

# TRAINING ON THE TEST TASK CONFOUNDS EVALUATION AND EMERGENCE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We study a fundamental problem in the evaluation of large language models that we call *training on the test task*. Unlike wrongful practices like training on the test data, leakage, or data contamination, training on the test task is not a malpractice. Rather, the term describes a growing set of techniques to include task-relevant data in the pretraining stage of a language model. We demonstrate that training on the test task confounds both relative model evaluations and claims about emergent capabilities. We argue that the seeming superiority of one model family over another may be explained by a different degree of training on the test task. To this end, we propose an effective method to adjust for the effect of training on the test task on benchmark evaluations. Put simply, to fine-tune each model under comparison on the same task-relevant data before evaluation. We then show that instances of emergent behavior disappear gradually as models train on the test task. Our work promotes a new perspective on the evaluation of large language models with broad implications for benchmarking and the study of emergent capabilities.

## 1 INTRODUCTION

The machine learning community has long recognized certain clear violations of the benchmarking protocol. Training on the test set is the most notorious among them (Duda & Hart, 1973; Hastie et al., 2017). Data leakage (Kapoor & Narayanan, 2022) and data contamination (Roberts et al., 2023; Jiang et al., 2024) are closely related problems linked to the rise of massive web-crawled training datasets. Researchers can all agree that test data should never appear in the training set.

But it’s been much less clear what to do about legitimate attempts to bring training closer to evaluation. There is an obvious a gap between next token prediction at training time and tasks, such as reasoning and question answering, at test time. Ongoing research and engineering efforts, in fact, aim to narrow precisely this gap (MetaAI, 2024). Why shouldn’t training be informed by knowledge about the downstream test tasks? What’s an unfair advantage of some may be the feature of others.

In this work, we group strategies to utilize task knowledge at training time under the umbrella term of *training on the test task*. Examples of training on the test task include the use of instruction-tuning data or question answering templates during pre-training (Bai et al., 2023; StabilityAI, 2023; Groeneveld et al., 2024; Zhang et al., 2024). Models may also implicitly train on the test task when their pretraining data is selected through benchmark ablations (Gemma et al., 2024; MetaAI, 2024). We work from the premise that training on the test task is acceptable—or at least, unavoidable.

In a nutshell, we show that training on the test task strongly confounds model comparisons across different scales and model families. Moreover, it significantly obscures the study of emergent capabilities. Rather than scrambling to detect and disallow various forms of training on the test task, we propose to “fight fire with fire”. We show that we can effectively level the playing field by giving each model the same, sufficient task-specific fine-tuning before evaluation. This adjustment restores cleaner log-linear scaling and makes capabilities predictable based on much smaller model scales.

### 1.1 OUR CONTRIBUTIONS

We introduce the term training on the test task to group a growing repertoire of practices that utilize knowledge about evaluation tasks at training time. We study its impact on present-day benchmark

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

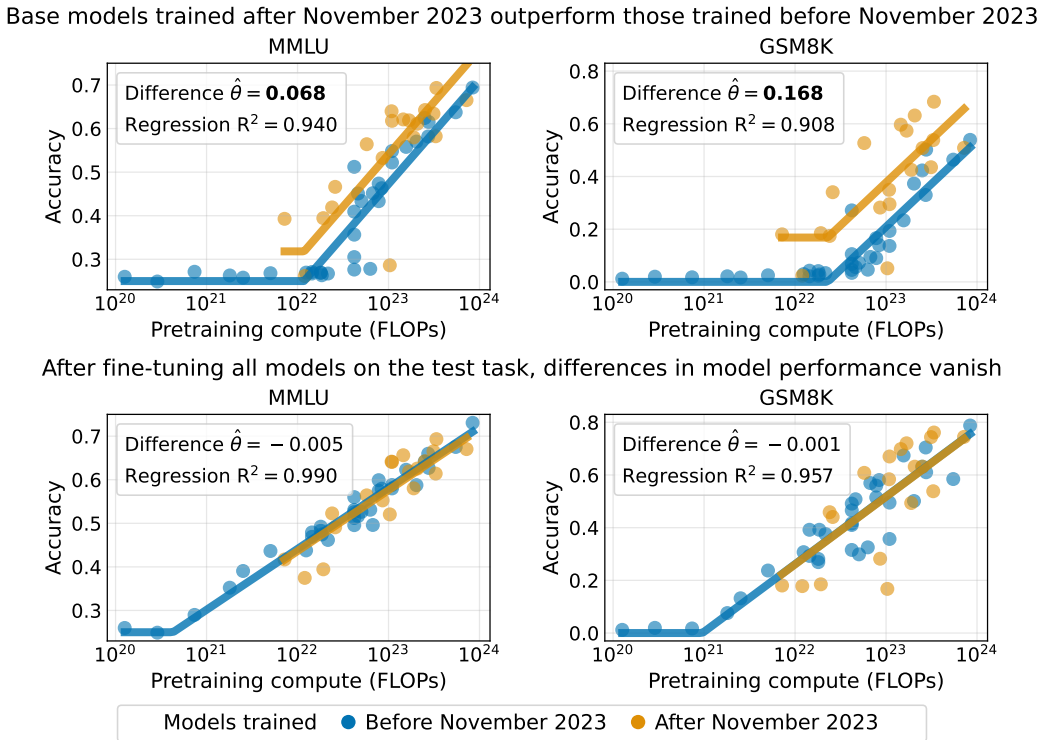


Figure 1: MMLU and GSM8K scores of 53 base models, with model sizes ranging from 70M to 70B parameters. Solid lines correspond to the regression fit of  $A = \alpha \max(0, \log C - c_e) + \theta N + r$ , where  $A$  is accuracy,  $C$  is pretraining compute,  $N$  is whether the model was trained after November 2023, and  $r$  is random chance accuracy. The coefficient  $\theta$  denotes the average improvement of models trained after November 2023 when controlling for pretraining compute. Bold indicates statistical significance with  $p$ -value  $< 0.05$ . (Top) We hypothesize that training on the test task confounds benchmark evaluations, resulting in newer base models substantially outperforming older ones. (Bottom) We propose to adjust for differences in test task training by fine-tuning all models on the same, sufficient amount of task-specific data before evaluation. After fine-tuning on the test task, differences in benchmark performance between older and newer models vanish.

evaluations by critically examining the performance improvements of recent language models. Our analysis spans 53 different language models and two major active benchmarks, MMLU and GSM8K.

We start in Section 2 by dividing models into those trained before November 2023 and those trained after. We find that for the same amount of pretraining compute, newer models strongly outperform older ones. We then fine-tune all models on the same amount of task-specific data before evaluation. After fine-tuning on the same task data, newer models no longer outperform older ones. Rather, their performance equalizes. See Figure 1. This outcome suggests that newer models outperform older ones on MMLU and GSM8K primarily because newer models trained more on the test task.

We propose a simple and effective method to adjust for the effect of training on the test task on benchmark evaluations. Put simply, to fine-tune each model on the same, sufficient amount of task-specific data before evaluation. To validate our method, we demonstrate its effectiveness in a controlled setting: we take the older models and fine-tune them on the test task. Remarkably this recreates the kind of performance differences observed between newer and older models. We then show that we can undo the advantage of the fine-tuned models over the other models by further fine-tuning all models on the test task (Section 3.1, Figure 3).

Next, we give evidence that training on the test task may be a more dominant factor in benchmark performance than data contamination. To argue this point, we consider ARC and HellaSwag. Here, at first there appears to be no sign of newer models have an unfair advantage over older models. But after reformulating these benchmarks as MMLU-style multiple choice question answering tasks,

we see the same confounded results as for MMLU (Section 3.2, Figure 4). This suggests that the improvements of newer models are not primarily because of memorization of specific testing data. Either way, our proposed adjustment recovers fair model comparisons.

Then, we show how training on the test task distorts model family comparisons. Certain model families appear markedly superior to others before adjusting for test task training but not after adjustment (Section 4.1). In fact, after adjustment, newer model families offer only modest improvements to the Pareto frontier of model performance relative to pre-training compute (Section 4.2).

