\hat{\theta}) = \mathbb{E}_{X \sim \mathcal{D}_{valid}} \sum_{i=1}^{|X|} - \log(P(x_i|x_{< i}, \hat{\theta}))$$
(9)

Typically, we pretrain language models via next token prediction (Radford, 2018) like Equation 10.

$$\hat{\theta} = \arg \max \mathbb{E}_{X \sim \mathcal{D}_{train}} \log P(X|\theta)$$

)

864

866

867

875 876 877

878

879

882 883

890 891

892 893

894

895

896

897

899

900

916

 $= \underset{\theta}{\arg\max} \mathbb{E}_{X \sim \mathcal{D}_{train}} \sum_{i=1}^{|X|} \log(P(x_i | x_{\leq i}, \theta))$ (10)

That is, we need to find a $\hat{\theta}$ that maximizes the expected probability of $X \sim D_{train}$, which is also known as Maximum Likelihood Estimation (MLE). When the data distributions of \mathcal{D}_{train} and \mathcal{D}_{valid} are aligned, the optimization target of language models pretraining process equals find a $\hat{\theta}$ that maximizes the expected probability of $X \sim \mathcal{D}_{valid}$. Therefore, we introduce the Distribution Alignment Assumption for language model pretaining, a novel method to find the optimal data mixture without training. After that, we propose to use the Meta-Domain Classifier to capture some underlying features of datasets which could helps modeling the data distribution of different datasets.

A.2 Algorithm

In Algorithm 1, we show the pseudo code for acquiring the domain vector for pretraining datasets.

In Algorithm 2 and 3, we show the pseudo code for how to use DOMAIN2VEC to find the optimal data mixture, including Distribution Alignment Assumption, and applying DOMAIN2VEC to RegMix (Liu et al., 2024).

Algorithm 1 DOMAIN2VEC

Require: Training Datasets $\mathcal{D}_{train} = {\mathcal{D}_1, \mathcal{D}_2, ..., \mathcal{D}_k}$, Validation Dataset \mathcal{D}_{valid} , Meta-Domain Classifier Classifier

901 1: 902 2: Domain Vectors $V_{train} = []$ 903 3: **for** i = 1 to k **do** 904 Sample N data points from \mathcal{D}_i 4: $v_i = \frac{1}{N} \sum_{j=0}^{N-1} \text{Classifier}(text_j)$, where $text_j \in \mathcal{D}_i$ \triangleright Get Domain Vectors of \mathcal{D}_{train} $V_{train} = [V_{train}, v_i]$ 905 5: 906 6: 907 7: **end for** 908 8: 909 9: Sample N data points from \mathcal{D}_{valid} 10: $q_{valid} = \frac{1}{N} \sum_{j=0}^{N-1} \text{Classifier}(text_j)$, where $text_j \in \mathcal{D}_{valid} \triangleright \text{Get Domain Vector of } \mathcal{D}_{valid}$ 910 911 11: 912 12: return $V_{train} = [v_1, v_2, ..., v_k], q_{valid}$ 913 914

915 A.3 DATA MIXTURE OF DIFFERENT METHODS

In this section, we will show the data mixture on The Pile (Gao et al., 2021) of different methods we used in this paper for reproduction. In Table 4, we show the optimal data mixture predicted by

