# DATA MIXING CAN INDUCE PHASE TRANSITIONS IN KNOWLEDGE ACQUISITION

Anonymous authors

Paper under double-blind review

### ABSTRACT

Large Language Models (LLMs) are typically trained on data mixtures: most data come from web scrapes, while a small portion is curated from high-quality sources with dense domain-specific knowledge. In this paper, we show that when training LLMs on such data mixtures, knowledge acquisition from knowledgedense datasets does not always follow a smooth scaling law but can exhibit phase transitions with respect to the mixing ratio and model size. First, through controlled experiments on a synthetic biography dataset mixed with web-scraped data, we demonstrate that: (1) as we increase the model size to a critical value, the model suddenly transitions from memorizing very few to most of the biographies; (2) below a critical mixing ratio, the model memorizes almost nothing even with extensive training, but beyond this threshold, it rapidly memorizes more biographies. We then adopt an information-theoretic perspective to understand and characterize the existence and value of the thresholds. Based on these insights, we identify two mitigation strategies that improve the efficiency of knowledge acquisition from knowledge-dense datasets, and validate their effectiveness on both synthetic and real-world Wikipedia datasets.

026 027 028

029

004

010 011

012

013

014

015

016

017

018

019

021

024

025

## 1 INTRODUCTION

Large Language Models (LLMs) are often trained on two types of datasets. The first is a large-scale corpus scraped from the web (Raffel et al., 2020; Penedo et al., 2024; Li et al., 2024), often spanning billions to trillions of tokens across diverse topics and styles. Due to the scale, it is inherently hard to ensure the information density of the dataset and its relevance to downstream tasks. Hence, a second type of data, smaller-scale datasets curated from high-quality sources, is incorporated. This type of data features very dense knowledge on a specific task or domain. Some of these datasets, such as Wikipedia and Stack Exchange, contain a wealth of world knowledge while others, such as OpenWebMath (Paster et al., 2024) and StarCoder (Li et al., 2023; Kocetkov et al., 2022), contribute domain expertise in mathematics and coding.

As the second-type data usually take a small fraction of the entire corpus, one may wonder: *How much knowledge can an LLM acquire from these knowledge-dense datasets?* Answering this question
 is crucial to understanding the ultimate performance that a model can achieve as we continue to
 improve data quality and scale up training compute.

One crucial aspect is the *model size*, since models cannot store more knowledge than their capacity. Allen-Zhu & Li (2024a) quantified the influence of model size on factual knowledge acquisition in a controlled setting, where the pre-training data only contains a synthetically generated biography dataset and no other data. They found that the amount of knowledge in a sufficiently trained model stores scales linearly with the model size. Similar scaling was observed for memorizing Wikidata fact triples by Lu et al. (2024), and theoretically analyzed by Nichani et al. (2025).

In this paper, we present a quantitative study on the amount of knowledge that a model can acquire from a knowledge-dense dataset under *data mixing*: the knowledge-dense dataset constitutes only a fraction of the pre-training data, denoted as the *mixing ratio r*, and the rest of the data is a large-scale corpus of web text. We show that knowledge acquisition from knowledge-dense datasets, when mixed with the web text, no longer follows a linear scaling law but instead exhibits a more intricate behavior with notable phase transitions with respect to mixing ratios and model sizes. More specifically, we study factual knowledge acquisition. We follow the approach of Allen-Zhu & Li (2024a) to curate a synthetic dataset of biographies with uniform data format and content, which enables us to quantify how much knowledge the model has stored simply by counting the number of memorized biographies. We pre-train models of different sizes on a mixture of this biography dataset and FineWeb-Edu (Penedo et al., 2024), a large-scale web corpus derived from Common Crawl.

Of course, setting r closer to 1 will make the model memorize more biographies, but it also hurts the model's general capabilities that are supposed to be learned from the more diverse web corpus. Therefore, the essence of our study is to understand whether models can still memorize a decent number of biographies for relatively small r (< 40%). Our experiments reveal two interesting findings (Section 3):

Finding 1: Phase Transition in Model Size (Figure 1). When varying the model size while keeping the mixing ratio r fixed, the number of biographies memorized by the model does not scale linearly with its size but instead exhibits a phase transition behavior. When the model size is smaller than a critical model size  $N_0$ , the number of memorized biographies can be nearly zero, and only until the model size reaches  $N_0$ , the model *suddenly* memorizes most of the biographies. The threshold  $N_0$  is higher for smaller mixing ratio r.

Finding 2: Phase Transition in Mixing Ratio (Figure 2). When varying the mixing ratio r while keeping the model size fixed, we find that below a critical mixing ratio  $r_0$ , the model memorizes almost nothing even after significantly longer training, during which each biography appears tens of times more (Figures 3(a) and 4), but beyond  $r_0$ , the number of memorized biographies grows rapidly with r. We further find that as we gradually decrease r, the number of steps needed to memorize a fixed number of biographies is initially growing linearly with 1/r (Figure 3(b)), but soon becomes exponential and even superexponential (Figure 3(c)), making it impossible or practically infeasible for the model to memorize a nontrivial number of biographies despite extensive training passes.

The above findings reveal a caution for practitioners that the mixing ratio should be set with care for
 the model: mixing in knowledge-dense datasets with small mixing ratios can be not beneficial at all,
 especially when training small LMs.

Theoretical Analysis (Section 4). We further present an information-theoretic explanation for the observed phase transitions. We show that these behaviors are not unique to LLM pre-training but can also arise in any learning algorithm that optimally minimizes the overall test loss with a bounded model capacity, which we refer to as *optimal bounded-capacity learner*. By assuming that the optimal test loss follows a power law in model size, we show how phase transition depends on the model size, mixing ratio, and the exposure frequency of each fact in a knowledge-dense dataset. We also derive a power law relationship between the threshold frequency of a fact and the model size.

Empirically, we estimate the threshold frequency of factual knowledge in PopQA (Mallen et al., 2023)
 across 32 open-source models, ranging from 1B to 70B parameters (Section 5). Our experiment demonstrates that our predicted power-law scaling approximately holds for popular LLMs.

**Mitigation Strategies.** Inspired by our theory, we further propose two mitigation strategies to improve the efficiency of knowledge acquisition by increasing exposure frequency:

- 1. Perhaps counter-intuitively, it is beneficial to randomly subsample the knowledge-dense dataset to a *smaller* size, which increases the exposure frequency of each biography to improve memorization (Section 6.2).
- 2. We propose rephrasing the knowledge in a more compact form and add rephrased data to the original dataset while keeping the overall mixing ratio fixed. This increases the exposure frequency of each biography, though represented in different forms. Perhaps surprisingly, LLMs can successfully memorize the knowledge and answer questions in natural language (Section 6.3). We call this method *Compact Knowledge Mixing* (CKM).

We validate on both our synthetic biographies and real-world Wikipedia biographies that our mitigation strategies significantly increase the number of memorized biographies while preserving models' general capability.

105 2 EXPERIMENTAL SETUP

092

094

095

096

097

098

099

100

106

**The SynBio Dataset.** We follow the approach of Allen-Zhu & Li (2024b) to curate a synthetic biography dataset with uniform data format and content. Specifically, each individual is characterized



Figure 1: Phase transition in model size. For each mixing ratio, as model size increases, accuracy initially remains zero. Once model size surpasses some threshold, accuracy rapidly grows to over 60%.

Figure 2: Phase transition in mixing ratio. For each model size, as mixing ratio r increases, accuracy initially remains zero. Only when r exceeds some threshold does accuracy quickly improve.

119 by five attributes: birth date, birth city, university, major, and employer. For each individual, the value 120 of each attribute is randomly and independently sampled from a predefined domain. These (name, 121 attribute, value) triplets are then converted into natural text descriptions using predefined sentence 122 templates. For instance, the triplet (Gracie Tessa Howell, birth city, St. Louis, MO) is converted 123 into the sentence: "Gracie Tessa Howell's birthplace is St. Louis, MO." Following (Allen-Zhu & Li, 124 2024b), each time the model encounters a biography, the five sentences is randomly shuffled, and a 125 new sentence template is selected for each attribute from a set of five possible templates. We denote the dataset containing N biographies as SynBio-N. See Appendix C.2.1 for full details. 126

127 **Evaluation.** Denote a knowledge triplet (name, attribute, value) as (n, a, v) and let |v| represent 128 the number of tokens in the value v. For evaluation, the model is prompted with the prefix of the 129 templated sentence containing n and a and is asked to generate |v| tokens using greedy decoding. 130 The triplet is successfully memorized if the generated text exactly matches v. For example, to test 131 whether the model has learned the triplet (Gracie Tessa Howell, birth city, St. Louis, MO), the prompt "Gracie Tessa Howell's birthplace is" is provided. The model is deemed to have memorized the fact if 132 it generates "St. Louis, MO." We report the accuracy averaged over all individuals, attributes, and 133 templates in the main text and defer the detailed results to Appendix B.2. 134

135 **Training Setup.** Our experiments use the GPT-NeoX library (Andonian et al., 2023) and the 136 Pythia model architecture (Biderman et al., 2023), with model sizes ranging from 14M to 1B and 137 a sequence length of 2048. The default setup involves pre-training from scratch on a mixture of FineWeb-Edu and SynBio, using a batch size of 512 and the Warmup-Stable-Decay (WSD) learning 138 rate schedule (Hu et al., 2024) with a peak learning rate of  $10^{-3}$ . We also investigate the continual 139 pre-training setup, which mimics data annealing phase where high-quality data are upweighted to 140 improve model performance (Dubey et al., 2024; Blakeney et al., 2024; Feng et al., 2024; OLMo 141 et al., 2025). Full details are provided in Appendix C. 142

# 3 PHASE TRANSITIONS OF KNOWLEDGE ACQUISITION WITHIN DATA MIXTURES

143

115

116

117

118

# 3.1 PHASE TRANSITION IN MODEL SIZE

148 We first investigate how knowledge acquisition within data mixtures is affected by model size at fixed 149 mixing ratios. For each  $r \in \{0.1, 0.2, 0.3, 0.4\}$ , we train models with sizes ranging from 14M to 150 410M on the mixture of FineWeb-Edu and SynBio-320k for 32B tokens, which is approximately four times the optimal computation for 410M models predicted by the Chinchilla scaling law (Hoffmann 151 et al., 2022). As shown in Figure 1, as model size increases, accuracy on SynBio initially remains 152 near zero. Once the model size surpasses some threshold, accuracy rapidly grows to above 60%. The 153 transition is consistently sharp across different mixing ratios while larger r leads to a smaller critical 154 point. 155

156 3.2 PHASE TRANSITION IN MIXING RATIO

We now study how knowledge acquisition in the data mixing scenario is affected by mixing ratios for fixed model sizes.

Accuracy on the knowledge dataset undergoes a phase transition as mixing ratio increases.

We begin by training models of the same size with different mixing ratios r. Specifically, we train 70M models on the mixture of FineWeb-Edu and SynBio-320K, varying r from 0.1 to 0.45 (stepsize



168 169

(a) Train until accuracy achieves 60% accuracy or until a total of 256B tokens are passed.

(b) Relationship between required training steps to achieve target ac- to attain 40% accuracy against 1/r. curacy and 1/r.

(c) Fitting required training steps

171 Figure 3: Training efficiency declines sharply as r decreases. We train 70M models on the mixture of FineWeb-172 Edu and SynBio-128k. (a) For each r, we train until accuracy achieves 60% accuracy or until a total of 256B tokens are passed. Notably, accuracy for r = 0.2 remains near zero even after extensive training up to 512B 173 tokens. (b) Required training steps to achieve target accuracy initially grows linearly with 1/r, but then escalate 174 rapidly. (c) Required training steps increase exponentially or even superexponentially with 1/r. 175

176

177 0.05), and 410M models on the mixture of FineWeb-Edu and SynBio-1.28M, varying r from 0.1 178 to 0.4 (stepsize 0.1). All models are trained for a total of 32B tokens. As shown in Figure 2(a), 179 for 70M models, as r increases from 0.1 to 0.25, its accuracy on SynBio remains near zero. Only when r exceeds 0.3 does the accuracy begin to steadily improve with increasing r. In Figure 2(b), the 181 accuracy for 410M models exhibit similar trends where it remains near zero for  $r \le 0.3$  and suddenly 182 attains 80% when r grows to 0.4.

183 Training longer barely helps for low mixing ratios. Having 184 identified that models struggle to memorize facts when mix-185 ing ratio falls below some threshold, we investigate whether extended training can mitigate this issue. We extend the train-187 ing horizon for r = 0.2 to 512B tokens for the 70M model 188 and 128B for the 410M model, increasing the model's expo-189 sure to the knowledge dataset by 16 and 4 times, respectively. 190 However, as shown in Figures 3(a) and 4, the accuracy on Syn-191 Bio remains near zero for both model sizes, even with these substantial extensions. 192

193 Required training steps to achieve target accuracy ini-194 tially grows linearly with 1/r, but then escalate rapidly. 195 We further examine how mixing ratio affects the train-196 ing efficiency in knowledge acquisition by measuring the total number of training steps required to reach a target 197 accuracy, denoted as T, across different mixing ratios r. Specifically, we train 70M models with mixing ratios from 199



Figure 4: For 410M models trained on the mixture of FineWeb-Edu and SynBio-1.28M, accuracy for r = 0.2remains near zero even when we extend the training by 4 times.

 $\{0.2, 0.25, 0.3, 0.4, 0.45, 0.5, 0.55, 0.6, 0.7, 0.8\}$ . For each mixing ratio, we evaluate 20 training 200 horizons, approximately evenly spaced on a logarithmic scale with a factor of 1.2, ranging from 0 to 201 256B tokens. For r > 0.4, where accuracy improves rapidly, we assess additional training horizons 202 for more precise estimations. Training continues until the model reaches 60% accuracy or exhausts 203 256B tokens. As shown in Figures 3(a) and 3(b), as we decrease r from 0.8, T initially increase 204 linearly with 1/r for r > 0.4 and quickly deviates from the linear trend as r falls below 0.4. 205

Quantifying the growth rate: required training steps increase exponentially or even superexpo-206 **nentially with** 1/r. To quantify the growth rate of T with respect to 1/r, we analyze the results 207 for target accuracy 40% and fit a scaling law for T against 1/r. Specifically, we use an exponential 208 function to fit T against 1/r for all  $r \ge 0.3$ , where 40% accuracy is achieved within 256B tokens. 209 Additionally, we use a power-law function to fit T against 1/r for  $r \in \{0.3, 0.4, 0.45, 0.5, 0.55\}$ . 210 Details on the fitting process can be found in Appendix C.3. To examine whether T follows the fitted 211 trend as r decreases further, we extend the training for r = 0.2 and 0.25 to 660B and 1024B tokens, 212 respectively. However, neither configurations attain 40% accuracy, even after such extended training. 213 As shown in Figure 3(c), the actual T is more than 2.9 times the power-law prediction for r = 0.25and more than 1.9 times for r = 0.2. Notably, the actual T for r = 0.25 is even more than twice the 214 exponential prediction. These significant deviations suggest exponential or even superexponential 215 growth of T with respect to 1/r as r decreases.



Figure 5: Ablation studies on hyperparameters. The models exhibit consistent trends in knowledge acquisition across different batch sizes, learning rate values and schedules. All experiments are conducted by training 70M models on the mixture of FineWeb-Edu and SynBio-320k.

### 3.3 ABLATION STUDIES

226

227 228

229

We now conduct ablation studies to demonstrate the robustness of our findings with respect to hyperparameters. We explore  $r \in \{0.2, 0.4, 0.8\}$  and train 70M models for a total of 64B, 32B, and 16B tokens, respectively, ensuring each configuration passes SynBio the same number of times.