Finally, we demonstrate that training on the test task has profound implications for the study of emergent capabilities. The phenomenon of emergence disappears gradually as the amount of training on the test task grows (Section 5). Specifically, we can make capabilities visible and predictable from much smaller model scales, recovering cleaner log linear-scaling.

Our work calls for a major reorientation of large language model evaluation. Model comparisons and claims of emergence are strongly confounded by the choice of training data relative to the test tasks. When comparing models with different pre-training data, our recommendation is to give each model the same sufficient amount of fine-tuning on task-relevant data prior to evaluation.

## 2 ADJUSTING FOR TRAINING ON THE TEST TASK

We choose MMLU (Hendrycks et al., 2020) and GSM8K (Cobbe et al., 2021) as a case study for investigating training on the test task in active benchmarks. MMLU tests for world knowledge, whereas GSM8K tests multi-step mathematical reasoning. These two benchmarks are arguably the two most prominent LLM benchmarks in recent times. We evaluate models using LM Evaluation Harness library (EleutherAI, 2024), in identical fashion to the HuggingFace leaderboard (5-shot). See Appendix C for results pertaining to the OpenLLM Leaderboard v2 (Fourrier et al., 2024a).

We evaluate 53 base models, ranging in size from 70M to 70B parameters. See Appendix A.1 for the full list. The HF leaderboard’s FAQ makes the distinction between “base pretrained models” and instruction-tuned or chat models, arguing that this is necessary to ensure fair model comparisons. We select models that are categorized as “pretrained”. We check that the technical report of each of the selected models makes no mention of the model being fine-tuned. We only consider models for which the number of training tokens is known. This allows us to estimate the total amount of pretraining compute in FLOPs as  $C \approx 6 \cdot N \cdot D$ , where  $C$  is pretraining compute,  $N$  is the number of model parameters, and  $D$  is the number of training tokens.

**Recent models outperform older ones given the same pretraining compute.** We evaluate models on MMLU and GSM8K, and plot benchmark accuracy against pretraining compute in Figure 1 top. We observe that performance correlates with pretraining compute for both benchmarks. However, on the surface it appears that more recent models better leverage pretraining compute. In other words, for a given compute budget newer models are able to attain better benchmark performance.

These improvements in benchmark performance coincide with the recent adoption of certain pre-training practices that may amount to training on the test task. For example, Qwen (Bai et al., 2023), and Olmo 1.7 (Groeneveld et al., 2024) include instruction data during pretraining. StableLM 2 (StabilityAI, 2023) reformulates some of its pretraining datasets to better resemble downstream tasks such as question-answering. More subtly, the pretraining data mixtures of Gemma (Gemma et al., 2024) and Llama 3 (MetaAI, 2024) were determined through ablations on benchmark evaluations.

This raises an important question: Do newer models outperform older ones mainly because newer models trained more on the test task? At first sight, an answer seems elusive. After all, the pretraining data of most recent models is not public. Retraining all model families with the same training data and compute budget would be both infeasible and cost prohibitive. Nevertheless, in the next section, we propose a way to get at the answer by adjusting for the effect of training on the test task.

### 2.1 ADJUSTING FOR TRAINING ON THE TEST TASK BY TRAINING ON THE TEST TASK

We propose to adjust for differences in test task training by fine-tuning all models on the same, sufficient amount of task-specific data before evaluation. To do so, we need a source of task-specific data for each of the tasks we consider. For multiple choice questioning answering, we use the

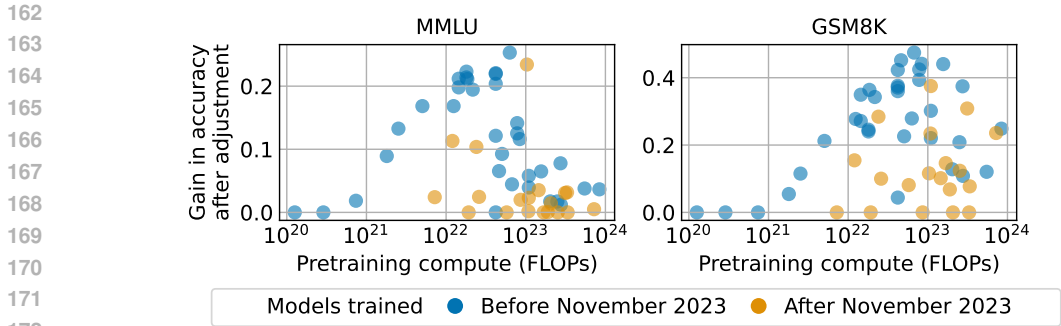


Figure 2: Older models tend to benefit much more from fine-tuning on task data.

auxiliary training set accompanying the HF MMLU repository<sup>1</sup>. It contains around 100,000 training examples and 30M tokens. For mathematical reasoning, we combine MetaMathQA (Yu et al., 2023) and Orca-Math (Mitra et al., 2024), totalling 600,000 training examples and 200M tokens. We fine-tune models for three epochs using standard hyperparameter choices, see Appendix A.2. The amount of compute required for fine-tuning is minimal compared to models’ pretraining compute.

We plot model scores on MMLU and GSM8K after fine-tuning in Figure 1 (bottom). We observe that after fine-tuning on task relevant data, both newer and older models follow remarkably similar scaling trends. That is, newer models no longer appear to outperform newer models.

Remarkably, we observe that older models tend to benefit much more from fine-tuning on task-relevant data, see Figure 2. The improvements in older models are striking, often leaping from random chance accuracy to double-digit gains in accuracy. In contrast, fine-tuning provides comparatively little benefit to newer models. This observation suggests that newer models have already been exposed to a substantial amount of task-relevant data, making additional fine-tuning less impactful.

## 2.2 QUANTIFYING PERFORMANCE DIFFERENCES BETWEEN NEWER AND OLDER MODELS

We draw inspiration from scaling laws (Kaplan et al., 2020) in how we model benchmark accuracy  $A$  to scale log-linearly with pretraining compute  $C$ . To account for emergence (Wei et al., 2022), we assume that models perform at the task’s random chance accuracy  $r$  up to scaling to some point of emergence  $c_e$ . We let the variable  $N$  denote whether a model was trained after November 2023, and regress the model

$$A = \alpha \max(0, \log C - c_e) + \theta N + r + \epsilon, \tag{1}$$

where  $\alpha$ ,  $\theta$  and  $c_e$  are the fit’s parameters, and  $\epsilon$  is random noise. We focus on the coefficient  $\theta$ , which corresponds to the average difference in benchmark performance between newer and older models after controlling for pretraining compute. We fit the model in Equation 1, and report the regression coefficient  $\theta$  in Figure 1. We obtain  $R^2 > 0.9$  for all model fits.

Before adjusting for test task training, the estimated difference in performance  $\hat{\theta}$  between newer and older models are statistically significant, positive, and large. Specifically, recent models outperform older ones on average by over 7 accuracy points in MMLU and 17 accuracy points in GSM8K. These are remarkably large differences in benchmark performance. However, after the adjustment, the estimated coefficient  $\hat{\theta}$  is both small and not statistically significant. See Figure 1 bottom. That is, conditioned on all models training on the same amount of task-specific data, we find no evidence for a significant difference in benchmark performance between newer and older models.

Therefore, the performance of newer and older models equalizes when all models are exposed to the same amount of task-relevant data. This suggests that the impressive benchmark improvements of newer models are primarily attributable to newer models training more on the test task. We present a causal interpretation of results in Appendix B, outlying the assumptions necessary to establish a causal link between training on the test task and the benchmark improvements of newer models.