1:		0	
2: Sa	ample K candidates data mixture r_i from Diric	$\operatorname{chlet}(\boldsymbol{a}_{train})$	
3:			. 1
4: 1	ne Optimal Data Mixture $r^* = r_1$	⊳ Iniu	lanze the optimal data n
5. 6. fo	$\mathbf{r} \mathbf{i} = 2$ to k do		
7:	if $\text{Dist}(V_{train} \cdot r, q_{valid}) < \text{Dist}(V_{train} \cdot r)$	*, \boldsymbol{q}_{valid}) then	▷ Updata the optim
m	ixture	, icana)	1 1
8:	$m{r}^*=m{r}_i$		
9:	end if		
10: e i	nd for		
11:	· · · · · · · · · · · · · · · · · · ·		
12: r e	eturn the optimal data mixture r		
Algor Requi da L 1:	ithm 3 DOMAIN2VEC+RegMix ire: Domain Vectors of Training Datasets V_{tro} ation Dataset q_{valid} , Token Distribution of Training inear Regression Model for Each Meta-Domain	$egin{aligned} & \mathbf{z}_{ain} = [oldsymbol{v}_1, oldsymbol{v}_2,, \mathbf{v}_{t}] \ & \mathbf{z}_{t} = [oldsymbol{v}_1, oldsymbol{v}_2,, oldsymbol{v}_{t}] \ & \mathbf{z}_{t} = [oldsymbol{v}_1, oldsym$	$[v_k]$, Domain Vectors o $v_{train} = [\alpha_1, \alpha_2,, \alpha_k],$
Algor Requir da L: 1: 2: Sa 3:	ithm 3 DOMAIN2VEC+RegMix ire: Domain Vectors of Training Datasets V_{tro} ation Dataset q_{valid} , Token Distribution of Trainine inear Regression Model for Each Meta-Domain ample K candidates data mixture r_i from Diric	$egin{aligned} & \mathbf{z}_{ain} = [oldsymbol{v}_1, oldsymbol{v}_2,, \mathbf{v}_{tining}] \ & ext{Datasets} \ oldsymbol{a}_t \ & ext{main} \ \mathcal{L}^{\mathcal{D}^*_i}(oldsymbol{p}). \ & ext{chlet}(oldsymbol{a}_{train}) \end{aligned}$	$[v_k]$, Domain Vectors o $v_{train} = [\alpha_1, \alpha_2,, \alpha_k],$
Algor Requi da L: 1: 2: Sa 3: 4: T	ithm 3 DOMAIN2VEC+RegMix ire: Domain Vectors of Training Datasets V_{tra} ation Dataset q_{valid} , Token Distribution of Traininear Regression Model for Each Meta-Domain ample K candidates data mixture r_i from Dirice he Optimal Data Mixture $r_* = r_1$	$a_{iin} = [\boldsymbol{v}_1, \boldsymbol{v}_2,, \mathbf{v}_{train}]$ ining Datasets \boldsymbol{a}_t in $\mathcal{L}^{\mathcal{D}^*_i}(\boldsymbol{p}).$ chlet (\boldsymbol{a}_{train}) \triangleright Initi	$[v_k]$, Domain Vectors o $v_{train} = [\alpha_1, \alpha_2,, \alpha_k]$, value the optimal data m
Algor Requi da L: 1: 2: Sa 3: 4: T 5: L	ithm 3 DOMAIN2VEC+RegMix ire: Domain Vectors of Training Datasets V_{trc} ation Dataset q_{valid} , Token Distribution of Traininear Regression Model for Each Meta-Domain ample K candidates data mixture r_i from Dirich the Optimal Data Mixture $r_* = r_1$ $C(r^*) = \sum_{i=0}^{n-1} q_i \cdot \mathcal{L}_{i}^{\mathcal{D}_i^*}(V_{train} \cdot r_1)$	$a_{iin} = [\boldsymbol{v}_1, \boldsymbol{v}_2,, \mathbf{v}_{t_i}]$ ining Datasets \boldsymbol{a}_t in $\mathcal{L}^{\mathcal{D}^*_i}(\boldsymbol{p}).$ chlet (\boldsymbol{a}_{train}) arphi Initi	$[v_k]$, Domain Vectors o $v_{rain} = [\alpha_1, \alpha_2,, \alpha_k]$,