233 **Consistent Trends Across Different Batch Sizes.** As shown in Figure 5(a), we evaluate three 234 batch sizes,  $B \in \{256, 512, 1024\}$ , for each r and observe consistent general trends across all batch 235 sizes. For r = 0.4 and r = 0.8, smaller batch sizes yield slightly higher accuracies, likely due to 236 the increased number of update steps. These experiments further distinguish between two types of 237 frequency at which the model encounters the knowledge dataset: per-token frequency and per-step 238 frequency. For a fixed mixing ratio, doubling the batch size doubles the occurrences of each biography 239 per step, while the occurrences per token remain unchanged. The results demonstrate that per-token frequency, rather than per-step frequency, determines training efficiency in knowledge acquisition. 240

Consistent trends across learning rate values and schedules. In Figure 5(b), we explore peak learning rates among  $\{2.5 \times 10^{-4}, 10^{-3}, 4 \times 10^{-3}\}$  using the WSD scheduler. We observe that the trends are consistent across these values, although the learning process slows down at the lowest value  $2.5 \times 10^{-4}$ . In Figure 5(c), results for both cosine and WSD schedulers show similar trends.

# 2452464 THEORETICAL ANALYSIS

In this section, we take an information-theoretic perspective and point out that the observed phenomena in the experiments are not unique to the current LLM architectures or optimization methods but can also happen for any model that is trained to maximally utilize its capacity to minimize the next-token prediction loss. We then formulate the notion of *optimal bounded-capacity learners* and show how they exhibit similar phase transitions as LLMs. In Sections 5 and 6, we will further study this type of model as a proxy for the real-world LLMs and conclude with several implications that indeed apply to LLMs.

254 4.1 PROBLEM FORMULATION

255 **Data distribution.** The essence of language modeling is to model the distribution of the next 256 token y for a given input context x consisting of all previous tokens. We take a Bayesian view and 257 assume that there is a latent variable  $\theta \in \Theta$  that determines the data distribution of (x, y), denoted 258 as  $(x, y) \sim \mathcal{D}_{\theta}$ . Conceptually,  $\theta$  should contain a lot of knowledge to be presented in the data. For 259 example, in our current universe someone may be born in 1996, but in a parallel universe, the same 260 person may be born in 1999. Or, in a different universe, popular Python libraries may have a very different set of functions. We assume that the universe first draws  $\theta$  from a prior distribution  $\mathcal{P}$  before 261 we observe the data distribution  $\mathcal{D}_{\theta}$ . 262

**Learning Algorithm.** A learning algorithm  $\mathcal{A}$  is a procedure that takes samples from a data distribution  $\mathcal{D}$  of (x, y) and outputs a predictor  $h = \mathcal{A}(\mathcal{D})$  in the end, where h is a function that maps x to a distribution over y. For a given predictor h, we measure its performance by the expected cross-entropy loss  $\mathcal{L}(h; \mathcal{D}) := \mathbb{E}_{(x,y)\sim\mathcal{D}}[-\log p(y \mid h, x)]$ , where  $p(y \mid h, x)$  denotes the predicted distribution of y given x by the predictor h, and log is in base 2 for convenience. We measure the performance of a learning algorithm  $\mathcal{A}$  by its expected loss over all data distributions  $\mathcal{D}_{\theta}$  with respect to the prior  $\mathcal{P}$ :

$$\bar{\mathcal{L}}_{\mathcal{P}}(\mathcal{A}) := \mathbb{E}_{\theta \sim \mathcal{P}}[\mathcal{L}(\mathcal{A}(\mathcal{D}_{\theta}); \mathcal{D}_{\theta})].$$
(1)

In practice, a predictor h can be a transformer, and A can be the pre-training algorithm.

Model Capacity and Mutual Information. We measure the "effective" model capacity—the 272 amount of information a model, produced by learning algorithm  $\mathcal{A}$ , stores about the data distribution 273  $\mathcal{D}_{\theta}$ —by the mutual information (MI) between the model and the data distribution  $\mathcal{D}_{\theta}$ , namely 274  $I(\mathcal{A}(\mathcal{D}_{\theta}); \mathcal{D}_{\theta})$ . Practical learning algorithms usually have a bounded capacity. Consider a learning 275 algorithm  $\mathcal{A}$  that always outputs a model h with at most N parameters, where each parameter is a 276 floating-point number with p bits of precision. Then information theory states that the MI is bounded 277 by pN, *i.e.*,  $I(\mathcal{A}(\mathcal{D}_{\theta}); \mathcal{D}_{\theta}) \leq pN$ . Empirically, Allen-Zhu & Li (2024a) measured the knowledge 278 capacity scaling laws of LLMs in a carefully controlled setting and found that  $I(\mathcal{A}(\mathcal{D}_{\theta}); \mathcal{D}_{\theta}) \approx 2N$ , 279 which hold consistently across various training setups.

**Optimal Bounded-Capacity Learner.** Now, imagine that we train a model under a capacity constraint  $I(\mathcal{A}(\mathcal{D}_{\theta}); \mathcal{D}_{\theta}) \leq M$ , where M is a constant. We optimize the training procedure with huge efforts and put massive computational resources into training so that the resulting model can minimize the loss as much as possible. How does the resulting model behave? This motivates us to define the following notion of optimal bounded-capacity learner.

**Definition 4.1.** For a given prior  $\mathcal{P}$  and M > 0, the best achievable loss under a capacity constraint M is defined as

290

291 292  $F_{\mathcal{P}}(M) := \inf_{\mathcal{A}} \left\{ \bar{\mathcal{L}}_{\mathcal{P}}(\mathcal{A}) : I(\mathcal{A}(\mathcal{D}_{\theta}); \mathcal{D}_{\theta}) \le M \right\},\tag{2}$ 

where the infimum is taken over all learning algorithms. An optimal M-bounded-capacity learner is a learning algorithm  $\mathcal{A}$  such that  $I(\mathcal{A}(\mathcal{D}_{\theta}); \mathcal{D}_{\theta}) \leq M$  and  $\overline{\mathcal{L}}_{\mathcal{P}}(\mathcal{A}) = F_{\mathcal{P}}(M)$ .

# 4.2 WARMUP: MIXTURE OF FACTS

We start with a simple case where the data distribution  $\mathcal{D}_{\theta}$  consists of K random facts. More formally, let  $X_1, \ldots, X_K$  be K disjoint sets of input contexts and  $y_1, \ldots, y_K$  be K target tokens. Each  $(X_i, y_i)$ represents a fact. The universe samples  $y_1, y_2, \ldots, y_K$  independently, where  $y_i$  is drawn from a fixed distribution  $\mathcal{Y}_i$ . Then the universe sets the latent variable  $\theta$  to be  $(y_1, y_2, \ldots, y_K)$  and  $\mathcal{D}_{\theta}(y \mid x_i)$  to be a point mass at  $y_i$  for all  $x_i \in X_i$ . There could be other inputs x that can occur in  $\mathcal{D}_{\theta}$ , but the target tokens of such inputs are independent of  $\theta$ .

Define the exposure frequency of each random fact as the total probability that an input  $x \in X_i$ occurs in  $\mathcal{D}_{\theta}$ . Despite that the K facts have different entropies, the following theorem shows that if their exposure frequencies are the same, then the optimal bounded-capacity learner reduces the expected loss linearly with the capacity M, thus no phase transition in capacity.

**Theorem 4.2.** For all 
$$M \ge 0$$
, if all the facts have the same exposure frequency  $p$ , then  
 $F_{\mathcal{P}}(M) = C + p \cdot \max\{H_{\text{tot}} - M, 0\},$ 
(3)  
where  $H_{\text{tot}} := \sum_{i=1}^{K} H(\mathcal{Y}_i)$  and  $C := F_{\mathcal{P}}(\infty)$ .

### 306 307 308

# 4.3 DATA MIXING INDUCES PHASE TRANSITIONS

What if we mix the random facts with another domain of data, say web text? More specifically, 310 imagine that the data distribution  $\mathcal{D}_{\theta}$  consists of two data domains. The first is the same as before, a 311 mixture of K random facts. The second is another data domain that can have a much more complex 312 structure. We assume that the latent variable  $\theta$  is structured as  $\theta = (\theta_1, \theta_2)$ , where  $\theta_1$  contains 313 information about the K random facts and  $\theta_2$  determines the target token distributions in the second 314 domain. When sampling  $\theta \sim \mathcal{P}$ , the universe draw  $\theta_1$  and  $\theta_2$  independently from priors  $\mathcal{P}_1$  and 315  $\mathcal{P}_2$ , respectively. The data distribution  $\mathcal{D}_{\theta}$  is then mixed together as  $\mathcal{D}_{\theta} = (1-r)\mathcal{D}_{\theta_1}^{(1)} + r\mathcal{D}_{\theta_2}^{(2)}$ , 316 where  $r \in (0, 1)$  is the mixing ratio, and  $\mathcal{D}_{\theta_1}^{(1)}$  and  $\mathcal{D}_{\theta_2}^{(2)}$  are the data distributions of the two domains. Same as before, we use p to denote the exposure frequency of any fact in the first domain, and 317 318  $H_{\text{tot}} := \sum_{i=1}^{K} H(\mathcal{Y}_i)$  to denote the total entropy of the target tokens in the first domain. 319

To understand how a model performs on the first domain after training with an algorithm  $\mathcal{A}$  on the data mixture, we define  $\overline{\mathcal{L}}_1(\mathcal{A}) := \mathbb{E}_{\theta \sim \mathcal{P}_1}[\mathcal{L}(\mathcal{A}(\mathcal{D}_\theta); \mathcal{D}_{\theta_1})]$ , which is the expected loss of  $\mathcal{A}$  on the first domain after learning from the data mixture. If  $\overline{\mathcal{L}}_1(\mathcal{A}) = F_{\mathcal{P}_1}(0)$ , then the learner's performance is the same as random guessing without seeing any data, namely the learner does not learn the facts at all. If  $\overline{\mathcal{L}}_1(\mathcal{A}) = F_{\mathcal{P}_1}(\infty)$ , the learner learns the facts perfectly. 324 The following theorem shows that the optimal bounded-capacity learner no longer learns the facts 325 linearly, but instead exhibits a phase transition in the capacity M. More specifically, the learner 326 sharply transitions between the two extremes as the model size increases. Two key functions for 327 characterizing the transition are

$$M_0^-(x) := \sup\{M \ge 0 : -F'_{\mathcal{P}_2}(M) > x\},\$$
  
$$M_0^+(x) := \inf\{M \ge 0 : -F'_{\mathcal{P}_2}(M) < x\},\$$

Due to convexity of  $F_{\mathcal{P}_2}(M)$ ,  $-F'_{\mathcal{P}_2}(M)$  is non-increasing. Imagine that we vary M from 0 to  $\infty$ . Then  $M_0^-(x)$  and  $M_0^+(x)$  can be interpreted as the last M such that  $-F'_{\mathcal{P}_2}(M)$  is larger than x and 332 the first M such that  $-F'_{\mathcal{P}_2}(M)$  is smaller than x, respectively. If  $F'_{\mathcal{P}_2}(M)$  is strictly decreasing, then  $M_0^-(x)$  and  $M_0^+(x)$  are the same.

**Theorem 4.3** (Phase Transition in Model Size). For any optimal M-bounded-capacity learner A,

1. if 
$$M \leq M_0^-(\frac{r}{1-r} \cdot p)$$
, then  $\bar{\mathcal{L}}_1(\mathcal{A}) = F_{\mathcal{P}_1}(0)$ ;

328

330

331

333

334 335

336 337

2. if  $M \geq M_0^+(\frac{r}{1-r} \cdot p) + H_{\text{tot}}$ , then  $\bar{\mathcal{L}}_1(\mathcal{A}) = F_{\mathcal{P}_1}(\infty)$ .

Key Example: When Web Data Loss Follows a Power Law in Model Size. Consider the case 341 where  $F_{\mathcal{P}_2}(M)$  is a power-law function of M, *i.e.*,  $F_{\mathcal{P}_2}(M) = C + A \cdot M^{-\alpha}$ . Here,  $\alpha \in (0,1)$  and 342 A is a large constant. This is a reasonable assumption since LLM pre-training usually exhibits such a 343 power-law scaling behavior in model size, as observed by many works (Kaplan et al., 2020; Hoffmann 344 et al., 2022). For example, Hoffmann et al. (2022) estimated that the power-law exponent  $\alpha$  is around 345 0.34 for their Chinchilla models. If the second domain is a large-scale dataset scraped from the web, 346 then we should expect that the best achievable loss on this domain is a power-law function of the 347 model capacity. In this case, taking the derivative of  $F_{\mathcal{P}_2}(M)$ , we have  $-F'_{\mathcal{P}_2}(M) = A \cdot \alpha \cdot M^{-\alpha-1}$ . Then,  $M_0^-(x) = M_0^+(x) = (\frac{A\alpha}{x})^{1/(\alpha+1)}$ . Plugging this into Theorem 4.3, we can see that the 348 349 transition point is around 350

$$M \sim \left(\frac{A}{rp}\right)^{1/(\alpha+1)}.$$
(4)

354 This implies that even if the model has the capacity to learn the first domain, it may still need to be 355 very large to acquire any knowledge from it, especially when r or p is small.

Arranging the terms in (4), we can also obtain the transition point in the mixing ratio r:

358 359 360

364

365

366

371

372

351 352

353

356

357

 $r \sim \frac{A}{p \cdot M^{\alpha + 1}},$ (5)

361 which aligns with the empirical observation that r has to be larger than a critical mixing ratio for the 362 model to learn any of the facts.

**Threshold Frequency.** At the fact level, its overall probability of being sampled is rp in the data mixture. Again arranging the terms in (5), we can further predict that for a fact to be learned by the model, its frequency of appearing in the pretraining corpus should be larger than a threshold frequency, which scales with the model size following a power law:

$$rp \sim \frac{A}{M^{\alpha+1}}.$$
 (6)

#### 5 POWER-LAW RELATIONSHIP OF THRESHOLD FREQUENCY AND MODEL SIZE

373 To validate our theoretical prediction of a power-law relationship between models size and threshold 374 frequency, we examine the threshold frequency of a set of knowledge extracted from Wikipedia. 375 Specifically, we evaluate models on PopQA (Mallen et al., 2023), which contains 14k QA pairs derived from Wikidata triplets, along with monthly page view for corresponding Wikipedia articles. 376 Since knowledge tested in PopQA can be structured as triplets, we consider them as homogeneous 377 and expect them to exhibit similar threshold frequencies for a given model size.



Figure 6: Model size and threshold frequency exhibit a power-law relationship, which aligns with our theoretical prediction.

389 **Estimating the Threshold Frequency.** Directly counting the frequency in the pre-training data can 390 be challenging due to the sheer data volume (Kandpal et al., 2023). To address this, we follow Mallen 391 et al. (2023) and use Wikipedia page views as a proxy for popularity, which is expected to be 392 roughly proportional to how frequently the knowledge appears in web data. To estimate the threshold 393 popularity  $P_{\rm thres}$ , we determine the smallest popularity P such that the model's accuracy on data 394 with popularity above P meets the target accuracy  $\alpha_{\text{target}}$ . In our experiments, we set  $\alpha_{\text{target}} = 60\%$ . See Appendix C.4 for details.