<sup>1</sup><https://huggingface.co/datasets/cais/mmlu>. This training set, far from being an i.i.d split of MMLU, compiles the training sets of other multiple-choice benchmarks, such as ARC (Clark et al., 2018), MCTest (Richardson et al., 2013), OpenBookQA (Mihaylov et al., 2018), and RACE (Lai et al., 2017)

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

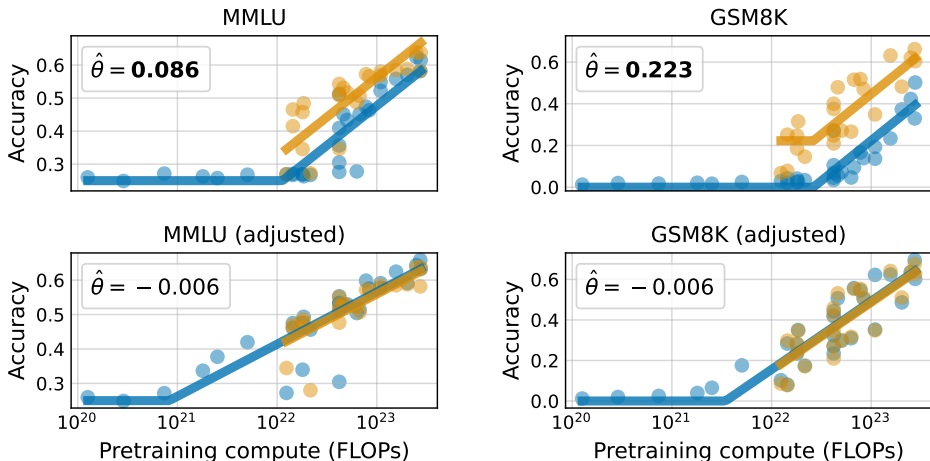


Figure 3: Models trained before November 2023 (●) without fine-tuning and (●) after fine-tuning on the test task. Their difference in benchmark performance  $\hat{\theta}$  resemble that between newer and older models (*top*). After adjusting by training on the test task, their difference vanishes (*bottom*).

### 3 RECREATING DIFFERENCES IN BENCHMARK PERFORMANCE

We have so far established that newer models strongly outperform older models for the same amount of pre-training compute. We now demonstrate how to recreate such differences in performance by actively manipulating how much models train on the test task. We do so in two ways. First, we fine-tune older models on task relevant data (Section 3.1). Second, we reformulate certain test tasks to use MMLU-style multiple choice prompts instead of “cloze” evaluations (Section 3.2). Both experiments recreate the kind of performance differences observed between newer and older models.

These results provides further evidence that the differences in performance between older and newer models are linked to test task training. They also demonstrates how test task training distorts benchmark evaluations. Fortunately, in both cases, we show that fine-tuning models on task-relevant data before evaluation is an effective mechanism for mitigating the bias introduced by training on the test task. In doing so, we systematically validate the proposed adjustment method.

#### 3.1 FINE-TUNING ON THE TEST TASK

For this section, we only consider models trained before November 2023. We split models into two cohorts: a control group and a treatment group. We take the older models as the control group. We then create a treatment group by fine-tuning the control group on the datasets of task-relevant data introduced in Section 2. We only fine-tune models with at least  $7 \cdot 10^{21}$  FLOPs, the pre-training compute of the smallest newer model, Qwen 1.5 0.5B. We fine-tune on each dataset independently, for a single epoch. We plot in Figure 3 top the benchmark performance of the two cohorts.

Qualitatively, the performance differences between the control and treatment groups mirror those observed earlier between newer and older models, contrast Figure 3 with Figure 1 Quantitatively, the estimated performance gain  $\hat{\theta}$  from fine-tuning is similar to the difference between newer and older models estimated in Section 2.2. That is, fine-tuning older models on the test task produces both qualitatively and quantitatively similar confounding to that observed between newer and older models. This results further supports our hypothesis that newer models are largely equivalent to older models that have trained on the test task. Furthermore, it demonstrates the large effect that training on the test task can have on benchmark evaluations.

We then apply our proposed adjustment by further fine-tuning both the control and treatment groups on the test task, see Figure 3 bottom. After the adjustment, the estimated difference in performance  $\hat{\theta}$  between the control and treatment group is both small and not statistically significant. We therefore validate a vital soundness property: after deliberately training some models on the test task, we can undo their advantage over other models by further training all models on the test task.

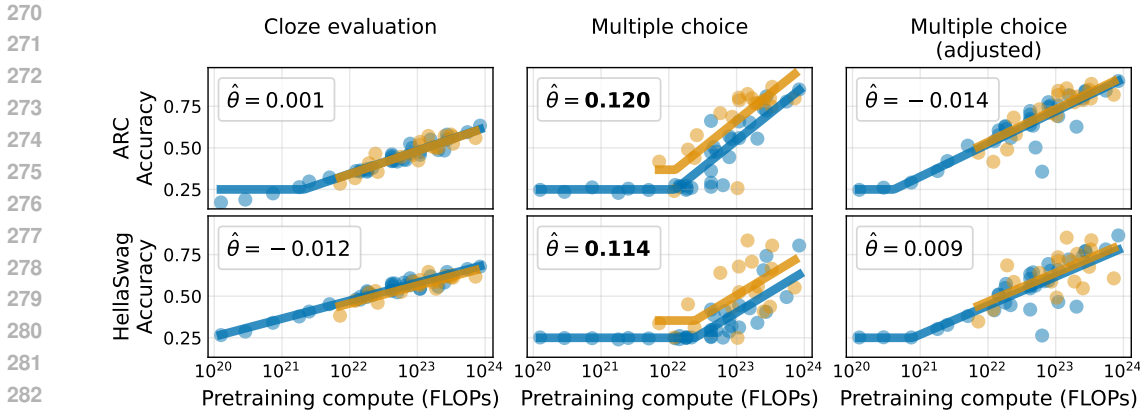


Figure 4: Reformulating ARC and HellaSwag as MMLU-style questions give rise to large differences  $\hat{\theta}$  between models trained (●) before November 2023 and (●) after November 2023 (*center*). After adjusting by fine-tuning on the test task, differences in performance vanish (*right*).

### 3.2 REFORMULATING THE TEST TASK

In this section, we show that reformulating other benchmarks as multiple-choice question answering tasks leads to similar differences in performance between older and newer models. We consider two additional benchmarks from the HF leaderboard v1: ARC Challenge (Clark et al., 2018) and HellaSwag (Zellers et al., 2019). Similarly to MMLU, ARC is comprised of grade-school level questions. HellaSwag instead tests for commonsense reasoning. Like MMLU, the questions in ARC and HellaSwag are accompanied by four possible answers. ARC and HellaSwag use “cloze” evaluations: a models’ answer is taken to be that with the largest completion likelihood given the input question. In contrast, MMLU formulates questions as multiple-choice: all four answer choices are listed, and the model is prompted to pick one of the answer choices.

We first evaluate on ARC and HellaSwag using the standard cloze evaluation, and plot their benchmark performance in Figure 4 left. We repeat the statistical analysis of Section 2.2. We find that the estimated difference in performance  $\hat{\theta}$  between newer and older models is small and not statistically significant. That is, newer models do not outperform older models on ARC and HellaSwag.

We then reformulate ARC and HellaSwag as MMLU-style multiple-choice questions, and plot the resulting benchmark performance in Figure 4 center. We observe large differences in performance between newer and older models. Specifically, we find the difference in performance  $\hat{\theta}$  between newer and older models to be significant, positive, and large, and to be roughly similar in magnitude to that estimated for MMLU in Section 2.2. That is, reformulating the test task as multiple choice question answering leads to similar confounding to that observed for MMLU. Therefore, newer models overperform on MMLU likely not because of memorization of specific testing data (i.e., due to data contamination), but rather due to an improved ability for multiple-choice question answering.

Lastly, we adjust for test task training by fine-tuning all models on the MMLU auxiliary training set, and plot their ARC Challenge and HellaSwag scores in Figure 4 right. We no longer find evidence of a large nor a significant difference in performance between newer and older models. Therefore, the proposed adjustment is effective in mitigating the bias introduced by evaluating models using multiple-choice question answering tasks. Notably, we achieve this using the same MMLU auxiliary training set, thus demonstrating that the adjustment data need not closely resemble the test data.