Algor Requi da L: 1: 2: 3: 4: 5: 6:	ithm 3 DOMAIN2VEC+RegMix ire: Domain Vectors of Training Datasets V_{tra} ation Dataset q_{valid} , Token Distribution of Training inear Regression Model for Each Meta-Domain ample K candidates data mixture r_i from Dirich the Optimal Data Mixture $r_* = r_1$ $f(r^*) = \sum_{i=0}^{n-1} q_i \cdot \mathcal{L}^{\mathcal{D}_i^*}(V_{train} \cdot r_1)$	$a_{in} = [\boldsymbol{v}_1, \boldsymbol{v}_2,, \mathbf{v}_t]$ ining Datasets \boldsymbol{a}_t n $\mathcal{L}^{\mathcal{D}^*_i}(\boldsymbol{p}).$ chlet (\boldsymbol{a}_{train}) \triangleright Initi	$[v_k]$, Domain Vectors o $v_{rain} = [\alpha_1, \alpha_2,, \alpha_k]$,
Algor Requi da L: 1: 2: 3: 4: T. 5: 6: 7: for	ithm 3 DOMAIN2VEC+RegMix ire: Domain Vectors of Training Datasets V_{tra} ation Dataset q_{valid} , Token Distribution of Training inear Regression Model for Each Meta-Domain ample K candidates data mixture r_i from Diric the Optimal Data Mixture $r_* = r_1$ $E(r^*) = \sum_{i=0}^{n-1} q_i \cdot \mathcal{L}^{\mathcal{D}_i^*}(V_{train} \cdot r_1)$ or $i = 2$ to k do	$a_{in} = [v_1, v_2,, v_{t_1}, v_{t_2},, v_{t_t}]$ ining Datasets a_t in $\mathcal{L}^{\mathcal{D}^*_i}(p).$ chlet (a_{train}) arphi Initi	$[v_k]$, Domain Vectors o $v_{rain} = [\alpha_1, \alpha_2,, \alpha_k]$,
Algor Requi da L: 1: 2: 3: 4: T 5: C: 7: 6: 7: 6: 7: 6: 8:	ithm 3 DOMAIN2VEC+RegMix ire: Domain Vectors of Training Datasets V_{trc} ation Dataset q_{valid} , Token Distribution of Traininear Regression Model for Each Meta-Domain ample K candidates data mixture r_i from Dirich the Optimal Data Mixture $r_* = r_1$ $f(r^*) = \sum_{i=0}^{n-1} q_i \cdot \mathcal{L}^{\mathcal{D}_i^*}(V_{train} \cdot r_1)$ or $i = 2$ to k do if $\mathcal{L}(r_i) < \mathcal{L}(r^*)$ then	$a_{iin} = [v_1, v_2,, v_{ining}]$ ining Datasets a_t in $\mathcal{L}^{\mathcal{D}_i^*}(p)$. chlet (a_{train}) \triangleright Initi	$[v_k]$, Domain Vectors o $v_{rain} = [\alpha_1, \alpha_2,, \alpha_k]$, alize the optimal data m pdata the optimal data m
Algor Requi da L: 1: 2: 3: 4: T 5: 6: 7: 6: 9: 10:	ithm 3 DOMAIN2VEC+RegMix ire: Domain Vectors of Training Datasets V_{tra} ation Dataset q_{valid} , Token Distribution of Training inear Regression Model for Each Meta-Domain ample K candidates data mixture r_i from Dirice the Optimal Data Mixture $r_* = r_1$ $\mathcal{L}(r^*) = \sum_{i=0}^{n-1} q_i \cdot \mathcal{L}^{\mathcal{D}_i^*}(V_{train} \cdot r_1)$ or $i = 2$ to k do if $\mathcal{L}(r_i) < \mathcal{L}(r^*)$ then $r^* = r_i$ $\mathcal{L}(r^*) = \mathcal{L}(r_i)$	$a_{in} = [\boldsymbol{v}_1, \boldsymbol{v}_2,, \mathbf{v}_{t_1}, \mathbf{v}_{t_2},, \mathbf{v}_{t_t}]$ in $\mathcal{L}^{\mathcal{D}_i^*}(\boldsymbol{p}).$ chlet (\boldsymbol{a}_{train}) \triangleright Initi	$[v_k]$, Domain Vectors o $v_{train} = [\alpha_1, \alpha_2,, \alpha_k]$, falize the optimal data modulate the optimal data m