396 Threshold frequency and model size follow a power law. We first examine base models from 397 three families: Llama-2 (Touvron et al., 2023), Qwen-2.5 (Qwen et al., 2024) and Gemma-2 (Team et al., 2024). According to their technical reports, models from the same family are likely trained 399 on the same data mixture. As shown in Figure 6(a),  $\log P_{\rm thres}$  generally decreases linearly as log 400 model size increases within each family. The slope may vary across model families, as the exponent 401 for model size in the loss scaling law can differ depending on model architecture and training data. Next, we relax the constraint of training on the same data mixture and investigate the overall trend 402 between model size and  $P_{\rm thres}$ . We add the Llama-3 (Dubey et al., 2024) family, and evaluate both 403 base and instruction-tuned models for all families, totaling 30 models. Interestingly, in Figure 6(b), 404 log model size and  $\log P_{\rm thres}$  also exhibit a linear relationship, with most models falling within the 405 95% confidence interval. We further use models from the OLMo (Groeneveld et al., 2024) family as 406 a validation set, where predictions of the fitted power law closely match the ground truth. 407

Potential Application: Inferring the Size of Proprietary Models. The identified power-law 408 relationship offers a potential method for estimating the size of proprietary models, such GPTs. As 409 a preliminary attempt, we estimate the threshold popularity for GPT-3.5-Turbo, GPT-4, GPT-40, 410 and GPT-4o-mini. Applying the fitted power law yields size predictions of 61B, 514B, 226B, and 411 24B, respectively. The 95% confidence intervals are 12–314B, 80–3315B, 39–1313B, and 5–118B, 412 respectively. 413

414 415

387

#### MITIGATION STRATEGIES 6





w/ subsampling, r=0.1 0

loss increase > 0.05

15000

learr

Num.

w/ CKM, r=0.1



(a) 410M, trained from scratch on the mixture of FineWeb-Edu and SynBio-1.28M.

(b) 410M, continually pre-trained on the mixture of the Pile and WikiBio.

(c) 1B, continually pre-trainined on the mixture of the Pile and SynBio-2.56M.



430

424

426

While previous sections focus on how many facts the model memorizes, we now consider a more 431 practical scenario where both factual accuracy and models' general capabilities are important. In this 432 context, the model's extremely slow acquisition of low-frequency facts presents a dilemma: using a 433 small r blows up the training steps to achieve the desired factual accuracy, while recklessly increasing 434 r degrades the model's general capabilities. One could also imagine when multiple knowledge-dense 435 datasets were mixed together to form a carefully balanced mixture, setting a large r could be even 436 more detrimental. Inspired by our theory, we propose two simple yet effective mitigation strategies:

- 1. **Random Subsampling**: Randomly subsample the knowledge dataset to accelerate knowledge acquisition.
- 2. **Compact Knowledge Mixing (CKM)**: Rephrase the knowledge in a more compact form and add rephrased data to the original dataset while keeping the overall mixing ratio fixed.

We validate on both SynBio and a new real-world knowledge dataset, WikiBio, that these strategies
significantly boost the accuracy on knowledge dataset. For example, applying subsampling and CKM
to WikiBio improve the number of learned facts by 4 and 20 times, respectively. This is particularly
surprising for subsampling, as it removes a significant proportion of the knowledge data but ends up
with higher accuracy. Below, we first introduce the real-world dataset that complements SynBio.

446 447 6.1 REAL-WORLD KNOWLEDGE DATA: WIKIBIO

437

438

439

440

448 The WikiBio Dataset. To extend our study to a more real-world scenario, we curate WikiBio, a 449 dataset containing Wikipedia biographies along with ten paraphrased versions of the first paragraph 450 for 275K individuals, totaling 453M tokens. We ensure that the key information—name, occupation, 451 and birth date—is mentioned within the first paragraph. Experimenting with this dataset reflects the case where one aims to guarantee that the model can generate accurate answers to factual inquires 452 about famous people. This task is more challenging as Wikipedia biographies comprise diverse texts 453 lack of uniform formats, requiring the model to generalize to prompts that rarely have exact matches 454 in the training data. See Appendix C.2.2 for full details. 455

- Evaluation. We assess whether the model can recall a person's birth date, using this as a proxy
  for how well it memorizes the person's information. Specifically, for a (name, occupation, birth
  date) triplet, we prompt the model with "The {occupation} {name} was born on" and consider
  the response correct if it accurately includes the birth year and month in the generated text. The
  occupation is included not only to create out-of-distribution prompts but also to provide additional
  context and assist in disambiguation.
- 6.2 STRATEGY 1: RANDOM SUBSAMPLING

The approach of random subsampling may seem counterintuitive at first glance, but it becomes reasonable if we consider the frequency of each fact, which is inversely proportional to the dataset size for fixed r. Instead of mixing a large knowledge dataset that dilutes each specific fact, subsampling allows models to focus on a smaller subset, facilitating faster learning. As a result, the model memorizes more facts within the same number of training steps. In the following text, we use  $\rho$  to represent the subsampling ratio.

**Experimental Setup.** We study both pre-training from scratch and continual pre-training setups. 470 To evaluate the model's general capabilities, we use its validation loss on the web data (the Pile or 471 FineWeb-Edu) and its zero-shot performance on downstream tasks. The selected downstream tasks 472 include LAMBADA (Paperno et al., 2016), ARC-E (Clark et al., 2018), PIQA (Bisk et al., 2020), 473 SciQ (Welbl et al., 2017), and HellaSwag (Zellers et al., 2019), covering core capabilities such as text 474 understanding, commonsense reasoning, and question answering. We compare the validation loss 475 and average downstream performance to the model trained with r = 0 (no knowledge-dense data) 476 in the pre-training-from-scratch setup and to the original Pythia model in the continual pre-training setup. Downstream performance drop of more than 2% is considered unacceptable. 477

478 Subsampling enables faster fact memorization while maintaining general capability. We train 479 410M models from scratch on the mixture of FineWeb-Edu and SynBio-1.28M using mixing ratios 480  $r \in \{0, 0.1, 0.2, 0.3\}$  for a total of 32B tokens. As shown in Figures 7(a) and 8(a), FineWeb-Edu 481 validation loss and downstream performance worsen as r increases, with the degradation becoming 482 unacceptable at r = 0.3, where downstream accuracy drops by 2.09% and loss increase exceeds 483 0.05. Despite this performance decline, SynBio accuracy remains near zero. Subsampling effectively mitigates this issue. Specifically, subsampling SynBio-1.28M to 25%, 50%, and 56.25% significantly 484 improves SynBio accuracy from near 0% to 23.53%, 37.46%, and 39.81%, respectively, while 485 maintaining downstream performance within the acceptable range. Note that further increasing  $\rho$  to

62.5% makes the frequency of each biography too low, resulting in SynBio accuracy dropping back to near zero. See training details in Appendix C.5 and detailed performance in Tables 1(b) and 2(a).

**Consistent Results for Continual Pre-training.** We extend our analysis to the continual pre-489 training setup, where we continually pre-train the 410M and 1B Pythia models from their 100k-step 490 checkpoints by mixing Pile with WikiBio and SynBio-2.56M, respectively. We train 410M models 491 for a total of 32B tokens and 1B models for 64B tokens. Since mixing in the knowledge-dense data 492 introduces a distribution shift, the Pile validation loss may increase with training due to catastrophic 493 forgetting (Ibrahim et al., 2024). To keep the model's general capabilities in the acceptable range, we 494 apply early stopping when Pile validation loss increases by 0.05 for 410M models and 0.03 for 1B 495 models, both corresponding to approximately 2% drop in downstream performance. As shown in 496 Figures 7(b) and 8(b), without subsampling, setting r = 0.1 and r = 0.15 results in slow learning of WikiBio. On the other hand, increasing r to 0.2 causes the Pile validation loss to grow during 497 training, leading to early stopping after 20B tokens, resulting in poor WikiBio performance. By 498 contrast, subsampling WikiBio to 25% or 50% significantly accelerates knowledge acquisition while 499 keeping Pile validation loss within the acceptable range. For example, when fixing r = 0.1, setting 500  $\rho$  to 50% improves the number of learned facts by 4 times. Further increasing  $\rho$  to 75% proves to 501 be too high, resulting in poor performance. Similar conclusions can be drawn from experiments 502 with the 1B models, where subsampling SynBio to 50% at r = 0.2 significantly outperforms both r = 0.2 and the early-stopped r = 0.4 without subsampling, achieving a margin of approximately 504 30% (see Figure 7(c)). We defer the training details to Appendix C.5 and detailed performance 505 to Tables 1(c), 2(b) and 3.

506

# 6.3 STRATEGY 2: COMPACT KNOWLEDGE MIXING (CKM)

508 The second strategy involves rephrasing knowledge in compact forms (e.g., tuples) and adding these 509 rephrased forms to the original dataset. This approach decreases the average number of tokens needed 510 to represent each fact, thereby exposing models to each specific fact more frequently given a fixed 511 overall mixing ratio. For our specific WikiBio dataset, we compress the key information-name, 512 birth date, and occupation—into a tuple format represented as "Bio: N {name} B {birth date} 513 O {occupation}". We keep adding these tuple-form data points to WikiBio until their token 514 count reaches  $\tau$  times the total token count of the original dataset. We name this ratio  $\tau$  as Compact 515 Knowledge Mixing (CKM) ratio.

Experimental Setup. We apply CKM to WikiBio to validate its effectiveness, with the same continual pre-training setup as in Section 6.2. Each time models encounter the tuple-form data point, the order of birth date and occupation is randomly flipped. We apply early stopping when Pile validation loss increases by 0.05.

CKM significantly improves knowledge acquisition efficiency while preserving general capa-521 **bility.** We explore CKM ratios  $\tau \in \{0.1, 0.3, 0.6\}$ , fixing the overall mixing ratio r = 0.1. As 522 shown in Figures 7(b) and 8(c), CKM keeps the general capability within the acceptable range and 523 consistently boosts knowledge acquisition. Notably, performance on WikiBio improves fourfold 524 even when the token count of the added tuple-form data points constitutes only 10% of the original 525 dataset. This aligns with the phase transition in frequency predicted by our theory. Increasing  $\tau$  to 526 30% further boosts the number of learned facts to 20 times compared with training without CKM. 527 WikiBio performance saturates as  $\tau$  reaches 90%, indicating that  $\tau$  should be carefully chosen to 528 balance memorization and generalization. Detailed downstream performance is provided in Table 4.

529 530

531

# 7 DISCUSSIONS AND FUTURE DIRECTIONS

This paper identifies two phase transitions for knowledge acquisition under data mixing and develops
a theory to explain observed phenomena. While our experiments focus on mixing factual knowledge
with web corpus, such transitions may also happen for other types of knowledge, such as math, coding,
and procedural knowledge (Ruis et al., 2024). We leave the extension to more types of knowledge to
future work. Another important future direction is to apply our theory-inspired mitigation strategies
to accelerate LLMs' knowledge acquisition, especially for small models with limited capacity.

- 538
- 539

#### 540 **BROADER IMPACT** 541

542 This paper identifies two phase transitions in knowledge acquisition within data mixtures and provides 543 a theoretical understanding of these phenomena. Building on our theory, we propose two mitigation 544 strategies to enhance the efficiency of knowledge acquisition. Our findings offer deeper insights into LLM behavior and can be applied to improve the factual accuracy of LLMs.

547 REFERENCES 548

546

567

569

572

576

580

581

582

583

- 549 Allen-Zhu, Z. and Li, Y. Physics of language models: Part 3.1, knowledge storage and extraction. 550 *arXiv preprint arXiv:2309.14316*, 2023.
- 551 Allen-Zhu, Z. and Li, Y. Physics of language models: Part 3.3, knowledge capacity scaling laws. 552 arXiv preprint arXiv:2404.05405, 2024a. 553
- 554 Allen-Zhu, Z. and Li, Y. Physics of language models: Part 3.2, knowledge manipulation, 2024b. 555
- 556 Andonian, A., Anthony, Q., Biderman, S., Black, S., Gali, P., Gao, L., Hallahan, E., Levy-Kramer, J., Leahy, C., Nestler, L., Parker, K., Pieler, M., Phang, J., Purohit, S., Schoelkopf, H., Stander, 558 D., Songz, T., Tigges, C., Thérien, B., Wang, P., and Weinbach, S. GPT-NeoX: Large Scale 559 Autoregressive Language Modeling in PyTorch, 9 2023. URL https://www.github.com/ eleutherai/gpt-neox. 560
- 561 Biderman, S., Schoelkopf, H., Anthony, Q. G., Bradley, H., O'Brien, K., Hallahan, E., Khan, M. A., 562 Purohit, S., Prashanth, U. S., Raff, E., et al. Pythia: A suite for analyzing large language models 563 across training and scaling. In International Conference on Machine Learning, pp. 2397–2430. 564 PMLR, 2023. 565
- 566 Biderman, S., Prashanth, U., Sutawika, L., Schoelkopf, H., Anthony, Q., Purohit, S., and Raff, E. Emergent and predictable memorization in large language models. Advances in Neural Information 568 Processing Systems, 36, 2024.
- Bisk, Y., Zellers, R., Gao, J., Choi, Y., et al. Piqa: Reasoning about physical commonsense in 570 natural language. In Proceedings of the AAAI conference on artificial intelligence, volume 34, pp. 571 7432-7439, 2020.
- 573 Blakeney, C., Paul, M., Larsen, B. W., Owen, S., and Frankle, J. Does your data spark joy? 574 performance gains from domain upsampling at the end of training. In First Conference on 575 Language Modeling, 2024.
- Carlini, N., Ippolito, D., Jagielski, M., Lee, K., Tramer, F., and Zhang, C. Quantifying memo-577 rization across neural language models. In The Eleventh International Conference on Learning 578 Representations, 2023. 579
  - Chang, H., Park, J., Ye, S., Yang, S., Seo, Y., Chang, D.-S., and Seo, M. How do large language models acquire factual knowledge during pretraining? In The Thirty-eighth Annual Conference on Neural Information Processing Systems, 2024.
- Clark, P., Cowhey, I., Etzioni, O., Khot, T., Sabharwal, A., Schoenick, C., and Tafjord, O. Think 584 you have solved question answering? try arc, the ai2 reasoning challenge. arXiv preprint 585 arXiv:1803.05457, 2018. 586
- Da, J., Bras, R. L., Lu, X., Choi, Y., and Bosselut, A. Analyzing commonsense emergence in few-shot 588 knowledge models. In 3rd Conference on Automated Knowledge Base Construction, 2021. 589
- Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., 591 Yang, A., Fan, A., et al. The llama 3 herd of models. arXiv preprint arXiv:2407.21783, 2024.
- Feldman, V. Does learning require memorization? a short tale about a long tail. In Proceedings of 593 the 52nd Annual ACM SIGACT Symposium on Theory of Computing, pp. 954–959, 2020.

- Feng, S., Prabhumoye, S., Kong, K., Su, D., Patwary, M., Shoeybi, M., and Catanzaro, B. Maximize your data's potential: Enhancing llm accuracy with two-phase pretraining. *arXiv preprint arXiv:2412.15285*, 2024.
- Gao, L., Tow, J., Abbasi, B., Biderman, S., Black, S., DiPofi, A., Foster, C., Golding, L., Hsu,
  J., Le Noac'h, A., Li, H., McDonell, K., Muennighoff, N., Ociepa, C., Phang, J., Reynolds, L.,
  Schoelkopf, H., Skowron, A., Sutawika, L., Tang, E., Thite, A., Wang, B., Wang, K., and Zou, A.
  A framework for few-shot language model evaluation, 07 2024.
- Ge, C., Ma, Z., Chen, D., Li, Y., and Ding, B. Data mixing made efficient: A bivariate scaling law for
   language model pretraining. *arXiv preprint arXiv:2405.14908*, 2024.
- Ghosal, G. R., Hashimoto, T., and Raghunathan, A. Understanding finetuning for factual knowledge
   extraction. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235
   of *Proceedings of Machine Learning Research*, pp. 15540–15558. PMLR, 21–27 Jul 2024.
- Groeneveld, D., Beltagy, I., Walsh, E., Bhagia, A., Kinney, R., Tafjord, O., Jha, A., Ivison, H., 608 Magnusson, I., Wang, Y., Arora, S., Atkinson, D., Authur, R., Chandu, K., Cohan, A., Dumas, 609 J., Elazar, Y., Gu, Y., Hessel, J., Khot, T., Merrill, W., Morrison, J., Muennighoff, N., Naik, A., 610 Nam, C., Peters, M., Pyatkin, V., Ravichander, A., Schwenk, D., Shah, S., Smith, W., Strubell, E., 611 Subramani, N., Wortsman, M., Dasigi, P., Lambert, N., Richardson, K., Zettlemoyer, L., Dodge, J., 612 Lo, K., Soldaini, L., Smith, N., and Hajishirzi, H. OLMo: Accelerating the science of language 613 models. In Ku, L.-W., Martins, A., and Srikumar, V. (eds.), Proceedings of the 62nd Annual 614 Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Bangkok, 615 Thailand, August 2024. Association for Computational Linguistics.
- Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E., Cai, T., Rutherford, E., de Las Casas, D., Hendricks, L. A., Welbl, J., Clark, A., Hennigan, T., Noland, E., Millican, K., van den Driessche, G., Damoc, B., Guy, A., Osindero, S., Simonyan, K., Elsen, E., Vinyals, O., Rae, J., and Sifre, L. An empirical analysis of compute-optimal large language model training. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 30016–30030. Curran Associates, Inc., 2022.
- Hu, S., Tu, Y., Han, X., He, C., Cui, G., Long, X., Zheng, Z., Fang, Y., Huang, Y., Zhao, W., et al. Minicpm: Unveiling the potential of small language models with scalable training strategies. *arXiv* preprint arXiv:2404.06395, 2024.
- Huang, J., Yang, D., and Potts, C. Demystifying verbatim memorization in large language models. In
   *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 10711–10732, 2024.
- Ibrahim, A., Thérien, B., Gupta, K., Richter, M. L., Anthony, Q., Lesort, T., Belilovsky, E., and Rish,
   I. Simple and scalable strategies to continually pre-train large language models. *arXiv preprint arXiv:2403.08763*, 2024.