**What does MMLU test for?** We evaluate MMLU using the cloze methodology instead of the usual multiple-choice prompts. We plot the results in Figure 5 center. With cloze evaluations, the difference in performance between newer and older models is both small and not statistically significant. This suggests that the standard MMLU evaluation conflates knowledge-testing with testing a models’ ability to answer multiple choice questions. Newer models therefore attain higher MMLU scores than older models largely because they are better at multiple-choice question answering, and not because they necessarily “know more”.



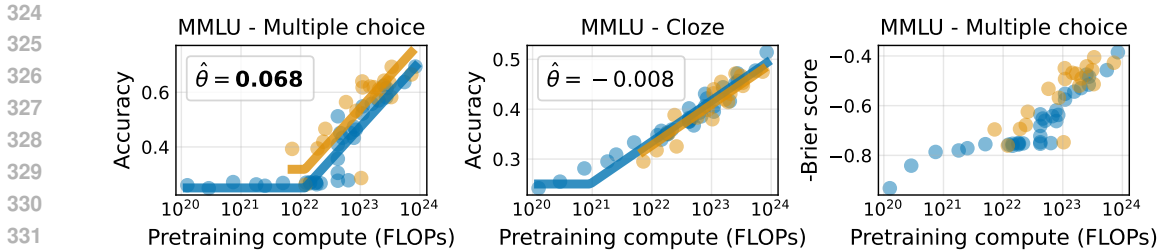


Figure 5: When evaluating MMLU using “cloze” prompts, models trained (●) after November 2023 no longer outperform those trained (●) before November 2023 (*middle*). When using Brier score as the metric, we still observe sharp improvements in performance around  $10^{22}$  FLOPs (*right*).

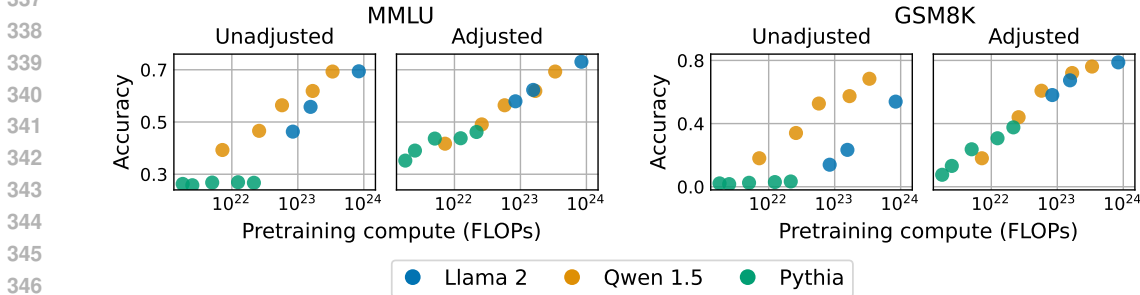


Figure 6: Training on the test task confounds relative comparisons between model families. After adjustment, none of the three model families appears to be superior beyond their compute.

## 4 IMPLICATIONS FOR MODEL COMPARISONS

So far, we have shown how training on the test task distorts benchmark evaluations. Next, we examine its impact on the relative comparison of model families (Section 4.1), as well as its implications for accurately measuring progress in model capabilities over time (Section 4.2).

### 4.1 COMPARING MODEL FAMILIES

We compare the performance of the Pythia, Llama 2, and Qwen 1.5 model families, which likely train on the test task to very different extents. Pythia was trained on the Pile (Gao et al., 2020), a collection of curated datasets that are unlikely to contain much test task data. Llama 2 was trained mostly on web data, which is reasonable to assume may contain more test task data. Lastly, Qwen 1.5 explicitly pre-trains on instruction data, thus likely training on the test task to a large extent.

We plot the MMLU and GSM8K scores of the three model families in Figure 6, as well as their adjusted scores (i.e., after fine-tuning on task relevant data). Without adjustment, Qwen 1.5 appears to be the superior model family: it Pareto dominates both the Llama 2 and Pythia models. In contrast, all Pythia models perform no better than random chance, making it unclear whether scaling Pythia offers any benefit at all. After adjustment, however, all three model families exhibit remarkably similar scaling trends. Therefore, after correcting for the confounding introduced by test task training, none of the model families appears superior to the others beyond their pre-training compute.

Training on the test task therefore profoundly confounds relative model comparisons. Base models are rarely used “as is” and are generally adapted in some way before deployment. Because of the confounding of training on the test task, performance before adaptation may not reliably predict performance after adaptation. It therefore makes little sense to compare base models at face value. For example, Llama 2 70B underperforms on GSM8K compared to Qwen 1.5 7B and StableLM 2 12B, but is the top-performing model after all models are adapted for mathematical reasoning.

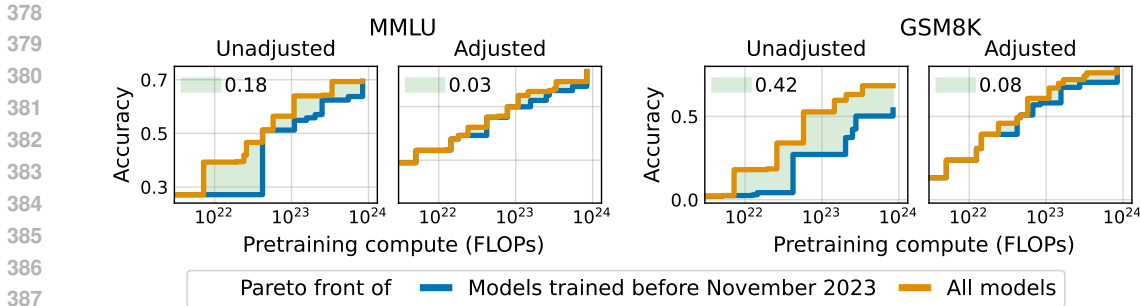


Figure 7: Training on the test task overestimates the improvements in performance-per-compute of recent models. After adjustment, the area of improvement (green) reduces by a sixfold.

#### 4.2 PROGRESS IN MODEL CAPABILITIES

Training on the test task substantially overestimates the progress in capabilities per unit of compute achieved by recent model families. In Figure 7 we plot the Pareto frontier of benchmark accuracy against pretraining compute, both for models trained before November 2023 and for all models. We measure progress by considering the area of improvement of the Pareto frontier since November 2023, shaded in green. Without adjustment, the difference between the two Pareto frontiers is large for both MMLU and GSM8K, indicating substantial progress since November 2023. After adjustment, however, the area of improvement reduces by a sixfold, showing only modest improvements.

On the other hand, recent models tend to be trained on more data than Chinchilla compute-optimal (Hoffmann et al., 2022). Given the Chinchilla scaling laws, it is remarkable that newer, smaller models match the performance of older, larger ones for the same amount of pretraining compute. For example, we find that Llama 3 8B closely matches the performance of Llama 2 70B.

### 5 IMPLICATIONS FOR EMERGENCE

Throughout our evaluations, we observe emergent behaviour for MMLU and GSM8K: models perform at near random chance up to a certain scale of pretraining compute, followed by relatively sharper improvements in performance at larger scale. After training on the test task, however, emergence for MMLU and GSM8K appears to occur at substantially lower scales. We dedicate this section to more closely investigate the relationship between training on the test task and emergence.

**Emergence arises at lower scales with increased test task training.** We consider only models trained before November 2023, as we have established that these models train on the test task to a lesser extent. We evaluate the models at intermediate checkpoints as they train on the datasets of task relevant data introduced in Section 2.1. We fit  $\alpha$  and  $c_e$  in Equation 1 to the different intermediate checkpoints, and report in Figure 8 top the corresponding points of emergence  $c_e$ . We find that emergence arises at increasingly lower compute regimes as models train on the test task. That is, the performance of models after training on the test task is predictable at substantially lower scales. For instance, for MMLU models exhibit emergence at around  $10^{22}$  FLOPs, roughly the scale of Pythia 6.9B. After training on 64,000 examples, emergence arises around around  $6 \cdot 10^{20}$  FLOPs, that is, roughly the scale of Pythia 410M. We find similar results for GSM8K, see Figure 17 in Appendix D.