Algor Requi da L: 1: 2: Sa 3: 4: T 5: \mathcal{L} 6: 7: 7: fo 9: 10: 11: 11:	ithm 3 DOMAIN2VEC+RegMix ire: Domain Vectors of Training Datasets V_{tra} ation Dataset q_{valid} , Token Distribution of Training inear Regression Model for Each Meta-Domain ample K candidates data mixture r_i from Dirice the Optimal Data Mixture $r_* = r_1$ $\mathcal{L}(r^*) = \sum_{i=0}^{n-1} q_i \cdot \mathcal{L}^{\mathcal{D}_i^*}(V_{train} \cdot r_1)$ or $i = 2$ to k do if $\mathcal{L}(r_i) < \mathcal{L}(r^*)$ then $r^* = r_i$ $\mathcal{L}(r^*) = \mathcal{L}(r_i)$ end if	$\mathbf{a}_{in} = [\mathbf{v}_1, \mathbf{v}_2,, \mathbf{u}_t]$ ining Datasets \mathbf{a}_t in $\mathcal{L}^{\mathcal{D}_i^*}(\mathbf{p}).$ chlet (\mathbf{a}_{train}) \triangleright Initi	$[v_k]$, Domain Vectors o $v_{rain} = [\alpha_1, \alpha_2,, \alpha_k]$, falize the optimal data modulate
Algor Requi da L: 2: 3: 4: T: 5: 6: 7: 6: 9: 10: 11: 12: et	ithm 3 DOMAIN2VEC+RegMix ire: Domain Vectors of Training Datasets V_{tra} ation Dataset q_{valid} , Token Distribution of Training inear Regression Model for Each Meta-Domain ample K candidates data mixture r_i from Dirice the Optimal Data Mixture $r_* = r_1$ $E(r^*) = \sum_{i=0}^{n-1} q_i \cdot \mathcal{L}^{\mathcal{D}_i^*}(V_{train} \cdot r_1)$ or $i = 2$ to k do if $\mathcal{L}(r_i) < \mathcal{L}(r^*)$ then $r^* = r_i$ $\mathcal{L}(r^*) = \mathcal{L}(r_i)$ end if and for	$a_{in} = [v_1, v_2,, v_n]$ ining Datasets a_t in $\mathcal{L}^{\mathcal{D}_i^*}(p)$. chlet (a_{train}) \triangleright Initi	$[v_k]$, Domain Vectors o $v_{rain} = [\alpha_1, \alpha_2,, \alpha_k]$, falize the optimal data modulate the optimal data m
Algor Requi da L: 2: 3: 4: T. 5: 6: 7: 6: 9: 10: 11: 12: et 13:	ithm 3 DOMAIN2VEC+RegMix ire: Domain Vectors of Training Datasets V_{tra} ation Dataset q_{valid} , Token Distribution of Training inear Regression Model for Each Meta-Domain ample K candidates data mixture r_i from Diric the Optimal Data Mixture $r_* = r_1$ $r_i(r^*) = \sum_{i=0}^{n-1} q_i \cdot \mathcal{L}^{\mathcal{D}_i^*}(V_{train} \cdot r_1)$ or $i = 2$ to k do if $\mathcal{L}(r_i) < \mathcal{L}(r^*)$ then $r^* = r_i$ $\mathcal{L}(r^*) = \mathcal{L}(r_i)$ end if nd for	$a_{in} = [v_1, v_2,, v_n]$ ining Datasets a_t in $\mathcal{L}^{\mathcal{D}_i^*}(p)$. chlet (a_{train}) \triangleright Initi	$[v_k]$, Domain Vectors o $v_{rain} = [\alpha_1, \alpha_2,, \alpha_k]$, falize the optimal data n pdata the optimal data n