633

634

- Kandpal, N., Deng, H., Roberts, A., Wallace, E., and Raffel, C. Large language models struggle to learn long-tail knowledge. In *International Conference on Machine Learning*, pp. 15696–15707. PMLR, 2023.
- Kang, F., Sun, Y., Wen, B., Chen, S., Song, D., Mahmood, R., and Jia, R. Autoscale: Automatic
   prediction of compute-optimal data composition for training llms. *arXiv preprint arXiv:2407.20177*, 2024.
- Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford, A.,
  Wu, J., and Amodei, D. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- Kocetkov, D., Li, R., Allal, L. B., Li, J., Mou, C., Ferrandis, C. M., Jernite, Y., Mitchell, M., Hughes, S., Wolf, T., et al. The stack: 3 tb of permissively licensed source code. *arXiv preprint arXiv:2211.15533*, 2022.
- Li, J., Fang, A., Smyrnis, G., Ivgi, M., Jordan, M., Gadre, S., Bansal, H., Guha, E., Keh, S., Arora, K.,
   et al. Datacomp-lm: In search of the next generation of training sets for language models. *arXiv* preprint arXiv:2406.11794, 2024.

684

- 648 Li, R., Allal, L., Zi, Y., Muennighoff, N., Kocetkov, D., Mou, C., Marone, M., Akiki, C., Li, J., Chim, 649 J., et al. Starcoder: May the source be with you! Transactions on machine learning research, 2023. 650
- Liu, Q., Zheng, X., Muennighoff, N., Zeng, G., Dou, L., Pang, T., Jiang, J., and Lin, M. Regmix: 651 Data mixture as regression for language model pre-training. arXiv preprint arXiv:2407.01492, 652 2024. 653
- 654 Lu, X., Li, X., Cheng, Q., Ding, K., Huang, X., and Qiu, X. Scaling laws for fact memorization of large 655 language models. In Al-Onaizan, Y., Bansal, M., and Chen, Y.-N. (eds.), Findings of the Association for Computational Linguistics: EMNLP 2024, pp. 11263–11282, Miami, Florida, USA, November 656 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp.658. 657
- 658 Mallen, A., Asai, A., Zhong, V., Das, R., Khashabi, D., and Hajishirzi, H. When not to trust language 659 models: Investigating effectiveness of parametric and non-parametric memories. In Rogers, A., 660 Boyd-Graber, J., and Okazaki, N. (eds.), Proceedings of the 61st Annual Meeting of the Association 661 for Computational Linguistics (Volume 1: Long Papers), pp. 9802–9822, Toronto, Canada, July 662 2023. Association for Computational Linguistics.
- Nichani, E., Lee, J. D., and Bietti, A. Understanding factual recall in transformers via associative 664 memories. In The Thirteenth International Conference on Learning Representations, 2025. 665
- 666 OLMo, T., Walsh, P., Soldaini, L., Groeneveld, D., Lo, K., Arora, S., Bhagia, A., Gu, Y., Huang, 667 S., Jordan, M., Lambert, N., Schwenk, D., Tafjord, O., Anderson, T., Atkinson, D., Brahman, F., Clark, C., Dasigi, P., Dziri, N., Guerquin, M., Ivison, H., Koh, P. W., Liu, J., Malik, S., Merrill, 668 W., Miranda, L. J. V., Morrison, J., Murray, T., Nam, C., Pyatkin, V., Rangapur, A., Schmitz, M., 669 Skjonsberg, S., Wadden, D., Wilhelm, C., Wilson, M., Zettlemoyer, L., Farhadi, A., Smith, N. A., 670 and Hajishirzi, H. 2 olmo 2 furious, 2025. 671
- 672 Paperno, D., Kruszewski, G., Lazaridou, A., Pham, Q. N., Bernardi, R., Pezzelle, S., Baroni, M., 673 Boleda, G., and Fernández, R. The lambada dataset: Word prediction requiring a broad discourse 674 context. arXiv preprint arXiv:1606.06031, 2016.
- 675 Paster, K., Santos, M. D., Azerbayev, Z., and Ba, J. Openwebmath: An open dataset of high-quality 676 mathematical web text. In The Twelfth International Conference on Learning Representations, 677 2024. 678
- Penedo, G., Kydlíček, H., allal, L. B., Lozhkov, A., Mitchell, M., Raffel, C., Werra, L. V., and Wolf, 679 T. The FineWeb datasets: Decanting the web for the finest text data at scale. In *The Thirty-eight* 680 Conference on Neural Information Processing Systems Datasets and Benchmarks Track, 2024. 681
- 682 Petroni, F., Rocktäschel, T., Lewis, P., Bakhtin, A., Wu, Y., Miller, A. H., and Riedel, S. Language 683 models as knowledge bases? arXiv preprint arXiv:1909.01066, 2019.
- Qwen, :, Yang, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Li, C., Liu, D., Huang, F., Wei, 685 H., Lin, H., Yang, J., Tu, J., Zhang, J., Yang, J., Yang, J., Zhou, J., Lin, J., Dang, K., Lu, K., Bao, 686 K., Yang, K., Yu, L., Li, M., Xue, M., Zhang, P., Zhu, Q., Men, R., Lin, R., Li, T., Xia, T., Ren, X., Ren, X., Fan, Y., Su, Y., Zhang, Y., Wan, Y., Liu, Y., Cui, Z., Zhang, Z., and Qiu, Z. Qwen2.5 688 technical report, 2024. 689
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. 690 Exploring the limits of transfer learning with a unified text-to-text transformer. Journal of machine 691 learning research, 21(140):1-67, 2020. 692
- 693 Roberts, A., Raffel, C., and Shazeer, N. How much knowledge can you pack into the parameters of a 694 language model? arXiv preprint arXiv:2002.08910, 2020.
- Ruis, L., Mozes, M., Bae, J., Kamalakara, S. R., Talupuru, D., Locatelli, A., Kirk, R., Rocktäschel, 696 T., Grefenstette, E., and Bartolo, M. Procedural knowledge in pretraining drives reasoning in large 697 language models. arXiv preprint arXiv:2411.12580, 2024. 698
- Sun, K., Xu, Y., Zha, H., Liu, Y., and Dong, X. L. Head-to-tail: How knowledgeable are large 699 language models (llms)? aka will llms replace knowledge graphs? In Proceedings of the 2024 700 Conference of the North American Chapter of the Association for Computational Linguistics: 701 Human Language Technologies (Volume 1: Long Papers), pp. 311-325, 2024.

702 Team, G., Riviere, M., Pathak, S., Sessa, P. G., Hardin, C., Bhupatiraju, S., Hussenot, L., Mesnard, T., 703 Shahriari, B., Ramé, A., Ferret, J., Liu, P., Tafti, P., Friesen, A., Casbon, M., Ramos, S., Kumar, R., 704 Lan, C. L., Jerome, S., Tsitsulin, A., Vieillard, N., Stanczyk, P., Girgin, S., Momchev, N., Hoffman, 705 M., Thakoor, S., Grill, J.-B., Neyshabur, B., Bachem, O., Walton, A., Severyn, A., Parrish, A., 706 Ahmad, A., Hutchison, A., Abdagic, A., Carl, A., Shen, A., Brock, A., Coenen, A., Laforge, A., Paterson, A., Bastian, B., Piot, B., Wu, B., Royal, B., Chen, C., Kumar, C., Perry, C., Welty, C., Choquette-Choo, C. A., Sinopalnikov, D., Weinberger, D., Vijaykumar, D., Rogozińska, D., 708 Herbison, D., Bandy, E., Wang, E., Noland, E., Moreira, E., Senter, E., Eltyshev, E., Visin, F., 709 Rasskin, G., Wei, G., Cameron, G., Martins, G., Hashemi, H., Klimczak-Plucińska, H., Batra, 710 H., Dhand, H., Nardini, I., Mein, J., Zhou, J., Svensson, J., Stanway, J., Chan, J., Zhou, J. P., 711 Carrasqueira, J., Iljazi, J., Becker, J., Fernandez, J., van Amersfoort, J., Gordon, J., Lipschultz, 712 J., Newlan, J., yeong Ji, J., Mohamed, K., Badola, K., Black, K., Millican, K., McDonell, K., 713 Nguyen, K., Sodhia, K., Greene, K., Sjoesund, L. L., Usui, L., Sifre, L., Heuermann, L., Lago, L., 714 McNealus, L., Soares, L. B., Kilpatrick, L., Dixon, L., Martins, L., Reid, M., Singh, M., Iverson, 715 M., Görner, M., Velloso, M., Wirth, M., Davidow, M., Miller, M., Rahtz, M., Watson, M., Risdal, 716 M., Kazemi, M., Moynihan, M., Zhang, M., Kahng, M., Park, M., Rahman, M., Khatwani, M., 717 Dao, N., Bardoliwalla, N., Devanathan, N., Dumai, N., Chauhan, N., Wahltinez, O., Botarda, P., Barnes, P., Barham, P., Michel, P., Jin, P., Georgiev, P., Culliton, P., Kuppala, P., Comanescu, R., 718 Merhej, R., Jana, R., Rokni, R. A., Agarwal, R., Mullins, R., Saadat, S., Carthy, S. M., Cogan, 719 S., Perrin, S., Arnold, S. M. R., Krause, S., Dai, S., Garg, S., Sheth, S., Ronstrom, S., Chan, S., 720 Jordan, T., Yu, T., Eccles, T., Hennigan, T., Kocisky, T., Doshi, T., Jain, V., Yadav, V., Meshram, V., 721 Dharmadhikari, V., Barkley, W., Wei, W., Ye, W., Han, W., Kwon, W., Xu, X., Shen, Z., Gong, 722 Z., Wei, Z., Cotruta, V., Kirk, P., Rao, A., Giang, M., Peran, L., Warkentin, T., Collins, E., Barral, 723 J., Ghahramani, Z., Hadsell, R., Sculley, D., Banks, J., Dragan, A., Petrov, S., Vinyals, O., Dean, 724 J., Hassabis, D., Kavukcuoglu, K., Farabet, C., Buchatskaya, E., Borgeaud, S., Fiedel, N., Joulin, 725 A., Kenealy, K., Dadashi, R., and Andreev, A. Gemma 2: Improving open language models at a 726 practical size, 2024. 727

- Tirumala, K., Markosyan, A., Zettlemoyer, L., and Aghajanyan, A. Memorization without overfitting:
   Analyzing the training dynamics of large language models. *Advances in Neural Information Processing Systems*, 35:38274–38290, 2022.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S.,
  Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv* preprint arXiv:2307.09288, 2023.
- Welbl, J., Liu, N. F., and Gardner, M. Crowdsourcing multiple choice science questions. *arXiv* preprint arXiv:1707.06209, 2017.
- Ye, J., Liu, P., Sun, T., Zhou, Y., Zhan, J., and Qiu, X. Data mixing laws: Optimizing data mixtures by predicting language modeling performance. *arXiv preprint arXiv:2403.16952*, 2024.
- Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., and Choi, Y. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.
- Zhu, Q., Guo, D., Shao, Z., Yang, D., Wang, P., Xu, R., Wu, Y., Li, Y., Gao, H., Ma, S., et al. Deepseek-coder-v2: Breaking the barrier of closed-source models in code intelligence. *arXiv* preprint arXiv:2406.11931, 2024.
- 745 746 747
- 748
- 749
- 750 751
- 752
- 753

754

#### 756 ADDITIONAL RELATED WORKS А

757

768

758 **Knowledge Capacity Scaling Law.** LLMs are typically trained on a vast amount of data that are 759 rich in knowledge, and extensive studies have investigated how much knowledge LLMs can acquire 760 from the training data. Pioneering studies (Petroni et al., 2019; Roberts et al., 2020; Da et al., 2021) demonstrate that LLMs can capture a substantial amount of knowledge, suggesting their potential as 761 762 knowledge bases. To quantify the relationship between model size and knowledge storage, Allen-Zhu & Li (2024a) and Lu et al. (2024) discover a linear relationship between models' knowledge capacity 763 and their parameter count by training LLMs on data only containing fixed-format knowledge for 764 sufficiently long horizons. Later, Nichani et al. (2025) formally proved this linear relationship. In 765 contrast, this paper examines the data mixing scenario and demonstrates that this linear scaling can 766 be disrupted when the knowledge-dense dataset is mixed with vast amounts of web-scraped data. 767 Another important factor is the frequency of occurrence for knowledge.

**Impact of Frequency on Knowledge Acquisition.** This paper identifies phase transitions in 769 knowledge acquisition within data mixtures with respect to model size and mixing ratio. Some 770 relevant observations can be found in previous papers, but we takes a more direct and systematic 771 approach. Kandpal et al. (2023); Mallen et al. (2023); Sun et al. (2024) find that LLMs can perform 772 poorly on low-frequency knowledge. Ghosal et al. (2024) show that frequency of knowledge in 773 the pre-training data determines how well the model encodes the knowledge, which influences its 774 extractability after QA fine-tuning. Taking a more microscopic view, Chang et al. (2024) insert 775 a few pieces of new knowledge during training and track their loss. By fitting a forgetting curve, 776 they conjecture that the model may fail to learn the knowledge if its frequency is lower than some 777 threshold.

778 Memorization and Forgetting. Our findings also relate to prior observations on the memorization 779 and forgetting behaviors of LLMs, but we explicitly characterize phase transitions in the context of data mixing. Carlini et al. (2023) show that memorization of training data follows a log-linear 781 relationship with model size, the number of repetitions, and prompt length. Biderman et al. (2024) 782 take a data point-level perspective and demonstrate that it is difficult to predict whether a given data 783 point will be memorized using a smaller or partially trained model. By injecting a few new sequences 784 into the training data, Huang et al. (2024) find that a sequence must be repeated a non-trivial number 785 of times to be memorized. By examining training dynamics, Tirumala et al. (2022) observe that memorization can occur before overfitting and that larger models memorize faster while forgetting 786 more slowly. From a theoretical perspective, Feldman (2020) prove that memorization of training 787 labels is necessary to achieve near-optimal generalization error for long-tailed data distributions. 788

789 Scaling laws for Data Mixing. LLM performance is significantly influenced by the mixing propor-790 tions of the training data from different domains. Our paper is related to a line of studies that optimize 791 the mixing proportions by modeling LLM performance as a function of the mixing proportions (Liu et al., 2024; Kang et al., 2024; Ye et al., 2024; Ge et al., 2024). However, their datasets can be highly 792 heterogeneous even within a single domain (e.g., OpenWebText, Pile-CC) while we focus on mixing 793 a uniform, knowledge-dense dataset into web-scraped data. 794

- 796
- 797
- 798
- 799 800
- 801
- 802
- 803
- 804
- 805
- 806
- 809

# 810 B ADDITIONAL EXPERIMENTAL RESULTS

### **B.1** ADDITIONAL PLOTS FOR MITIGATION STRATEGIES



# 864 B.2 DETAILED PERFORMANCE ON SYNBIO

In Table 1(a), we detail the accuracy of each attribute for 70M models trained on the mixture of FineWeb-Edu and SynBio-320k with  $r \in \{0.2, 0.4, 0.8\}$ , trained for 64B, 32B, and 16B tokens respectively. We notice that the accuracy for birth date is lower than other attributes. This can be attributed to the complexity of precisely recalling the combined elements of day, month, and year information, which together form a much larger domain than other attributes. To maintain clarity and conciseness, we omit the detailed performance in other 70M experiments, as this pattern persists across them.