**Training on the test task yields increasingly better log-linear fits.** The log-linear relationship between pretraining loss and compute is well-established (Kaplan et al., 2020). We observe that, for the compute ranges that we consider, training on the test task increasingly recovers log-linear scaling between pretraining compute and benchmark accuracy. Similarly to the earlier section, we evaluate intermediate checkpoints but instead fit log-linear functions in Figure 8 bottom. We observe that the  $R^2$  of the fit improves substantially as the models train on more task-relevant data. For MMLU, the  $R^2$  value jumps from 0.63 to 0.95 after training on 64,000 examples. Therefore, after training on the test task almost all of the variation in benchmark accuracy can be explained by log-linear scaling of pre-training compute. We find similar results for GSM8K, see Figure 17 in Appendix D.



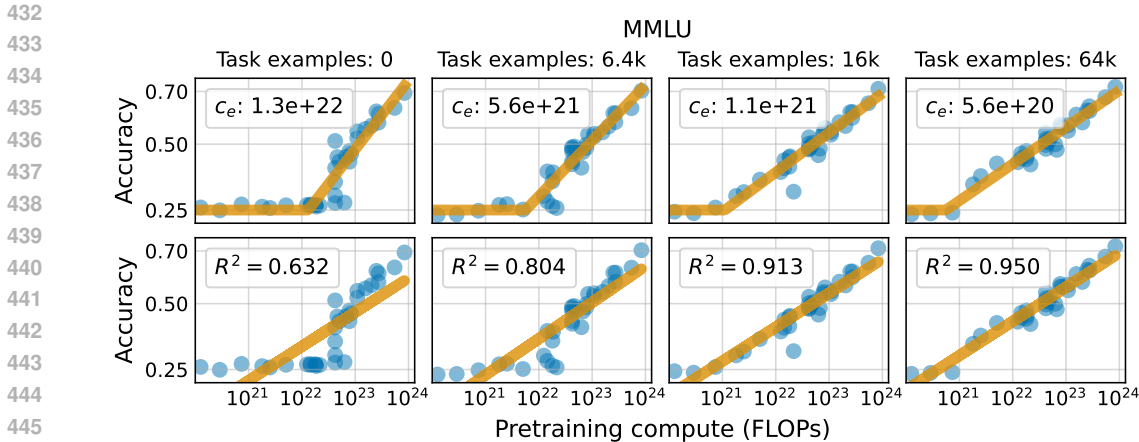


Figure 8: Scaling on MMLU as models increasingly train on the test task. The point of emergence  $c_e$  arises at lower scales (*top*). Training on the test task yields cleaner log-linear scaling fits (*bottom*).

**Recommendations.** Schaeffer et al. (2024a) argue that emergence appears due to the choice of metric. To mitigate emergence, they suggest to consider Brier score instead of accuracy. We observe, however, that emergence for MMLU does not disappear when using the Brier score. We discuss two practical solutions to obtain predictive scaling while maintaining accuracy as the evaluation metric.

For MMLU and multiple-choice benchmarks more broadly, cloze evaluations consistently yield smoother and more predictable scaling even when using accuracy as the evaluation metric. Since the purpose of these benchmarks is knowledge-testing more so than testing multiple-choice answering ability, cloze evaluations should be preferable insofar predictive scaling is an important consideration. Our recommendation aligns with the concurrent work by Gu et al. (2024).

More broadly, if sufficient task relevant data is available, then training on the test task can result in much more predictable scaling by shifting emergence to smaller compute scales. That is, by considering “adjusted” scaling laws where models across scales are fine-tuned on the same, sufficient task-relevant data prior to evaluation. Such scaling laws potentially correspond to those of “specialist” models, which for some domains –such as the legal domain (Dominguez-Olmedo et al., 2024)– or purposes –e.g., safety– might be preferable to the scaling laws of generalist models.

## 6 RELATED WORK

Benchmarks have played a central role in both machine learning (Hardt & Recht, 2022) and natural language processing (Storks et al., 2019). Classically, benchmarks comprised both a test set and a reasonably large training set (LeCun et al., 1998; Deng et al., 2009). Models were trained on the same training set, and then evaluated on the accompanying test set. The success of unsupervised language modelling (Peters et al., 2018; Kenton & Toutanova, 2019; Radford et al., 2019), however, has changed this paradigm. Firstly, present-day language models differ in their training data, which is not standardized but rather treated as a design choice (Raffel et al., 2020; Albalak et al., 2024; Li et al., 2024). Secondly, language models are a priori not trained with the explicit objective of maximizing any single benchmark score. Rather, language models are expected to be able to perform a broad range of tasks (Wang et al., 2018; Brown et al., 2020).

**Data contamination.** Data contamination or test-set contamination refers to any overlap between the training and the test data such that test results overestimate a model’s generalization performance. The scale and often little curation of present-day pretraining corpora exacerbates data contamination concerns in language model evaluations (Sainz et al., 2023; Magar & Schwartz, 2022; Jiang et al., 2024). However, detecting and preventing data contamination is currently an open problem (Yang et al., 2023b; Golchin & Surdeanu, 2023). Roberts et al. (2023) and Li & Flanigan (2024) find that models often perform better on datasets that were publicly available during model training. While all models that we consider were released at least a year and a half after MMLU and GSM8K, we

486 nonetheless find that, controlling for compute, more recent models perform better. These perfor-  
 487 mance gains are unlikely to be driven solely by test set leakage and require additional explanation.  
 488 In Section 3.2, we find evidence that that training on the test task may be a more dominant factor in  
 489 benchmark performance than data contamination.  
 490

491 **Training on the test task.** The effectiveness of fine-tuning on the training set accompanying LLM  
 492 benchmarks is well-known (Wei et al., 2021; Wang et al., 2022; Chung et al., 2024). Consequently,  
 493 many influential instruction-tuning datasets contain or are partly derived from benchmark train  
 494 data (Wei et al., 2021; Honovich et al., 2022; Mukherjee et al., 2023). Li & Flanigan (2024) identify  
 495 small amounts of benchmark-specific data in the publicly available Alpaca (Taori et al., 2023) and  
 496 Vicuna (Chiang et al., 2023) instruction-tuning sets. Zhou et al. (2023b) empirically analyze the ef-  
 497 fects of fine-tuning on benchmark-specific data and warn about its impacts on benchmark validity. In  
 498 contrast, we find evidence for training on the test task without the need for explicitly identifying spe-  
 499 cific data points used at training time, or modifying tasks. In addition, our proposed method of fine-  
 500 tuning on task data prior to evaluation allows us to correct for its effect on benchmark performance.

501 **Emergent abilities of language models.** Emergent capabilities (Wei et al., 2022; Ganguli et al.,  
 502 2022) refer to levels of model performance at large scales that cannot be easily predicted by extrap-  
 503 olating from smaller scales. Wei et al. (2022) report emergent capabilities for various benchmarks  
 504 including MMLU and GSM8K (Srivastava et al., 2022). However, Srivastava et al. (2022); Schaeffer  
 505 et al. (2024b) find that the log-probability of the correct answer often improves smoothly, even when  
 506 other metrics seem to show emergence. Schaeffer et al. (2024a) argue that emergent capabilities are  
 507 mostly an artifact of non-linear and discontinuous evaluation metrics like accuracy. In contrast, we  
 508 find signs of emergence on MMLU even when using continuous metrics like the Brier score. We  
 509 additionally show that increasingly fine-tuning on the test task yields more predictive scaling by  
 510 shifting the point of emergence to smaller compute scales, recovering cleaner log-linear scaling.