A.4 EXPERIMENTAL RESULTS OF PILOT STUDY

mixture predicted may change as model sizes change.

960 In this section, we report the validation loss on various datasets Arxiv, C4, Book3, PG19 from Red-961 Pajama (Computer, 2023), and BookCorpus2, DM Mathematics, Enron Emails, FreeLaw, Hack-962 erNews, NIH ExPorter, OpenSubtitles, OpenWebText2, PhilPapers, PubMed Abstracts, PubMed 963 Central, USPTO Backgrounds, Ubuntu IRC, Youtube Subtitles from The Pile (Gao et al., 2021) in 964 Figure 4, Figure 9 and Figure 8. According to the experimental results, we find that 1) for different 965 validation sets, the ranking of mixture ratios varies significantly. 2) for the same validation set, 966 the data mixture ranking of validation loss on identical validation dataset does not change with the variation in model parameters. We hope our experimental results and findings could provide 967 some insights to the community about efficiently finding the optimal data mixture. 968

- 969
- 970

Table 4: The data mixture of The Pile (Gao et al., 2021) from different baselines, which aligns with the data mixture used in Liu et al. (2024).

975	75						
976		Data Mixture	Human	DoReMi	Pile-CC Only	RegMix	
977		ArXiv	0.134	0.004	0.0	0.001	
978		FreeLaw	0.049	0.005	0.0	0.001	
070		NIH ExPorter	0.007	0.008	0.0	0.001	
979		PubMed Central	0.136	0.006	0.0	0.003	
980		Wikipedia (en)	0.117	0.086	0.0	0.016	
981		DM Mathematics	0.025	0.002	0.0	0.0	
082		Github	0.054	0.022	0.0	0.0	
502		PhilPapers	0.003	0.034	0.0	0.0	
983		Stack Exchange	0.118	0.019	0.0	0.0	
984		Enron Emails	0.004	0.009	0.0	0.002	
085		Gutenberg (PG-19)	0.025	0.009	0.0	0.002	
505		Pile-CC	0.142	0.743	1.0	0.87	
986		Ubuntu IRC	0.009	0.011	0.0	0.064	
987		EuroParl	0.005	0.008	0.0	0.0	
088		HackerNews	0.01	0.016	0.0	0.012	
300		PubMed Abstracts	0.107	0.014	0.0	0.024	
989		USPTO Backgrounds	0.053	0.004	0.0	0.002	

Table 5: The optimal data mixture predicted by $DOMAIN2VEC + DA^2$ and DOMAIN2VEC + Reg-Mix. To avoid the over-fitting problem, any subset of The Pile (Gao et al., 2021) will be trained at most one epoch. And we adopt rejection sampling to filter out certain unreasonable data mixtures. Thus, the data mixture predicted may change as model sizes change.

Data Mixture	DOMAIN2VEC+DA ²			Do	DOMAIN2VEC+RegMix			
	106M	290M	595M	1B	106M	290M	595M	1B
ArXiv	0.0131	0.0131	0.0389	0.0431	0.0152	0.0070	0.0114	0.010
FreeLaw	0.0076	0.0076	0.0316	0.0305	0.0395	0.0267	0.0339	0.026
NIH ExPorter	0.0008	0.0008	0.0028	0.0023	0.0000	0.0199	0.0000	0.000
PubMed Central	0.0773	0.0773	0.0519	0.0704	0.0343	0.0576	0.0099	0.051
Wikipedia (en)	0.2970	0.2970	0.2049	0.2126	0.0847	0.0101	0.1014	0.257
DM Mathematics	0.0003	0.0003	0.0056	0.0026	0.0177	0.0018	0.0011	0.000
Github	0.0096	0.0096	0.0290	0.0298	0.0034	0.0538	0.0500	0.013
PhilPapers	0.0018	0.0018	0.0093	0.0025	0.0118	0.0005	0.0333	0.040
Stack Exchange	0.0464	0.0464	0.0661	0.0585	0.0698	0.0430	0.1199	0.026
Enron Emails	0.0000	0.0000	0.0009	0.0000	0.0018	0.0000	0.0000	0.000
Gutenberg (PG-19)	0.0217	0.0217	0.0484	0.0370	0.0467	0.0223	0.0007	0.025
Pile-CC	0.4338	0.4338	0.3191	0.3814	0.5370	0.6323	0.5546	0.470
Ubuntu IRC	0.0022	0.0022	0.0063	0.0072	0.1019	0.0123	0.0161	0.006
EuroParl	0.0003	0.0003	0.0042	0.0040	0.0070	0.0037	0.0116	0.000
HackerNews	0.0154	0.0154	0.0521	0.0199	0.0028	0.0551	0.0170	0.067
PubMed Abstracts	0.0596	0.0596	0.0739	0.0532	0.0259	0.0102	0.0190	0.001
USPTO Backgrounds	0.0130	0.0130	0.0549	0.0449	0.0004	0.0438	0.0201	0.001