Furthermore, we present the detailed performance of 410M models on SynBio-1.28M corresponding to Figure 7(a) in Table 1(b). We also provide the detailed performance of 1B models on SynBio-2.56M corresponding to Figure 7(c) in Table 1(c).

Table 1: Detailed performance on SynBio. We report the accuracy (%) for each attribute averaged over five templates.

879

882 883

885 886 887

899

900

876

(a) 70M model, pre-trained from scratch on the mixture of FineWeb-Edu and SynBio-320k.

r	Birth date	Birth city	University	Major	Employer	Avg.
Random guess	0.00	0.50	0.33	1.00	0.38	0.44
0.2	0.00	0.63	0.43	1.12	0.38	0.51
0.4	16.96	45.67	41.03	50.78	43.93	39.68
0.8	79.76	88.64	88.55	90.10	88.30	87.07

(b) 410M model, pre-trained from scratch on the mixture of FineWeb-Edu and SynBio-1.28M.

N	$\rho$ (%)	r	Birth date	Birth city	University	Major	Employer	Avg.
Ra	ndom gue	ess	0.00	0.50	0.33	1.00	0.38	0.44
-	-	0.00	0.00	0.00	0.00	0.00	0.00	0.00
1.28M	100	0.1	0.00	0.42	0.33	1.01	0.21	0.39
1.28M	100	0.2	0.00	0.45	0.34	1.09	0.22	0.42
1.28M	100	0.3	0.00	0.49	0.35	1.14	0.25	0.45
320k	25	0.2	22.34	23.98	23.64	24.03	23.65	23.53
640k	50	0.2	27.97	39.66	38.51	41.50	39.68	37.46
720k	56.25	0.2	28.02	42.94	42.15	44.07	41.88	39.81
800k	62.5	0.2	0.01	1.16	0.85	3.19	0.89	1.22

(c) 1B model, continually pre-trained on the mixture of the Pile and SynBio-2.56M. Note that r = 0.4 is early stopped due to its Pile validation loss increasing beyond the acceptable range.

N	$\rho$ (%)	r	Training tokens (B)	Birth date	Birth city	University	Major	Employer	Avg.
	Rande Pythia-11	om gue B-100l	ess k-ckpt	$\begin{array}{c} 0.00 \\ 0.00 \end{array}$	$\begin{array}{c} 0.50 \\ 0.00 \end{array}$	$\begin{array}{c} 0.33 \\ 0.00 \end{array}$	$\begin{array}{c} 100 \\ 0.00 \end{array}$	$\begin{array}{c} 0.38 \\ 0.00 \end{array}$	$\begin{array}{c} 0.44 \\ 0.00 \end{array}$
2.56M 2.56M	100 100	$0.2 \\ 0.4$	64 24	$0.01 \\ 0.05$	$\begin{array}{c} 0.46 \\ 10.95 \end{array}$	$0.33 \\ 3.90$	$\begin{array}{c} 0.98\\ 4.74\end{array}$	$0.21 \\ 3.64$	$0.39 \\ 4.66$
1.28M	50	0.2	64	23.95	34.55	35.05	35.96	35.19	32.94

909 910 911

912

# **B.3** DETAILED DOWNSTREAM PERFORMANCE

We employ the lm-eval-harness (Gao et al., 2024) codebase to evaluate the zero-shot performance on five downstream tasks, including LAMBADA, ARC-E, Sciq, PIQA, and HellaSwag. We compute the validation loss on about 50M tokens on a holdout set from the Pile or FineWeb-Edu. The detailed downstream performance and validation loss for applying the random subsampling strategy to SynBio and WikiBio are presented in Tables 2 and 3, respectively. Additionally, we report the detailed downstream results for applying CKM to WikiBio in Table 4.

Table 2: Detailed downstream performance and validation loss for applying the random subsampling strategy to SynBio. We report the accuracy (%) and standard deviation (%) in the format acc. (std. dev.) for each downstream task.

				(a) 41	0M model, tra	in from scrate	ch.		
Ν	$ ho\left(\% ight)$	r	LAMBADA	ARC-E	Sciq	PIQA	HellaSwag	Avg.	FineWeb-Edu val. loss
-	-	0	$38.25_{(0.68)}$	$61.83_{(1.00)}$	83.60(1.17)	68.01 <sub>(1.09)</sub>	35.04 <sub>(0.48)</sub>	57.35	2.667
1.28M 1.28M 1.28M	100 100 100	$0.1 \\ 0.2 \\ 0.3$	$\begin{array}{c} 34.56_{(0.66)} \\ 34.43_{(0.67)} \\ 33.94_{(0.66)} \end{array}$	$\begin{array}{c} 62.33_{(0.99)} \\ 62.13_{(0.99)} \\ 60.77_{(1.00)} \end{array}$	$\begin{array}{c} 83.50_{(1.17)} \\ 83.80_{(1.17)} \\ 80.80_{(1.25)} \end{array}$	$\begin{array}{c} 68.34_{(1.09)} \\ 68.12_{(1.09)} \\ 66.54_{(1.10)} \end{array}$	$\begin{array}{c} 35.13_{(0.48)} \\ 35.39_{(0.48)} \\ 34.23_{(0.47)} \end{array}$	$\begin{array}{c} 56.77(\downarrow 0.58) \\ 56.77(\downarrow 0.58) \\ 55.26(\downarrow 2.09) \end{array}$	$\begin{array}{c} 2.668(\uparrow 0.001\\ 2.668(\uparrow 0.001\\ 2.722(\uparrow 0.054\end{array}$
320k 640k 720k 800k	$25 \\ 50 \\ 56.25 \\ 62.5$	$0.2 \\ 0.2 \\ 0.2 \\ 0.2 \\ 0.2$	$\begin{array}{c} 36.70_{(0.67)} \\ 36.58_{(0.67)} \\ 35.61_{(0.67)} \\ 35.20_{(0.67)} \end{array}$	$\begin{array}{c} 60.35_{(1.00)} \\ 60.61_{(1.00)} \\ 60.94_{(1.00)} \\ 60.48_{(1.00)} \end{array}$	$\begin{array}{c}(82.70_{1.20})\\83.30_{(1.18)}\\83.00_{(1.19)}\\83.40_{(1.20)}\end{array}$	$\begin{array}{c} 67.74_{(1.09)} \\ 66.65_{(1.10)} \\ 67.14_{(1.10)} \\ 66.54_{(1.10)} \end{array}$	$\begin{array}{r} 34.76_{(0.48)}\\ 34.53_{(0.47)}\\ 34.54_{(0.47)}\\ 34.45_{(0.47)}\end{array}$	$\begin{array}{c} 56.45(\downarrow\ 0.90)\\ 56.33(\downarrow\ 1.02)\\ 56.25(\downarrow\ 1.10)\\ 56.01(\downarrow\ 1.34)\end{array}$	$\begin{array}{c} 2.686(\uparrow\ 0.019\\ 2.688(\uparrow\ 0.021\\ 2.687(\uparrow\ 0.020\\ 2.688(\uparrow\ 0.021\\ \end{array})$

(b) 1B model, continually pre-trained. Note that r = 0.4 is early stopped due to its Pile validation loss increasing beyond the acceptable range.

Ν	ho (%)	r	Training tokens (B)	LAMBADA	ARC-E	Sciq	PIQA	HellaSwag	Avg.	Pile val. loss
	Pythia-1	B-100k-	ckpt	$55.66_{(0.69)}$	54.50 <sub>(1.02)</sub>	83.00(1.19)	$70.78_{(1.06)}$	$36.97_{(0.48)}$	60.18	2.168
2.56M 2.56M	$\begin{array}{c} 100 \\ 100 \end{array}$	$0.2 \\ 0.4$	$\begin{array}{c} 64 \\ 24 \end{array}$	$\frac{53.68_{(0.69)}}{52.38_{(0.70)}}$	$51.47_{(1.03)}$ $51.47_{(1.03)}$	$\frac{81.00_{(1.24)}}{80.70_{(1.25)}}$	${}^{68.77_{(1.08)}}_{68.17_{(1.09)}}$	$35.91_{(0.48)}$ $34.95_{(0.48)}$	$58.17(\downarrow 2.01) \\ 57.53(\downarrow 2.65)$	$\begin{array}{c} 2.184(\uparrow 0.016) \\ 2.198(\uparrow 0.030) \end{array}$
1.28M	50	0.2	64	$54.71_{(0.69)}$	$52.86_{(1.02)}$	81.30 <sub>(1.23)</sub>	$68.99_{(1.08)}$	35.48 <sub>(0.48)</sub>	$58.67(\downarrow 1.51)$	$2.189(\uparrow 0.022)$

Table 3: Detailed downstream performance for applying the random subsampling strategy to WikiBio. We use  $\rho$  to denote the subsampling ratio. We report the accuracy (%) and standard deviation (%) in the format acc. (std. dev.) for each downstream task. Note that r = 0.2 is early stopped due to its Pile validation loss increasing beyond the acceptable range.

3	N	$ ho\left(\% ight)$	r	Training tokens (B)	LAMBADA	ARC-E	Sciq	PIQA	HellaSwag	Avg.	Pile val. loss
)		Pythia	-1B-100k-	-ckpt	$50.86_{(0.70)}$	$52.10_{(1.03)}$	83.70 <sub>(1.17)</sub>	$67.14_{(1.10)}$	$34.09_{(0.47)}$	57.58	2.255
	277k 277k 277k 277k	100 100 100	$0.1 \\ 0.15 \\ 0.2$	32 32 20	$50.77_{(0.70)} \\ 49.12_{(0.70)} \\ 49.37_{(0.70)}$	$\begin{array}{c} 48.95_{(1.03)} \\ 49.66_{(1.03)} \\ 49.87_{(1.03)} \end{array}$	$\begin{array}{c} 80.80_{(1.25)} \\ 81.80_{(1.22)} \\ 79.70_{(1.27)} \end{array}$	$\begin{array}{c} 66.43_{(1.10)} \\ 66.38_{(1.10)} \\ 65.40_{(1.11)} \end{array}$	$\begin{array}{c} 33.16_{(0.47)} \\ 32.84_{(0.47)} \\ 33.08_{(0.47)} \end{array}$	$56.02(\downarrow 1.56) \\ 55.96(\downarrow 1.62) \\ 55.48(\downarrow 2.10)$	$\begin{array}{c} 2.286(\uparrow 0.031) \\ 2.292(\uparrow 0.037) \\ 2.306(\uparrow 0.051) \end{array}$
	69k 137k 208k	25 50 75	$0.1 \\ 0.1 \\ 0.1$	32 32 32	$\begin{array}{c} 48.63_{(0.70)} \\ 50.30_{(0.70)} \\ 50.34_{(0.70)} \end{array}$	$\begin{array}{c} 50.59_{(1.03)} \\ 50.38_{(1.03)} \\ 49.20_{(1.03)} \end{array}$	$\begin{array}{c} 81.00_{(1.24)} \\ 78.80_{(1.29)} \\ 80.10_{(1.26)} \end{array}$	$\begin{array}{c} 66.49_{(1.10)} \\ 66.27_{(1.10)} \\ 66.97_{(1.10)} \end{array}$	$\begin{array}{c} 33.16_{(0.47)} \\ 33.16_{(0.47)} \\ 33.19_{(0.47)} \end{array}$	$55.97(\downarrow 1.54) \\ 55.78(\downarrow 1.80) \\ 55.96(\downarrow 1.62)$	$\begin{array}{c} 2.286(\uparrow 0.031) \\ 2.285(\uparrow 0.030) \\ 2.286(\uparrow 0.031) \end{array}$

Table 4: Detailed downstream performance for applying the compact knowledge mixing strategy on WikiBio. We use  $\tau$  to denote the CKM ratio. We report the accuracy (%) and standard deviation (%) in the format  $\operatorname{acc.(std. dev.)}$  for each downstream task. Note that r = 0.2 is early stopped due to its Pile validation loss increasing beyond the acceptable range.

r	$\tau$ (%)	Training tokens (B)	LAMBADA	ARC-E	Sciq	PIQA	HellaSwag	Avg.	Pile val. loss
P	ythia-1B-1	00k-ckpt	$50.86_{(0.70)}$	$52.10_{(1.03)}$	$83.70_{(1.17)}$	$67.14_{(1.10)}$	$34.09_{(0.47)}$	57.58	2.255
0.1 0.15 0.2	0 0 0	$32 \\ 32 \\ 20$	$50.77_{(0.70)}$ $49.12_{(0.70)}$ $49.37_{(0.70)}$	$48.95_{(1.03)} \\ 49.66_{(1.03)} \\ 49.87_{(1.03)}$	$80.80_{(1.25)}$ $81.80_{(1.22)}$ $79.70_{(1.27)}$	$66.43_{(1.10)}$ $66.38_{(1.10)}$ $65.40_{(1.11)}$	$33.16_{(0.47)}$ $32.84_{(0.47)}$ $33.08_{(0.47)}$	$56.02(\downarrow 1.56)$ $55.96(\downarrow 1.62)$ $55.48(\downarrow 2.10)$	$2.286(\uparrow 0.031) \\ 2.292(\uparrow 0.037) \\ 2.306(\uparrow 0.051)$
0.1 0.1 0.1	10 30 60	32 32 32	$\begin{array}{c} 49.70_{(0.70)} \\ 50.11_{(0.70)} \\ 49.99_{(0.70)} \end{array}$	$\begin{array}{c} 49.54_{(1.03)} \\ 49.12_{(1.03)} \\ 49.41_{(1.03)} \end{array}$	$\frac{80.40_{(1.26)}}{80.20_{(1.26)}}$ $\frac{80.00_{(1.27)}}{80.00_{(1.27)}}$	$\begin{array}{c} 66.32_{(1.10)} \\ 66.54_{(1.10)} \\ 65.78_{(1.11)} \end{array}$	$\begin{array}{c} 33.11_{(0.47)}\\ 33.11_{(0.47)}\\ 32.99_{(0.47)}\end{array}$	$55.81(\downarrow 1.77) 55.82(\downarrow 1.76) 55.63(\downarrow 1.76)$	$\begin{array}{c} 2.287(\uparrow 0.032)\\ 2.285(\uparrow 0.030)\\ 2.286(\uparrow 0.031)\end{array}$

### 972 C EXPERIMENTAL DETAILS

### 974 C.1 GENERAL SETUP

976 Hyperparameters. For each training setup, we specify the batch size and learning rate in the
977 respective sections, while maintaining other hyperparameters consistent with those used in Pythia.
978 For both WSD and cosine schedulers, we use 160 steps for linear warmup. When employing the
979 WSD scheduler, we allocate the final 10% steps for cooldown.