## 511 7 DISCUSSION

512 The 1968 Olympics took place in Mexico City at the significant altitude of 2340 meters, higher than  
 513 Australia’s tallest peak. Runners who had trained at altitude in their home countries were better  
 514 prepared to compete in Mexico City’s conditions, as it turned out. But the hotly debated results  
 515 of the Games did not lead the organizers to prohibit training at natural altitude. Instead, they let  
 516 everyone do it; and athletes came to consider altitude training an excellent way to train.  
 517

518 The anecdote holds a lesson for the evaluation of large language models half a century later. Knowl-  
 519 edge about the evaluation conditions necessarily influences training practices under competitive  
 520 pressure. It may be a fool’s errand to prohibit the practice. Instead, we propose to adjust for it  
 521 by giving every model the same task-specific preparation before evaluation. We work from the as-  
 522 sumption that training on the test task, in general, cannot be effectively detected, disallowed, or  
 523 disincentivized. Detecting what training data a model has seen is a notoriously difficult problem.  
 524 Researchers routinely acknowledge the futility of fighting data contamination. Moreover, we antic-  
 525 ipate that the ways to effectively train on the test task will only grow in scope and adoption.  
 526

527 Our work demonstrates that comparisons of different models are confounded by the choice of  
 528 training data and training practices. Different model families vary in the degree that they were—  
 529 implicitly or explicitly—trained on various test tasks. It therefore makes little sense to compare  
 530 model performance at face value without accounting for how the training data relate to the test task.  
 531 A small amount of task data can have a disproportional large effect on benchmark performance.

532 We can apply the same principles to emergent behavior. After training on the test task, model  
 533 capabilities become predictable at smaller model size and grow continuously with scale. Training  
 534 on the test task greatly reduces the unpredictability associated with emergence, notably without any  
 535 change in the metric, thus largely disarming the ominous nature of emergence.

536 Despite the daunting challenges that training on the test task poses for the fair evaluation of language  
 537 models, it’s also its own best remedy. Giving each model the same sufficient task-specific fine-  
 538 tuning harmonizes model comparisons and linearizes the relationship between model capabilities  
 539 and pretraining compute. We hope that our work informs stronger evaluation standards that address  
 central challenges in the current evaluation ecosystem.

## REFERENCES

- 540  
541  
542 Alon Albalak, Yanai Elazar, Sang Michael Xie, Shayne Longpre, Nathan Lambert, Xinyi Wang,  
543 Niklas Muennighoff, Bairu Hou, Liangming Pan, Haewon Jeong, et al. A survey on data selection  
544 for language models. *arXiv preprint arXiv:2402.16827*, 2024.
- 545 Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Co-  
546 jocar, Mérouane Debbah, Étienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic,  
547 et al. The falcon series of open language models. *arXiv preprint arXiv:2311.16867*, 2023.
- 548 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,  
549 Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- 551 Marco Bellagente, Jonathan Tow, Dakota Mahan, Duy Phung, Maksym Zhuravinskyi, Reshith  
552 Adithyan, James Baicoianu, Ben Brooks, Nathan Cooper, Ashish Datta, et al. Stable lm 2 1.6 b  
553 technical report. *arXiv preprint arXiv:2402.17834*, 2024.
- 554 Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric  
555 Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al.  
556 Pythia: A suite for analyzing large language models across training and scaling. In *International  
557 Conference on Machine Learning*, pp. 2397–2430. PMLR, 2023.
- 558 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
559 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
560 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 561 Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui  
562 Chen, Zhi Chen, Pei Chu, et al. Internlm2 technical report. *arXiv preprint arXiv:2403.17297*,  
563 2024.
- 564 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng,  
565 Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot  
566 impressing gpt-4 with 90%\* chatgpt quality. *See https://vicuna.lmsys.org (accessed 14 April  
567 2023)*, 2(3):6, 2023.
- 568 Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li,  
569 Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned lan-  
570 guage models. *Journal of Machine Learning Research*, 25(70):1–53, 2024.
- 571 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and  
572 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.  
573 *arXiv preprint arXiv:1803.05457*, 2018.
- 574 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
575 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to  
576 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 577 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hier-  
578 archical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*,  
579 pp. 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- 580 Ricardo Dominguez-Olmedo, Vedant Nanda, Rediet Abebe, Stefan Bechtold, Christoph Engel, Jens  
581 Frankenreiter, Krishna Gummadi, Moritz Hardt, and Michael Livermore. Lawma: The power of  
582 specialization for legal tasks. *arXiv preprint arXiv:2407.16615*, 2024.
- 583 Richard O. Duda and Peter E. Hart. *Pattern Classification and Scene Analysis*. Wiley New York,  
584 1973.
- 585 EleutherAI. Language model evaluation harness. [https://github.com/EleutherAI/  
586 lm-evaluation-harness](https://github.com/EleutherAI/lm-evaluation-harness), 2024. Accessed: 2024-05-20.
- 587 Clémentine Fourrier, Nathan Habib, Alina Lozovskaya, Konrad Szafer, and Thomas Wolf. Open  
588 llm leaderboard v2. [https://huggingface.co/spaces/open-llm-leaderboard/  
589 open\\_llm\\_leaderboard](https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard), 2024a. Accessed: 2024-07-08.

- 594 Clémentine Fourrier, Nathan Habib, Alina Lozovskaya, Konrad Szafer, and Thomas Wolf. Performances are plateauing, let's make the leaderboard steep again. <https://huggingface.co/spaces/open-llm-leaderboard/blog>, 2024b. Accessed: 2024-07-08.
- 595  
596  
597
- 598 Ruyi Gan, Ziwei Wu, Renliang Sun, Junyu Lu, Xiaojun Wu, Dixiang Zhang, Kunhao Pan, Ping  
599 Yang, Qi Yang, Jiaying Zhang, et al. Ziya2: Data-centric learning is all llms need. *arXiv preprint*  
600 *arXiv:2311.03301*, 2023.
- 601  
602 Deep Ganguli, Danny Hernandez, Liane Lovitt, Amanda Askell, Yuntao Bai, Anna Chen, Tom  
603 Conerly, Nova Dassarma, Dawn Drain, Nelson Elhage, et al. Predictability and surprise in large  
604 generative models. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and*  
605 *Transparency*, pp. 1747–1764, 2022.
- 606  
607 Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason  
608 Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. The Pile:  
609 An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- 610  
611 Team Gemma, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya  
612 Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open  
613 models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
- 614  
615 Shahriar Golchin and Mihai Surdeanu. Time travel in LLMs: Tracing data contamination in large  
616 language models. *arXiv preprint arXiv:2308.08493*, 2023.
- 617  
618 Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord,  
619 Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson,  
620 Russell Authur, Khyathi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack  
621 Hessel, Tushar Khot, William Merrill, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik,  
622 Crystal Nam, Matthew E. Peters, Valentina Pyatkin, Abhilasha Ravichander, Dustin Schwenk,  
623 Saurabh Shah, Will Smith, Nishant Subramani, Mitchell Wortsman, Pradeep Dasigi, Nathan Lambert,  
624 Kyle Richardson, Jesse Dodge, Kyle Lo, Luca Soldaini, Noah A. Smith, and Hannaneh  
625 Hajishirzi. Olmo: Accelerating the science of language models. *Preprint*, 2024.
- 626  
627 Yuling Gu, Oyvind Tafjord, Bailey Kuehl, Dany Haddad, Jesse Dodge, and Hannaneh Hajishirzi.  
628 Olmes: A standard for language model evaluations. *arXiv preprint arXiv:2406.08446*, 2024.
- 629  
630 Moritz Hardt and Benjamin Recht. *Patterns, predictions, and actions: Foundations of machine*  
631 *learning*. Princeton University Press, 2022.
- 632  
633 Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning: Data*  
634 *Mining, Inference, and Prediction (Corrected 12th printing)*. Springer, 2017.
- 635  
636 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob  
637 Steinhardt. Measuring massive multitask language understanding. In *International Conference*  
638 *on Learning Representations*, 2020.
- 639  
640 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn  
641 Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset.  
642 In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks*  
643 *Track*, 2021.
- 644  
645 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza  
646 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training  
647 compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- 648  
649 Or Honovich, Thomas Scialom, Omer Levy, and Timo Schick. Unnatural instructions: Tuning  
650 language models with (almost) no human labor. *arXiv preprint arXiv:2212.09689*, 2022.
- 651  
652 Team InternLM. Internlm: A multilingual language model with progressively enhanced capabilities,  
653 2023.