Table 6: The parameters of different models we used in Section 4.1 and Section 4.2. When calculat-ing the model parameters, we do not take into account the embedding layer and the language model head layer.

Parameter	Text Generation		Downstream Task			
	83M	1.6B	106M	290M	595M	1B
Hidden Size	768	2048	768	1280	1536	2048
FFN Hidden Size	2048	5504	2048	3392	4096	5440
Num of Layers	12	24	15	15	21	21
Num of Heads	12	16	12	10	12	32
Max Seq Length	4096	4096	4096	4096	4096	4096
Vocab Size	128256	128256	151936	151936	151936	151930
RoPE Base	10000	10000	10000	10000	10000	10000

1020	Table 7: Downstream Task Performance of different data mixture on 106M Model. Similar to Liu
1027	et al. (2024). Human refers the original data mixture from The Pile. Pile-CC Only refers only
1028	training on the Pile-CC subset. The data mixture and estimated flops of DoReMi and RegMix are
1029	from Liu et al. (2024).

Benchmark	Human	DoReMi	Pile-CC Only	RegMix	$DOMAIN2VEC + DA^2$	DOMAIN2VEC + RegMix
			106M Model Pr	etrained on 2B	Tokens	
Social IQA	0.340	0.349	0.353	0.356	0.339	0.342
HellaSwag	0.268	0.268	0.269	0.269	0.267	0.264
PiQA	0.573	0.584	0.580	0.586	0.579	0.583
OpenBookQA	0.245	0.251	0.249	0.242	0.245	0.249
Lambada	0.065	0.099	0.102	0.091	0.091	0.090
SciQ	0.550	0.520	0.509	0.537	0.549	0.518
ARC Easy	0.329	0.339	0.335	0.337	0.334	0.331
COPA	0.525	0.570	0.572	0.585	0.578	0.557
RACE	0.236	0.254	0.246	0.251	0.240	0.244
LogiQA	0.282	0.280	0.271	0.274	0.268	0.286
WinoGrande	0.516	0.516	0.502	0.508	0.506	0.499
MultiRC	0.539	0.520	0.515	0.533	0.541	0.544
Average Performance	0.372	0.379	0.375	0.381	0.378	0.376
Estimated EL OPs	0	3.7×10^{19}	0	3.5×10^{18}	9.66×10^{16}	9.66×10^{16}
Estimated FLOFS	0	(100%)	0	(9.46%)	(0.26%)	(0.26%)
Part C. Public Control Antion			10 C 10 C 10 C			100 C 100
FreiLav Stellenarge		1. S. L.				
Miled Raditation Senders (96-19) Wildpedie Ionio Wildpedies					2 - 1 - 1 - 1 - 1 - 1 - 1 - 1 - 1 - 1 -	
Userla IIC Eurofini Hackithan					1.1	
Wit before free frame the tar						and the set of the set

Figure 7: The Domain Vector of each sub-dataset of The Pile (Gao et al., 2021), where each row corresponds to a sub-dataset and each column corresponds to a Meta-Domain. The higher the proportion of data belonging to a particular Meta-Domain, the closer the color of the corresponding cell is to blue). Additionally, since The Pile primarily consists of English texts, we only display the distribution on English Meta-Domains for clarity.



Figure 8: The validation loss on different dataset of models trained using data mixture in Table 1.



Figure 9: The validation loss on different dataset of models trained using data mixture in Table 1.