Hardware. We train models of sizes 70M and 160M using 8 NVIDIA RTX 6000 Ada GPUs, while models of sizes 410M and 1B are trained using either 16 NVIDIA RTX 6000 Ada GPUs or 8 NVIDIA A100 GPUs. The estimated runtime required to train each model size on 1B tokens is detailed in Table 5. Consequently, a typical run training a 410M model on 32B tokens takes approximately 32 hours, whereas the longest run, which trains a 1B model on 64B tokens, exceeds five days.

Table 5: Estimated runtime required to train each model size on 1B tokens on our hardware.

Model size	Hardware	Runtime (h) per billion tokens	
70M 160M	8xNVIDIA RTX 6000 Ada	0.25 0.70	
410M 1B	16xNVIDIA RTX 6000 Ada or 8xNVIDIA A100	1.0 2.0	

**Implementation of data mixing.** During training, when loading each sample, the model flips a biased coin: with probability r, it loads from the SynBio dataset, and with probability 1 - r, it loads from the web-scraped data (FineWeb-Edu or the Pile).

1026 C.2 DETAILS OF DATASET CONSTRUCTION

# 1028 C.2.1 CONSTRUCTING THE SYNBIO DATASET

1030 To generate names, we collect a list of 400 common first names, 400 common middle names, and 1031 1000 common last names, resulting in  $1.6 \times 10^8$  unique names. To generate SynBio-N, we sample N 1032 names from this set without replacement. For each individual, the value for each attribute is randomly 1033 assigned as follows: birth date (1-28 days × 12 months × 100 years spanning 1900-2099), birth 1034 city (from 200 U.S. cities), university (from 300 institutions), major (from 100 fields of study), and employer (from 263 companies). Each attribute is paired with five sentence templates, which are used 1035 to convert (name, attribute, value) triplets into natural text descriptions. A complete list of sentence 1036 templates is provided in Table 6, and an example of a synthetic biography can be found in Table 7. 1037

- 1038
- 1039 1040

1073

1074 1075

1076

1077

1078

1079

Table 6: Sentence templates to generate the SynBio Dataset.

Attribute	Template
Birth date	<pre>{name} was born on {birth date}. {name} came into this world on {birth date}. {name}'s birth date is {birth date}. {name}'s date of birth is {birth date}. {name} celebrates {possessive pronoun} birthday on {birth date}. </pre>
Birth city	<pre>{name} spent {possessive pronoun} early years in {birth city}. {name} was brought up in {birth city}. {name}'s birthplace is {birth city}. {name} originates from {birth city}. {name} was born in {birth city}.</pre>
University	<pre>{name} received mentorship and guidance from faculty members a {university}. {name} graduated from {university}. {name} spent {possessive pronoun} college years at {university} {name} completed {possessive pronoun} degree at {university} {name} completed {possessive pronoun} academic journey a {university}.</pre>
Major	<pre>{name} completed {possessive pronoun} education with a focus o {major}. {name} devoted {possessive pronoun} academic focus to {major}. {name} has a degree in {major}. {name} focused {possessive pronoun} academic pursuits on {major}] {name} specialized in the field of {major}.</pre>
Employer	<pre>{name} is employed at {employer}. {name} a staff member at {employer}. {name} is associated with {employer}. {name} is engaged in work at {employer}. {name} is part of the team at {employer}.</pre>

Table 7: An example of a synthetic biography. The values that we expect the model to recall during evaluation are underlined.

Gracie Tessa Howell's birth date is August 09, 1992. Gracie Tessa Howell's birthplace is <u>St. Louis, MO</u>. Gracie Tessa Howell received mentorship and guidance from faculty members at <u>Santa Clara University</u>. Gracie Tessa Howell has a degree in <u>Robotics</u>. Gracie Tessa Howell is engaged in work at <u>Truist Financial</u>.

# 1080 C.2.2 CONSTRUCTING THE WIKIBIO DATASET

To create the WikiBio dataset, we query Wikidata to gather names and birth dates of individuals from 16 common occupations, identify each person's biography on Wikipedia by matching the name with the Wikipedia document title, and retain only those biographies where the person's occupation and birth date are mentioned within the first paragraph. The detailed composition is listed in Table 8. Inspired by Allen-Zhu & Li (2023), we employ Llama-3.1-70B-Instruct to paraphrase the first paragraph of each biography ten times as a form of data augmentation. See Table 9 for the prompt for paraphrasing. An example of the original text and the paraphrased versions are presented below.

Original text: Rebecca Jo Budig (born June 26, 1973) is an American actress and television presenter. Her career began in 1993, and in 1995, she was cast in the role of Michelle Bauer on the CBS soap opera Guiding Light. In 1999, she was cast as Greenlee Smythe on the ABC soap opera All My Children; she held the role off-and-on until the network series finale in 2011. In 2015, she was cast in the role of Hayden Barnes on General Hospital. In 2019, Budig was cast on L.A.'s Finest, as Carlene Hart, the drug trafficking-soccer mom of two young school-age children, a boy and a girl.

# 6 • Paraphrased versions:

- Born on June 26, 1973, Rebecca Jo Budig is an American television presenter and actress with a career spanning nearly three decades. Her professional journey began in 1993 and has included notable roles such as Michelle Bauer on Guiding Light and Greenlee Smythe on All My Children. After playing the latter role on-and-off until 2011, she went on to portray Hayden Barnes in General Hospital in 2015 and Carlene Hart in the 2019 television series L.A.'s Finest.
- 2. With a diverse career in television, Rebecca Jo Budig, born June 26, 1973, has established herself as a talented actress and presenter. Her career milestones include her roles as Michelle Bauer in the CBS soap opera Guiding Light, and Greenlee Smythe in All My Children. Her portrayal of Greenlee spanned several years, concluding with the show's finale in 2011. Budig continued to expand her repertoire with roles in General Hospital and L.A.'s Finest.
- 11083. Rebecca Jo Budig is a versatile American actress and television host, born on June 26, 1973.1109Since her career began in 1993, she has landed prominent roles in several television series. One1110of her earliest notable roles was Michelle Bauer in Guiding Light, followed by her portrayal1111of Greenlee Smythe in All My Children, a character she played until the series' conclusion in11122011. In the years that followed, she appeared in General Hospital and L.A.'s Finest.
- 1112 2011. In the years that followed, she appeared in Octobal Hospital and E.A. s Finlest.
  4. Rebecca Jo Budig, an American actress and television presenter, was born on June 26, 1973. She began her career two decades later, securing the role of Michelle Bauer on Guiding Light. Budig's subsequent roles have included Greenlee Smythe on All My Children, a part she played intermittently until the series ended in 2011. Her later appearances include a role in General Hospital and as Carlene Hart in the series L.A.'s Finest.
- 5. American actress Rebecca Jo Budig was born on June 26, 1973. Her television career, which began in 1993, encompasses multiple notable roles, such as Michelle Bauer on the soap opera Guiding Light and Greenlee Smythe on All My Children. She portrayed the latter character until the series finale in 2011. Budig later appeared as Hayden Barnes in General Hospital and took on the role of Carlene Hart in L.A.'s Finest.
- 6. Since launching her career in 1993, Rebecca Jo Budig has established herself as a talented actress and television presenter in the United States. Born on June 26, 1973, she has appeared in a range of notable roles, including Michelle Bauer on Guiding Light and Greenlee Smythe on All My Children. The latter role spanned several years, concluding with the show's finale in 2011. Her subsequent appearances include General Hospital and L.A.'s Finest.
- 7. Rebecca Jo Budig, born on June 26, 1973, has enjoyed a successful career in American television as an actress and presenter. Her breakout role came in 1995 when she was cast as Michelle Bauer on Guiding Light. Later, she played the character Greenlee Smythe on All My Children, a part she held intermittently until the show's finale in 2011. Her more recent roles include appearances in General Hospital and as Carlene Hart in L.A.'s Finest.
- 8. Born on June 26, 1973, Rebecca Jo Budig is a talented American actress and television presenter. Since her career began in 1993, she has appeared in various television series. Notable roles include her portrayal of Michelle Bauer on the soap opera Guiding Light, as

well as Greenlee Smythe on All My Children. Budig continued to expand her acting repertoire with roles in General Hospital and L.A.'s Finest, including her portrayal of Carlene Hart.

- 9. As an American actress and television host, Rebecca Jo Budig has had a diverse career spanning nearly three decades. Born on June 26, 1973, she began her professional journey in 1993. Her notable roles include Michelle Bauer on Guiding Light and Greenlee Smythe on All My Children, a character she played until the series finale in 2011. Her subsequent appearances include General Hospital and the series L.A.'s Finest, where she portrayed Carlene Hart.
- 10. With a career in television that began in 1993, Rebecca Jo Budig, born June 26, 1973, has established herself as a versatile actress and presenter. Her early roles include Michelle Bauer on Guiding Light, while her breakout role came as Greenlee Smythe on All My Children. She continued to portray Greenlee intermittently until the show's finale in 2011. Her later roles include appearances in General Hospital and L.A.'s Finest, where she took on the role of Carlene Hart.

Occupation	Num. Wikipedia biographies
Singer	18,482
Actor	31,846
Politician	38,653
Businessperson	8,068
Mathematician	5,093
Physicist	4,296
Writer	26,746
Football player	56,547
Basketball player	16,956
Sport shooter	3,156
Tennis plater	7,602
Swimmer	9,108
Painter	12,927
Volleyball player	3,556
Composer	13,719
Athlete	18,013
Total	274,768

Table 8: Detailed Composition of WikiBio.

Table 9: The prompt for paraphrasing the first paragraph of Wikipedia documents.

1172	I am creating the training data for an LLM. I
1173	→ would like to teach it to flexibly extract
1174	ightarrow knowledge from a Wikipedia paragraph.
1175	$\hookrightarrow$ Therefore, I want to diversify the Wikipedia
1176	$\hookrightarrow$ paragraphs as much as possible so that the
1177	$\hookrightarrow$ model can learn the actual relationships
1178	$ \hookrightarrow $ between entities, rather than just memorizing
1170	ightarrow the text. Please assist with the
1175	ightarrow paraphrasing task. Paraphrase the following
1180	→ Wikipedia paragraph about {Wikipedia document
1181	title} 10 times. Aim to make the paraphrased
1182	→ versions as varied as possible. Ensure all
1183	→ essential information is retained,
1184	ightarrow particularly the information about the
1185	$\hookrightarrow$ birthday and the occupuation.
1186	

# 1188 C.3 DETAILS OF THE FITTING PROCESS

1109	
1190	
1191	
1192	
1193	We use T to denote the required training steps to reach 40% accuracy and r to denote the mixing
1194	ratio.
1195	<b>Existing the exponential function</b> We fit T with respect to $\pi$ for all $\pi > 0.2$ using the function
1196	Fitting the exponential function. We fit T with respect to T for all $T \ge 0.3$ using the function $T(r) = A \exp(B/r)$ where A and B are coefficients to be fitted. Taking logarithmic on both sides
1197	$T(r) = A \exp(D/r)$ , where A and D are coefficients to be fitted. Taking logarithmic on both sides, we obtain a linear function log $T = \log A + B/r$ . By fitting log T against $1/r$ with linear regression
1198	we obtain a mean function $\log T = \log T + D/T$ . By number of against $1/T$ with mean regression, we obtain $\log A \approx -0.25512$ $B \approx 1.5137$ with goodness-of-fit $B^2 = 0.9980$
1199	we obtain $\log 11 \approx -0.20012$ , $D \approx 1.0101$ with goodness of in $11^{\circ} = 0.0000$ .
1200	Fitting the power-law function. We fit T with respect to r for all $r \in \{0.3, 0.4, 0.45, 0.5, 0.55\}$
1200	using the function $T(r) = Cr^{-D}$ , where C and D are coefficients to be fitted. Taking logarithmic on
1201	both sides, we obtain a linear function $\log T = \log C - D \log r$ . By fitting $\log T$ against $\log r$ with
1202	linear regression, we obtain $C \approx 0.098158$ , $D \approx 3.83878$ with goodness-of-fit $R^2 = 0.9853$ .
1203	
1204	
1200	
1200	
1207	
1208	
1209	
1210	
1211	C.4 DETAILS OF ESTIMATING THE THRESHOLD POPULARITY
1212	
1213	
1214	
1215	
1216	Following Mallen et al. (2023), we evaluate models using 15-shot prompting. We use the prompt
1217	presented in Table 10 for evaluation and allow models to generate up to 128 tokens with greedy
1218	we instruct the Liama 3.1.8B Instruct model to evaluate the semantic similarity between the model
1219	generated answer and the reference answer provided in $PopOA$ . The prompt used for the Llama judge
1220	is detailed in Table 11
1221	
1222	After judging the correctness of each answer, we use Algorithm 1 to estimate the popularity threshold.
1223	In our experiments, we set the target accuracy $\alpha_{\text{target}} = 60\%$ and the fault tolerance level $N_{\text{fail}} = 5$ .
1224	
1225	
1226	
1227	
1228	
1229	
1230	
1231	Table 10: The prompt for evaluating models on PopQA.
1232	
1233	
1234	You are a helpful assistant. I want to test vour
1235	$\rightarrow$ knowledge level. Here are a few examples.
1236	
1237	{few shot examples text with templates}
1238	
1239	Now, I have a question for you. Please respond in just a
1240	ightarrow few words, following the style of the examples
1241	ightarrow provided above.

Table 11: The prompt for testing synonym.

```
1244
1245
                 <|begin\_of\_text|><|start_header_id|>system<|</pre>
1246

→ end_header_id|>

1247
                 Cutting Knowledge Date: December 2023
1248
                 Today Date: 19 Dec 2024
1249
1250
                 You are a linguistic expert specializing in synonyms.
                 \leftrightarrow Your task is to determine whether two given English
1251
                 \leftrightarrow words are synonyms or not. A synonym is a word that
1252
                 \leftrightarrow has a very similar meaning to another word and can
1253
                 ↔ often replace it in sentences without significantly
1254
                 \rightarrow changing the meaning.
1255
1256
                 For each pair of words provided:
1257
                 1. Analyze their meanings and typical usage.
1258
                 2. Decide whether they are synonyms (Yes/No).
1259
                 3. Provide a brief explanation for your decision.
1260
1261
                 Here are some examples to guide you:
1262
                 Words: "happy" and "joyful"
1263
                 Yes
1264
                 Explanation: Both words describe a state of being
1265
                 \rightarrow pleased or content and are often interchangeable in
1266
                 → most contexts.
1267
1268
                 Words: "run" and "jog"
1269
                 No
1270
                 Explanation: While both refer to forms of movement, "run"
1271
                 → typically implies a faster pace than "jog."
1272
                 Words: "angry" and "frustrated"
1273
                 No
1274
                 Explanation: Although both express negative emotions, "
1275
                 → angry" implies strong displeasure or rage, while "
1276
                 → frustrated" conveys annoyance due to obstacles or
1277
                 \rightarrow failure.
1278
1279
                 <|eot id|><|start header id|>user<|end header id|>
1280
1281
                 Words: {} and {}
1282
                 <|eot_id|><|start_header_id|>assistant<|end_header_id|>
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
```

Igorithm 1: Estimate Threshold Popularity
Input:
- x: A list of popularity values for each data point, where $x_i$ represents the popularity of the
<i>i</i> -th data point.
- y: A list of binary values indicating the correctness of the model's response, where $y_i = 1$ if
the model answers the <i>i</i> -the question correctly, and $y_i = 0$ otherwise.
- $\alpha_{target}$ : The target accuracy.
- $N_{\text{fail}}$ : The maximum number of failures before termination, denoting the fault tolerance
level.
Output:
- $P_{\rm thres}$ : The threshold popularity.
Initialize correct count: sum_correct $\leftarrow 0$
Initialize error count: $e \leftarrow 0$
Sort $(x, y)$ by x in ascending order and store the indices in a list I.
Initialize loop variable $j \leftarrow \text{len}(x) - 1$
Initialize flag counter flag $\leftarrow 0$
while $j \ge 0$ do
$k \leftarrow j$
while $k \ge 0$ and $x_{I_k} = x_{I_j}$ do
$k \leftarrow k - 1$
end while
$ \lim_{k \to \infty} l = k + 1 \text{ to } j \text{ do} $
$i \leftarrow I_l$
sum_contect $\leftarrow$ sum_contect + $y_i$
if sum_correct < set threshold then
$\lim_{x \to \infty} \frac{1}{ en(x)-k-1 } \leq \operatorname{set-uncestord}$ then
$e \leftarrow e + 1$
the line is $a = N_{\rm exp}$ where $b = N_{\rm exp}$
<b>Return</b> $T_{\text{tail}}$ <b>Return</b> the threshold popularity
end if
$i \leftarrow k$
end while

1341

1342

1337

# D PROOFS OF THEORETICAL RESULTS

We follow the notations in Section 4. We use  $H(\cdot)$  to denote the entropy and  $I(\cdot; \cdot)$  to denote the mutual information.

corresponds to the original learning rate used at step 100k in the Pythia model training.

without experiencing extreme loss spikes. Specifically, we continually pre-train 410M and 1B Pythia

models from their respective 100k-step checkpoint with a constant learning rate of  $8.7 \times 10^{-5}$ , which

1345 1346 We define a data distribution  $\mathcal{D}$  as a distribution over (x, y), where x is an input and y is a token. A data universe  $\mathcal{U} = (\mathcal{P}, \mathcal{D}_{\theta})$  is defined by a prior  $\mathcal{P}$  over a latent variable  $\theta$  and a family of data distributions  $\mathcal{D}_{\theta}$  indexed by  $\theta$ .