- 648 Minhao Jiang, Ken Liu, Ming Zhong, Rylan Schaeffer, Siru Ouyang, Jiawei Han, and Sanmi Koyejo.  
649 Does data contamination make a difference? insights from intentionally contaminating pre-  
650 training data for language models. In *ICLR 2024 Workshop on Navigating and Addressing Data  
651 Problems for Foundation Models*, 2024.
- 652 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child,  
653 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language  
654 models. *arXiv preprint arXiv:2001.08361*, 2020.
- 655 Sayash Kapoor and Arvind Narayanan. Leakage and the reproducibility crisis in ml-based science,  
656 2022.
- 657 Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. Bert: Pre-training of deep  
658 bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pp. 4171–  
659 4186, 2019.
- 660 Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. Race: Large-scale reading  
661 comprehension dataset from examinations. In *Proceedings of the 2017 Conference on Empirical  
662 Methods in Natural Language Processing*, pp. 785–794, 2017.
- 663 Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database. *ATT Labs [Online].  
664 Available: <http://yann.lecun.com/exdb/mnist>*, 2, 1998.
- 665 Changmao Li and Jeffrey Flanigan. Task contamination: Language models may not be few-shot any-  
666 more. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 18471–  
667 18480, 2024.
- 668 Jeffrey Li, Alex Fang, Georgios Smyrnis, Maor Ivgi, Matt Jordan, Samir Gadre, Hritik Bansal, Etash  
669 Guha, Sedrick Keh, Kushal Arora, et al. Datacomp-lm: In search of the next generation of training  
670 sets for language models. *arXiv preprint arXiv:2406.11794*, 2024.
- 671 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Confer-  
672 ence on Learning Representations*, 2018.
- 673 Inbal Magar and Roy Schwartz. Data contamination: From memorization to exploitation. *arXiv  
674 preprint arXiv:2203.08242*, 2022.
- 675 MetaAI. Llama 3: Advancing open foundation models, 2024. URL [https://ai.meta.com/  
676 blog/meta-llama-3/](https://ai.meta.com/blog/meta-llama-3/).
- 677 Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct elec-  
678 tricity? a new dataset for open book question answering. In *Proceedings of the 2018 Conference  
679 on Empirical Methods in Natural Language Processing*, pp. 2381–2391, 2018.
- 680 Arindam Mitra, Hamed Khanpour, Corby Rosset, and Ahmed Awadallah. Orca-math: Unlocking  
681 the potential of slms in grade school math. *arXiv preprint arXiv:2402.14830*, 2024.
- 682 Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and  
683 Ahmed Awadallah. Orca: Progressive learning from complex explanation traces of gpt-4. *arXiv  
684 preprint arXiv:2306.02707*, 2023.
- 685 OpenLlama. Openllama, 2023. URL [https://github.com/openlm-research/open\\_  
686 llama](https://github.com/openlm-research/open_llama).
- 687 Judea Pearl. *Causality*. Cambridge university press, 2009.
- 688 Judea Pearl. Linear models: A useful “microscope” for causal analysis. *Journal of Causal Inference*,  
689 1(1):155–170, 2013.
- 690 Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and  
691 Luke Zettlemoyer. Deep contextualized word representations. *NAACL*, 2018.
- 692 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language  
693 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.

- 702 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi  
703 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text  
704 transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- 705  
706 Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: Memory optimizations  
707 toward training trillion parameter models. In *SC20: International Conference for High Perform-*  
708 *ance Computing, Networking, Storage and Analysis*, pp. 1–16. IEEE, 2020.
- 709 David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Di-  
710 rani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a bench-  
711 mark. *arXiv preprint arXiv:2311.12022*, 2023.
- 712  
713 Matthew Richardson, Christopher JC Burges, and Erin Renshaw. Mctest: A challenge dataset for the  
714 open-domain machine comprehension of text. In *Proceedings of the 2013 conference on empirical*  
715 *methods in natural language processing*, pp. 193–203, 2013.
- 716 Manley Roberts, Himanshu Thakur, Christine Herlihy, Colin White, and Samuel Dooley. Data  
717 contamination through the lens of time. *arXiv preprint arXiv:2310.10628*, 2023.
- 718  
719 Oscar Sainz, Jon Ander Campos, Iker García-Ferrero, Julen Etxaniz, Oier Lopez de Lacalle, and  
720 Eneko Agirre. NLP evaluation in trouble: On the need to measure LLM data contamination for  
721 each benchmark. *arXiv preprint arXiv:2310.18018*, 2023.
- 722 Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. Are emergent abilities of large language  
723 models a mirage? *Advances in Neural Information Processing Systems*, 36, 2024a.
- 724  
725 Rylan Schaeffer, Hailey Schoelkopf, Brando Miranda, Gabriel Mukobi, Varun Madan, Adam  
726 Ibrahim, Herbie Bradley, Stella Biderman, and Sanmi Koyejo. Why has predicting downstream  
727 capabilities of frontier ai models with scale remained elusive? *arXiv preprint arXiv:2406.04391*,  
728 2024b.
- 729 Zayne Rea Sprague, Xi Ye, Kaj Bostrom, Swarat Chaudhuri, and Greg Durrett. Musr: Testing the  
730 limits of chain-of-thought with multistep soft reasoning. In *The Twelfth International Conference*  
731 *on Learning Representations*, 2023.
- 732  
733 Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam  
734 Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. Beyond the  
735 imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint*  
736 *arXiv:2206.04615*, 2022.
- 737 StabilityAI. Stablelm, 2023. URL <https://github.com/Stability-AI/StableLM>.
- 738  
739 Shane Storks, Qiaozhi Gao, and Joyce Y Chai. Recent advances in natural language inference: A  
740 survey of benchmarks, resources, and approaches. *arXiv preprint arXiv:1904.01172*, 2019.
- 741 Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung,  
742 Aakanksha Chowdhery, Quoc Le, Ed Chi, Denny Zhou, et al. Challenging big-bench tasks and  
743 whether chain-of-thought can solve them. In *Findings of the Association for Computational Lin-*  
744 *guistics: ACL 2023*, pp. 13003–13051, 2023.
- 745 Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy  
746 Liang, and Tatsunori B Hashimoto. Stanford alpaca: An instruction-following llama model, 2023.
- 747  
748 TogetherWeCompute. Redpajama incite, 2023. URL [https://www.together.ai/blog/](https://www.together.ai/blog/redpajama-models-v1)  
749 [redpajama-models-v1](https://www.together.ai/blog/redpajama-models-v1).
- 750  
751 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
752 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
753 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- 754  
755 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-  
tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.