1349 A predictor h is a function that maps x to a distribution over y. A learning algorithm  $\mathcal{A}$  is a procedure that takes samples from a data distribution  $\mathcal{D}$  of (x, y) and outputs a predictor  $h \sim \mathcal{A}(\mathcal{D})$  in the end.

1350 1351	For a given predictor $h$ , we measure its performance by the expected cross-entropy loss $\mathcal{L}(h; \mathcal{D}) := \mathbb{E}_{(x,y)\sim\mathcal{D}}[-\log p(y \mid h, x)],  (7)$
1352 1353 1354 1355	where $p(y \mid h, x)$ denotes the predicted distribution of $y$ given $x$ by the predictor $h$ , and log is in base 2 for convenience. For a data universe $\mathcal{U} = (\mathcal{P}, \mathcal{D}_{\theta})$ , we measure the performance of a learning algorithm $\mathcal{A}$ by its expected loss over all data distributions $\mathcal{D}_{\theta}$ with respect to the prior $\mathcal{P}$ : $\bar{\mathcal{L}}_{\mathcal{P}}(\mathcal{A}) := \mathbb{E}_{\theta \sim \mathcal{P}} \mathbb{E}_{h \sim \mathcal{A}(\mathcal{D}_{\theta})}[\mathcal{L}(h; \mathcal{D}_{\theta})].$ (8)
1356 1357 1358	We use the mutual information $I(\mathcal{A}(\mathcal{D}_{\theta}); \mathcal{D}_{\theta})$ as a measure of the <i>effective model capacity</i> for the predictor picked by $\mathcal{A}$ on $\mathcal{D}_{\theta}$ , where $\theta$ is sampled from $\mathcal{Q}$ .
1359 1360	Same as Definition 4.1, for a data universe $\mathcal{U} = (\mathcal{P}, \mathcal{D}_{\theta})$ and $M > 0$ , we define the best achievable loss under a capacity constraint $M$ as
1361	$F_{\mathcal{P}}(M) := \inf_{\mathcal{A}} \left\{ \mathcal{L}_{\mathcal{P}}(\mathcal{A}) : I(\mathcal{A}(D_{\theta}); D_{\theta}) \le M \right\},\tag{9}$
1362 1363 1364	where the infimum is taken over all learning algorithms. An optimal $M$ -bounded-capacity learner is a learning algorithm $\mathcal{A}$ such that $I(\mathcal{A}(\mathcal{D}_{\theta}); \mathcal{D}_{\theta}) \leq M$ and $\overline{\mathcal{L}}_{\mathcal{P}}(\mathcal{A}) = F_{\mathcal{P}}(M)$ .
1365 1366	D.1 CONVEXITY OF THE BEST ACHIEVABLE LOSS
1367 1368 1369	It is easy to see that $F_{\mathcal{P}}(M)$ is non-negative and non-increasing in $M$ . A classic result in rate distortion theory is that the rate distortion function is convex. This further implies that $F_{\mathcal{P}}(M)$ is convex in $M$ . Here we present it as a lemma for completeness.
1370 1371	<b>Lemma D.1.</b> For any data universe $\mathcal{U} = (\mathcal{P}, \mathcal{D}_{\theta})$ , $F_{\mathcal{P}}(M)$ is convex in $M$ .
1372 1373 1374	<i>Proof.</i> Let $\epsilon > 0$ be any positive number. Let $\mathcal{A}_1$ be a learning algorithm that achieves a loss $\leq F_{\mathcal{P}}(M_1) + \epsilon$ with mutual information $\leq M_1$ and $\mathcal{A}_2$ be a learning algorithm that achieves a loss $F_{\mathcal{P}}(M_2) + \epsilon$ with mutual information $\leq M_2$ .
1375 1376 1377 1378	Let $\mathcal{A}$ be a new learning algorithm that outputs the same as $\mathcal{A}_1$ with probability $1 - p$ and the same as $\mathcal{A}_2$ with probability $p$ . Then the mutual information between $\mathcal{A}(\mathcal{D}_{\theta})$ and $\mathcal{D}_{\theta}$ is $I(\mathcal{A}(\mathcal{D}_{\theta}); \mathcal{D}_{\theta}) = (1 - p)I(\mathcal{A}_1(\mathcal{D}_{\theta}); \mathcal{D}_{\theta}) + pI(\mathcal{A}_2(\mathcal{D}_{\theta}); \mathcal{D}_{\theta})$
1379	$\leq (1-p)M_1 + pM_2.$
1380 1381 1382	By linearity of expectation, the expected loss of $\mathcal{A}$ can be bounded as $\mathbb{E}_{\theta \sim \mathcal{P}(\theta)}[\mathcal{L}(\mathcal{A}(\mathcal{D}_{\theta}); \mathcal{D}_{\theta})] = (1 - p)\mathbb{E}_{\theta \sim \mathcal{P}(\theta)}[\mathcal{L}(\mathcal{A}_{1}(\mathcal{D}_{\theta}); \mathcal{D}_{\theta})] + p\mathbb{E}_{\theta \sim \mathcal{P}(\theta)}[\mathcal{L}(\mathcal{A}_{2}(\mathcal{D}_{\theta}); \mathcal{D}_{\theta})]$ $\leq (1 - p)F_{\mathcal{P}}(M_{1}) + pF_{\mathcal{P}}(M_{2}) + 2\epsilon.$
1383 1384	Therefore, we have $F_{\mathcal{P}}((1-p)M_1 + pM_2) \leq \mathbb{E}_{\theta \sim \mathcal{P}(\theta)}[\mathcal{L}(\mathcal{A}(\mathcal{D}_{\theta}); \mathcal{D}_{\theta})] \leq (1-p)F_{\mathcal{P}}(M_1) + pF_{\mathcal{P}}(M_2) + 2\epsilon,$
1385 1386	taking $\epsilon \to 0$ finishes the proof.
1387 1388	D.2 PROOFS FOR THE WARMUP CASE
1389 1390 1391	<b>Definition D.2</b> (Factual Data Universe). We define a fact as a pair $(X, y)$ , where X is a set of inputs and y is a target token. A factual data universe is a data universe $\mathcal{U} = (\mathcal{P}, \mathcal{D}_{\theta})$ containing K random facts $(X_1, y_1), \ldots, (X_K, y_K)$ in the following way:
1392 1393	1. $X_1, \ldots, X_K$ are K disjoint sets of inputs, and $y_1, \ldots, y_K$ are random tokens;
1394 1395	2. $\theta$ is structured as $(y_1, \ldots, y_K)$ . Given $\theta = (y_1, \ldots, y_K)$ , the data distribution $\mathcal{D}_{\theta}$ satisfies that for all $x \in X_i$ , $\mathcal{D}_{\theta}(y \mid x_i)$ is a point mass at $y_i$ ;
1396 1397	3. For all $\theta$ , the input distribution $\mathcal{D}_{\theta}(x)$ is the same;
1398 1399	4. For all $\theta$ , the target distribution $\mathcal{D}_{\theta}(y \mid x)$ is the same for all $x \notin \bigcup_{i=1}^{K} X_i$ ;
1400 1401 1402	5. The prior distribution $\mathcal{P}$ over $\theta$ is given by the product distribution $\mathcal{P}(y_1, y_2, \dots, y_K) = \prod_{k=1}^{K} \mathcal{Y}_k(y_k)$ , where $\mathcal{Y}_k$ is a fixed prior distribution over $y_k$ .
1-104	

The exposure frequency of each random fact is defined as the total probability that an input  $x \in X_i$  occurs in  $\mathcal{D}_{\theta}$ .

**Theorem D.3** (Theorem 4.2, restated). For a factual data universe  $\mathcal{U} = (\mathcal{P}, \mathcal{D}_{\theta})$  with K random facts, if all the facts have the same exposure frequency p, then  $E_{\theta}(M) = C + \pi \max \left[ \frac{W}{W} - \frac{W}{W} \right]$  (10)

$$F_{\mathcal{P}}(M) = C + p \cdot \max\{H_{\text{tot}} - M, 0\},$$
(10)

1407 1408

where  $H_{\text{tot}} := \sum_{i=1}^{K} H(\mathcal{Y}_i)$  and  $C := F_{\mathcal{P}}(\infty)$ .

Proof. First, we prove a lower bound for  $F_{\mathcal{P}}(M)$ . For any learning algorithm  $\mathcal{A}$  with  $I(\mathcal{A}(\mathcal{D}_{\theta}); \mathcal{D}_{\theta}) \leq M$ ,

$$\mathcal{L}_{\mathcal{P}}(\mathcal{A}(\mathcal{D}_{\theta})) = \mathbb{E}_{\theta \sim \mathcal{P}} \mathbb{E}_{(x,y) \sim \mathcal{D}_{\theta}} \mathbb{E}_{h \sim \mathcal{A}(\mathcal{D}_{\theta})} [-\log p(y \mid h, x)]$$

$$= \mathbb{E}_{x} \mathbb{E}_{\theta \sim \mathcal{P}} \mathbb{E}_{y \sim \mathcal{D}_{\theta}}(\cdot | x) \mathbb{E}_{h \sim \mathcal{A}(\mathcal{D}_{\theta})} \mathbb{E}[-\log p(y \mid h, x)$$

1412 1413

$$\geq \underbrace{\mathbb{E}_{x} \left[ \mathbb{1}_{\{x \in \bigcup_{i=1}^{K} X_{i}\}} H_{\theta \sim \mathcal{P}}(\mathcal{D}_{\theta}(\cdot \mid x)) \right]}_{=:C_{0}} + p \left[ \sum_{i=1}^{K} \left( H(\mathcal{Y}_{i}) - I(\mathcal{A}(\mathcal{D}_{\theta}); y_{i}) \right) \right]$$

1420

1425 1426

1436

1438

1441

1442

1443

1444 1445

1451

1452

1453 1454 1455

$$\geq C_0 + p \left[ H_{\text{tot}} - I(\mathcal{A}(\mathcal{D}_{\theta}); \mathcal{D}_{\theta}) \right]_+$$

1419 
$$\geq C_0 + p [H_{\text{tot}} - M]_+$$
.

For upper bounds, we first show that  $F_{\mathcal{P}}(M) \leq C_0$  for all  $M \geq H(\theta)$ . Let  $\mathcal{A}_1$  be the learning algorithm that inputs  $\mathcal{D}_{\theta}$  and outputs the predictor h that always outputs the token  $y_i$  for the input  $x \in X_i$ . For all the other inputs x, the predictor just outputs  $h(y \mid h, x) = \mathbb{E}_{\theta \sim \mathcal{P}}[\mathcal{D}_{\theta}(y \mid x)]$ . Both  $\mathcal{A}_1(\mathcal{D}_{\theta})$  and  $\mathcal{D}_{\theta}$  can be transformed from  $\theta$  with a reversible function, so

$$I(\mathcal{A}_1(\mathcal{D}_\theta); \mathcal{D}_\theta) = H(\theta) = \sum_{i=1}^{K} H(\mathcal{Y}_i) = H_{\text{tot}}$$

1427 It is easy to see that  $\overline{\mathcal{L}}_{\mathcal{P}}(\mathcal{A}_1) = C_0$ . This implies that  $F_{\mathcal{P}}(M) \leq C_0$  for all  $M \geq H(\theta)$ .

1429 Now, if  $M < H(\theta)$ , we construct a learning algorithm  $\mathcal{A}_q$  that outputs the same as  $\mathcal{A}_1$  with probability 1430 q and outputs  $h(y \mid h, x) = \mathbb{E}_{\theta \sim \mathcal{P}}[\mathcal{D}_{\theta}(y \mid x)]$  with probability 1 - q. Setting  $q = \frac{M}{H(\theta)}$ , we have 1431  $I(\mathcal{A}_q(\mathcal{D}_{\theta}); \mathcal{D}_{\theta}) = q \cdot H(\theta) = M$ .

By linearity of expectation, we also have  $\bar{\mathcal{L}}_{\mathcal{P}}(\mathcal{A}_q(\mathcal{D}_\theta)) = \bar{\mathcal{L}}_{\mathcal{P}}(\mathcal{A}_1) + (1-q) \cdot p \sum_{i=1}^K H(\mathcal{Y}_i)$ . This implies that  $F_{\mathcal{P}}(M) \leq \bar{\mathcal{L}}_{\mathcal{P}}(\mathcal{A}_1) + p \cdot \max\{H_{\text{tot}} - M, 0\}$  for all  $M < H_{\text{tot}}$ .

1435 Putting all the pieces together finishes the proof.

+

# 1437 D.3 PROOFS FOR THE DATA MIXING CASE

**Definition D.4** (Mixture of Data Universes). Let  $U_1 = (\mathcal{P}_1, \mathcal{D}_{\theta_1})$  and  $U_2 = (\mathcal{P}_2, \mathcal{D}_{\theta_2})$  be two data universes. We mix them together to form a new data universe  $\mathcal{U} = (\mathcal{P}, \mathcal{D}_{\theta})$ :

- 1.  $\theta$  is structured as  $(\theta_1, \theta_2)$ . Given  $\theta = (\theta_1, \theta_2)$ , the data distribution  $\mathcal{D}_{\theta}$  is formed as  $\mathcal{D}_{\theta} = r\mathcal{D}_{\theta_1} + (1-r)\mathcal{D}_{\theta_2}$ , where r is called *the mixing ratio*;
- 2. The prior distribution  $\mathcal{P}$  over  $\theta$  is a joint distribution of  $\mathcal{P}_1$  and  $\mathcal{P}_2$ .

In the reality, mixing two datasets can be seen as mixing two data universes first and then sampling a data distribution from the mixed data universe. Here we consider the simplified case where the two data universes are so different from each other that they convey orthogonal information.

**1449 Definition D.5** (Orthogonal Mixture of Data Universes). We say that  $\mathcal{U}$  is an orthogonal mixture of  $\mathcal{U}_1$  and  $\mathcal{U}_2$  if

1. For any x that is in the both supports of  $\mathcal{D}_{\theta_1}$  and  $\mathcal{D}_{\theta_2}$ , we have  $\mathcal{D}_{\theta_1}(y \mid x) = \mathcal{D}_{\theta_2}(y \mid x)$  for all  $\theta_1$  and  $\theta_2$ ;

2. 
$$\mathcal{P}(\theta_1, \theta_2) = \mathcal{P}_1(\theta_1) \cdot \mathcal{P}_2(\theta_2)$$
, i.e.,  $\theta_1$  and  $\theta_2$  are independent.

Below, we first establish two lemmas that provide conditions for when the loss on the first domain is
very low or very high for an optimal *M*-bounded-capacity learner running on an orthogonal mixture of two data universes. Then we use these lemmas to prove the main theorem we stated in Section 4.3.

We use  $D^{-}F(x)$  and  $D^{+}F(x)$  to denote the left and right derivatives of a function F at a point x, respectively.

$$\frac{r}{1-r} < \frac{D \ F_{\mathcal{P}_2}(M)}{D^+ F_{\mathcal{P}_1}(0)},\tag{11}$$

then for any optimal *M*-bounded-capacity learner  $\mathcal{A}$  on  $\mathcal{U}$ ,  $\mathbb{E}_{\theta \sim \mathcal{P}} \mathbb{E}_{h \sim \mathcal{A}(\mathcal{D}_{\theta})}[\mathcal{L}(h; \mathcal{D}_{\theta_1})] = F_{\mathcal{P}_1}(0).$ 

*Proof.* Let h be the predictor picked by  $\mathcal{A}$  on  $\mathcal{D}_{\theta}$ . Let  $\mathcal{X}_1$  and  $\mathcal{X}_2$  be the supports of x in  $\mathcal{D}_{\theta_1}$  and 1468  $\mathcal{D}_{\theta_2}$ , respectively. Let  $m_1 := I(h|_{\mathcal{X}_1}; \mathcal{D}_{\theta_1})$  and  $m_2 := I(h|_{\mathcal{X}_2}; \mathcal{D}_{\theta_2})$ . By data processing inequality, 1469 we have

$$m_1 = I(h|_{\mathcal{X}_1}; \mathcal{D}_{\theta_1}) \le I(h; \mathcal{D}_{\theta_1})$$

$$m_2 = I(h|_{\mathcal{X}_2}; \mathcal{D}_{\theta_2}) \le I(h; \mathcal{D}_{\theta_2}),$$

Further noticing that  $I(h; \mathcal{D}_{\theta}) = I(h; \mathcal{D}_{\theta_1}; \mathcal{D}_{\theta_2}) \ge I(h; \mathcal{D}_{\theta_1}) + I(h; \mathcal{D}_{\theta_2})$ , we have  $m_1 + m_2 \le I(h; \mathcal{D}_{\theta}) \le M$ .

1474 Since  $h|_{\chi_1}$  and  $h|_{\chi_2}$  are valid predictors on  $\mathcal{D}_{\theta_1}$  and  $\mathcal{D}_{\theta_2}$ , respectively, we have 1475  $\mathbb{E}[\mathcal{L}(h;\mathcal{D}_2)] = \mathbb{E}[\mathcal{L}(h|_1;\mathcal{D}_2)] \ge E_{\mathcal{D}_2}(m_1)$ 

$$\mathbb{E}[\mathcal{L}(h; \mathcal{D}_{\theta_1})] = \mathbb{E}[\mathcal{L}(h|_{\mathcal{X}_1}; \mathcal{D}_{\theta_1})] \ge F_{\mathcal{P}_1}(m_1)$$

$$\mathbb{E}[\mathcal{L}(h;\mathcal{D}_{\theta_2})] = \mathbb{E}[\mathcal{L}(h|_{\mathcal{X}_2};\mathcal{D}_{\theta_2})] \ge F_{\mathcal{P}_2}(m_2) \ge F_{\mathcal{P}_2}(M-m_1).$$

1478 Adding the two inequalities with weights r and 1 - r, we have

$$\bar{\mathcal{L}}_{\mathcal{P}}(\mathcal{A}) = \mathbb{E}[\mathcal{L}(h; \mathcal{D}_{\theta})] \ge rF_{\mathcal{P}_1}(m_1) + (1-r)F_{\mathcal{P}_2}(M-m_1).$$

1480 By convexity (Lemma D.1), we have

$$F_{\mathcal{P}_1}(m_1) \ge F_{\mathcal{P}_1}(0) + \mathrm{D}^+ F_{\mathcal{P}_1}(0)m_1, \qquad F_{\mathcal{P}_2}(M - m_1) \ge F_{\mathcal{P}_2}(M) - \mathrm{D}^- F_{\mathcal{P}_2}(M)m_1.$$
  
Plugging these into the previous inequality, we have

 $\mathbb{E}[\mathcal{L}(h; \mathcal{D}_{\theta})] \ge rF_{\mathcal{P}_1}(0) + (1-r)F_{\mathcal{P}_2}(M) + \left(rD^+F_{\mathcal{P}_1}(0) - (1-r)D^-F_{\mathcal{P}_2}(M)\right)m_1.$ 

1484 By (11) and the fact that  $D^+F_{\mathcal{P}_1}(0) \leq 0$ , we have  $rD^+F'_{\mathcal{P}_1}(0) > (1-r)D^-F'_{\mathcal{P}_2}(M)$ . So the right-hand side is strictly increasing in  $m_1$ .

Now we claim that  $m_1 = 0$ . If not, then the following learning algorithm  $\mathcal{A}'$  is better than  $\mathcal{A}$ . Let  $\mathcal{A}_1$ be an optimal 0-bounded-capacity learner on  $\mathcal{U}_1$  and  $\mathcal{A}_2$  be an optimal M-bounded-capacity learner on  $\mathcal{U}_2$ . Run the algorithms to obtain  $h_1 \sim \mathcal{A}_1(\mathcal{D}_{\theta}|_{\mathcal{X}_1})$  and  $h_2 \sim \mathcal{A}_2(\mathcal{D}_{\theta}|_{\mathcal{X}_2})$ . Then, whenever seeing an input x from  $\mathcal{X}_1$ , output  $h_1(x)$ ; otherwise output  $h_2(x)$ . This algorithm achieves the expected loss  $rF_{\mathcal{P}_1}(0) + (1-r)F_{\mathcal{P}_2}(M)$ , which is strictly less than  $\mathcal{L}_{\mathcal{P}}(\mathcal{A})$  and contradicts the optimality of  $\mathcal{A}$ .

1492 Therefore, for the algoritm  $\mathcal{A}$ ,  $\mathbb{E}[\mathcal{L}(h; \mathcal{D}_{\theta_1})] \geq F_{\mathcal{Q}_1}(m_1) = F_{\mathcal{Q}_1}(0)$ . In fact,  $\mathbb{E}[\mathcal{L}(h; \mathcal{D}_{\theta_1})] = F_{\mathcal{Q}_1}(m_1)$  must hold because otherwise  $\mathcal{A}'$  is sitll better than  $\mathcal{A}$ .

**Lemma D.7.** Let  $\mathcal{U} = (\mathcal{P}, \mathcal{D}_{\theta})$  be an orthogonal mixture of  $\mathcal{U}_1 = (\mathcal{P}_1, \mathcal{D}_{\theta_1})$  and  $\mathcal{U}_2 = (\mathcal{P}_2, \mathcal{D}_{\theta_2})$ 1495 with mixing ratio r. For all  $r \in (0, 1)$ ,  $M \ge 0$  and  $\beta \ge 0$ , if the following inequality holds, 1496  $r = D^+ F_{\mathcal{P}_2}(M - \beta)$ 

$$\frac{r}{1-r} > \frac{D^{+}F_{\mathcal{P}_{2}}(M-\beta)}{D^{-}F_{\mathcal{P}_{1}}(\beta)},$$
(12)

then for any optimal *M*-bounded-capacity learner  $\mathcal{A}$  on  $\mathcal{U}$ ,  $\mathbb{E}_{\theta \sim \mathcal{P}} \mathbb{E}_{h \sim \mathcal{A}(\mathcal{D}_{\theta})}[\mathcal{L}(h; \mathcal{D}_{\theta_1})] \leq F_{\mathcal{P}_1}(\beta)$ .

1501 Proof. Similar to the previous proof, letting  $m_1 := I(h|_{\chi_1}; \mathcal{D}_{\theta_1})$  and  $m_2 := I(h|_{\chi_2}; \mathcal{D}_{\theta_2})$ , we have 1502  $m_1 + m_2 \leq I(h; \mathcal{D}_{\theta}) \leq M$ ,

$$\mathbb{E}[\mathcal{L}(h; \mathcal{D}_{\theta_1})] = \mathbb{E}[\mathcal{L}(h|_{\mathcal{X}_1}; \mathcal{D}_{\theta_1})] \ge F_{\mathcal{P}_1}(m_1)$$

$$\mathbb{E}[\mathcal{L}(h;\mathcal{D}_{\theta})] \ge rF_{\mathcal{P}_1}(m_1) + (1-r)F_{\mathcal{P}_2}(M-m_1).$$

First, we show that  $m_1 \ge \beta$ . If not, then by convexity (Lemma D.1), we have

1507 First, we show that  $m_1 \leq \beta$ . If not, then by contextly (Definite D17), we have  $F_{\mathcal{P}_1}(m_1) \geq F_{\mathcal{P}_1}(\beta) - D^- F_{\mathcal{P}_1}(\beta) \cdot (\beta - m_1), \quad F_{\mathcal{P}_2}(M - m_1) \geq F_{\mathcal{P}_2}(M - \beta) + D^+ F_{\mathcal{P}_2}(M - \beta) \cdot (\beta - m_1).$ Plugging these into the previous inequality, we have

<sup>1509</sup>  $\mathbb{E}[\mathcal{L}(h; \mathcal{D}_{\theta})] \ge rF_{\mathcal{P}_{1}}(\beta) + (1-r)F_{\mathcal{P}_{2}}(M-\beta) + (-rD^{-}F_{\mathcal{P}_{1}}(\beta) + (1-r)D^{+}F_{\mathcal{P}_{2}}(M-\beta))(\beta-m_{1}).$ 

 $\mathbb{E}[\mathcal{L}(h; \mathcal{D}_{\theta_2})] = \mathbb{E}[\mathcal{L}(h|_{\mathcal{X}_2}; \mathcal{D}_{\theta_2})] \ge F_{\mathcal{P}_2}(m_2) \ge F_{\mathcal{P}_2}(M - m_1),$ 

1511 By (12) and the fact that  $D^-F_{\mathcal{P}_1}(\beta) \leq 0$ , we have  $rD^-F_{\mathcal{P}_1}(\beta) < (1-r)D^+F_{\mathcal{P}_2}(M-\beta)$ . So the right-hand side is strictly decreasing in  $m_1$ .

Then we claim that the following learning algorithm  $\mathcal{A}'$  is better than  $\mathcal{A}$ , leading to contradition. Let  $\mathcal{A}_1$  be an optimal  $\beta$ -bounded-capacity learner on  $\mathcal{U}_1$  and  $\mathcal{A}_2$  be an optimal  $(M - \beta)$ -boundedcapacity learner on  $\mathcal{U}_2$ . Run the algorithms to obtain  $h_1 \sim \mathcal{A}_1(\mathcal{D}_{\theta}|_{\mathcal{X}_1})$  and  $h_2 \sim \mathcal{A}_2(\mathcal{D}_{\theta}|_{\mathcal{X}_2})$ . Then, whenever seeing an input x from  $\mathcal{X}_1$ , output  $h_1(x)$ ; otherwise output  $h_2(x)$ . This algorithm achieves the expected loss  $rF_{\mathcal{P}_1}(\beta) + (1 - r)F_{\mathcal{P}_2}(M - \beta)$ , which is strictly less than  $\overline{\mathcal{L}_{\mathcal{P}}}(\mathcal{A})$ .

Therefore, we have  $m_1 \geq \beta$  for the algorithm  $\mathcal{A}$ . Now we prove that  $\mathbb{E}[\mathcal{L}(h; \mathcal{D}_{\theta_1})] \leq F_{\mathcal{P}_1}(\beta)$ . If not, then the following learning algorithm  $\mathcal{A}''$  is better than  $\mathcal{A}$ . Construct  $\mathcal{A}''$  similarly as  $\mathcal{A}'$ , but with  $\mathcal{A}_1$  and  $\mathcal{A}_2$  replaced by the optimal  $m_1$ -bounded-capacity learner on  $\mathcal{U}_1$  and the optimal  $m_2$ -bounded-capacity learner on  $\mathcal{U}_2$ , respectively. If  $\mathbb{E}[\mathcal{L}(h; \mathcal{D}_{\theta_1})] > F_{\mathcal{P}_1}(\beta)$ , then  $\mathcal{A}''$  achieves a lower expected loss than  $\mathcal{A}$ , which contradicts the optimality of  $\mathcal{A}$ .

1522

1523 Now we consider the case where  $U_1$  is a factual data universe and  $U_2$  is an arbitrary data universe. 1524 For any learning algorithm  $\mathcal{A}$ , define  $\overline{\mathcal{L}}_1(\mathcal{A}) := \mathbb{E}_{\theta \sim \mathcal{P}_1}[\mathcal{L}(\mathcal{A}(\mathcal{D}_\theta); \mathcal{D}_{\theta_1})]$ , which is the expected loss 1525 of  $\mathcal{A}$  on the first domain after learning from the data mixture.

**Theorem D.8.** Let  $U_1$  be a factual data universe with K random facts, each with the same exposure frequency p, and the entropies of their target tokens sum to  $H_{\text{tot}} := \sum_{i=1}^{K} H(\mathcal{Y}_i)$ . Let  $U_2$  be an arbitrary data universe. Let  $\mathcal{U} = (\mathcal{P}, \mathcal{D}_{\theta})$  be an orthogonal mixture of  $\mathcal{U}_1$  and  $\mathcal{U}_2$  with mixing ratio r. For all  $r \in (0, 1)$  and  $M \ge 0$ ,

1534

1537

1539 1540

1544

1546

1. if  $\frac{r}{1-r} \cdot p < -D^{-}F_{\mathcal{P}_{2}}(M)$ , then  $\bar{\mathcal{L}}_{1}(\mathcal{A}) = F_{\mathcal{P}_{1}}(0)$ ;

1533 2. if 
$$\frac{r}{1-r} \cdot p > -D^+ F_{\mathcal{P}_2}(M - H_{\text{tot}})$$
, then  $\bar{\mathcal{L}}_1(\mathcal{A}) = F_{\mathcal{P}_1}(\infty)$ 

1535 *Proof.* By Theorem D.3,  $D^+F_{\mathcal{P}_1}(0) = D^-F_{\mathcal{P}_1}(H_{tot}) = p$ . Plugging this into Lemma D.6 1536 and Lemma D.7 with  $\beta = H_{tot}$  finishes the proof.

1538 Now we are ready to prove the main theorem we stated in Section 4.3.

$$M_0^-(x) := \sup\{M \ge 0 : -F'_{\mathcal{P}_2}(M) > x\},\$$

$$M_0^+(x) := \inf\{M \ge 0 : -F'_{\mathcal{P}_2}(M) < x\}$$

**Theorem D.9** (Theorem 4.3, restated). *For any optimal M-bounded-capacity learner A*, 1542

1543 *I.* if 
$$M \le M_0^-(\frac{r}{1-r} \cdot p)$$
, then  $\bar{\mathcal{L}}_1(\mathcal{A}) = F_{\mathcal{P}_1}(0)$ ;

1545 2. if 
$$M \ge M_0^+(\frac{r}{1-r} \cdot p) + H_{\text{tot}}$$
, then  $\bar{\mathcal{L}}_1(\mathcal{A}) = F_{\mathcal{P}_1}(\infty)$ .

1547 *Proof.* This is a direct consequence of Theorem D.8 by noting that  $(1) - D^- F_{\mathcal{P}_2}(M)$  is left continuous 1548 and non-increasing in M;  $(2) - D^+ F_{\mathcal{P}_2}(M)$  is right continuous and non-increasing in M;  $(3) F_{\mathcal{P}_2}(M)$ 1549 is almost everywhere differentiable.

1550 1551

1552

1553 1554

1555

1556

1558

1559

1560

1561

1562

1563

1564