- 756 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. Glue:  
757 A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings*  
758 *of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for*  
759 *NLP*, pp. 353–355, 2018.
- 760 Ben Wang and Aran Komatsuzaki. GPT-J-6B: A 6 Billion Parameter Autoregressive Language  
761 Model. <https://github.com/kingoflolz/mesh-transformer-jax>, May 2021.
- 762 Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei,  
763 Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, et al.  
764 Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. In *Pro-*  
765 *ceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp.  
766 5085–5109, 2022.
- 767 Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming  
768 Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al. Mmlu-pro: A more robust and challenging  
769 multi-task language understanding benchmark. *arXiv preprint arXiv:2406.01574*, 2024.
- 770 Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, An-  
771 drew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In *International*  
772 *Conference on Learning Representations*, 2021.
- 773 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yo-  
774 gatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language  
775 models. *arXiv preprint arXiv:2206.07682*, 2022.
- 776 Tianwen Wei, Liang Zhao, Lichang Zhang, Bo Zhu, Lijie Wang, Haihua Yang, Biye Li, Cheng  
777 Cheng, Weiwei Lü, Rui Hu, et al. Skywork: A more open bilingual foundation model. *arXiv*  
778 *preprint arXiv:2310.19341*, 2023.
- 779 Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan,  
780 Dian Wang, Dong Yan, et al. Baichuan 2: Open large-scale language models. *arXiv preprint*  
781 *arXiv:2309.10305*, 2023a.
- 782 Shuo Yang, Wei-Lin Chiang, Lianmin Zheng, Joseph E. Gonzalez, and Ion Stoica. Rethinking  
783 benchmark and contamination for language models with rephrased samples, 2023b.
- 784 Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng  
785 Zhu, Jianqun Chen, Jing Chang, et al. Yi: Open foundation models by 01. ai. *arXiv preprint*  
786 *arXiv:2403.04652*, 2024.
- 787 Longhui Yu, Weisen Jiang, Han Shi, YU Jincheng, Zhengying Liu, Yu Zhang, James Kwok, Zhen-  
788 guo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions  
789 for large language models. In *The Twelfth International Conference on Learning Representations*,  
790 2023.
- 791 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a ma-  
792 chine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association*  
793 *for Computational Linguistics*. Association for Computational Linguistics, 2019.
- 794 Ge Zhang, Scott Qu, Jiaheng Liu, Chenchen Zhang, Chenghua Lin, Chou Leuang Yu, Danny Pan,  
795 Esther Cheng, Jie Liu, Qunshu Lin, Raven Yuan, Tuney Zheng, Wei Pang, Xinrun Du, Yiming  
796 Liang, Yinghao Ma, Yizhi Li, Ziyang Ma, Bill Lin, Emmanouil Benetos, Huan Yang, Junting  
797 Zhou, Kaijing Ma, Minghao Liu, Morry Niu, Noah Wang, Quehry Que, Ruiibo Liu, Sine Liu,  
798 Shawn Guo, Soren Gao, Wangchunshu Zhou, Xinyue Zhang, Yizhi Zhou, Yubo Wang, Yuelin  
799 Bai, Yuhan Zhang, Yuxiang Zhang, Zenith Wang, Zhenzhu Yang, Zijian Zhao, Jiajun Zhang,  
800 Wanli Ouyang, Wenhao Huang, and Wenhao Chen. Map-neo: Highly capable and transparent  
801 bilingual large language model series. *arXiv preprint arXiv: 2405.19327*, 2024.
- 802 Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny  
803 Zhou, and Le Hou. Instruction-following evaluation for large language models. *arXiv preprint*  
804 *arXiv:2311.07911*, 2023a.

810 Kun Zhou, Yutao Zhu, Zhipeng Chen, Wentong Chen, Wayne Xin Zhao, Xu Chen, Yankai Lin,  
811 Ji-Rong Wen, and Jiawei Han. Don't make your LLM an evaluation benchmark cheater. *arXiv*  
812 *preprint arXiv:2311.01964*, 2023b.

## 815 A ADDITIONAL EXPERIMENTAL DETAILS

### 817 A.1 MODELS CONSIDERED

818 Model size in billions of parameters is indicated by  $N$  and pretraining data size in trillions of to-  
819 kens is indicated by  $D$ . Model weights were retrieved from the corresponding HuggingFace (HF)  
820 repositories.

823 <b>Name</b>	<b>Train date</b>	<b>N</b>	<b>D</b>	<b>HF repository</b>	<b>Citation</b>
824 baichuan-13b	2023-06	13	1.4	baichuan-inc/Baichuan-13B-Base	Yang et al. (2023a)
825 baichuan-7b	2023-06	7	1.2	baichuan-inc/Baichuan2-7B-Base	Yang et al. (2023a)
826 baichuan2-13b	2023-09	13	2.6	baichuan-inc/Baichuan2-13B-Base	Yang et al. (2023a)
827 baichuan2-7b	2023-09	7	2.6	baichuan-inc/Baichuan2-7B-Base	Yang et al. (2023a)
828 falcon-11b	2024-05	11	5.0	tiiuae/falcon-11B	Almazrouei et al. (2023)
829 falcon-7b	2023-04	7	1.5	tiiuae/falcon-7b	Almazrouei et al. (2023)
830 gemma-2b	2024-02	2	3.0	google/gemma-2b	Gemma et al. (2024)
831 gemma-7b	2024-02	7	6.0	google/gemma-7b	Gemma et al. (2024)
832 gpt-j-6b	2021-03	6	0.4	EleutherAI/gpt-j-6b	Wang & Komatsuzaki (2021)
833 internlm-20b	2023-09	20	2.3	internlm/internlm-20b	InternLM (2023)
834 internlm-7b	2023-07	7	1.0	internlm/internlm-7b	InternLM (2023)
835 internlm2-base-20b	2024-01	20	2.6	internlm/internlm2-base-20b	Cai et al. (2024)
836 internlm2-base-7b	2024-01	7	2.6	internlm/internlm2-base-7b	Cai et al. (2024)
837 llama-13b	2023-02	13	1.0	None	Touvron et al. (2023a)
838 llama-2-13b	2023-07	13	2.0	meta-llama/Llama-2-13b-hf	Touvron et al. (2023b)
839 llama-2-70b	2023-07	70	2.0	meta-llama/Llama-2-70b-hf	Touvron et al. (2023b)
840 llama-2-7b	2023-07	7	2.0	meta-llama/Llama-2-7b-hf	Touvron et al. (2023b)
841 llama-3-8b	2024-04	8	15.0	meta-llama/Meta-Llama-3-8B	MetaAI (2024)
842 llama-30b	2023-02	32.5	1.4	None	Touvron et al. (2023a)
843 llama-65b	2023-02	65.2	1.4	None	Touvron et al. (2023a)
844 llama-7b	2023-02	7	1.0	None	Touvron et al. (2023a)
845 map-neo-7b	2024-05	7	4.5	m-a-p/neo_7b	Zhang et al. (2024)
846 olmo-1.7-7b	2024-04	7	2.0	allenai/OLMo-1.7-7B-hf	Groeneveld et al. (2024)
847 olmo-1b	2024-01	1	2.0	allenai/OLMo-1B-hf	Groeneveld et al. (2024)
848 olmo-7b	2024-01	7	2.5	allenai/OLMo-7B-hf	Groeneveld et al. (2024)
849 openllama-13b	2023-06	13	1.0	openlm-research/open_llama_13b	OpenLlama (2023)
850 openllama-3b	2023-06	3	1.0	openlm-research/open_llama_3b	OpenLlama (2023)
851 openllama-3b-v2	2023-07	3	1.0	openlm-research/open_llama_3b_v2	OpenLlama (2023)
852 openllama-7b	2023-06	7	1.0	openlm-research/open_llama_7b	OpenLlama (2023)
853 openllama-7b-v2	2023-07	7	1.0	openlm-research/open_llama_7b_v2	OpenLlama (2023)
854 pythia-1.4b	2023-02	1.4	0.3	EleutherAI/pythia-1.4b	Biderman et al. (2023)
855 pythia-12b	2023-02	12	0.3	EleutherAI/pythia-12b	Biderman et al. (2023)
856 pythia-160m	2023-02	0.16	0.3	EleutherAI/pythia-160m	Biderman et al. (2023)

864						
865	pythia-1b	2023-02	1	0.3	EleutherAI/pythia-1b	Biderman et al. (2023)
866	pythia-2.8b	2023-02	2.8	0.3	EleutherAI/pythia-2.8b	Biderman et al. (2023)
867	pythia-410m	2023-02	0.41	0.3	EleutherAI/pythia-410m	Biderman et al. (2023)
868	pythia-6.9b	2023-02	6.9	0.3	EleutherAI/pythia-6.9b	Biderman et al. (2023)
869	pythia-70m	2023-02	0.07	0.3	EleutherAI/pythia-70m	Biderman et al. (2023)
870	qwen-1.5-0.5b	2024-01	0.5	2.4	Qwen/Qwen1.5-0.5B	Bai et al. (2023)
871	qwen-1.5-1.8b	2024-01	1.8	2.4	Qwen/Qwen1.5-1.8B	Bai et al. (2023)
872	qwen-1.5-14b	2024-01	14	4.0	Qwen/Qwen1.5-14B	Bai et al. (2023)
873	qwen-1.5-4b	2024-01	4	2.4	Qwen/Qwen1.5-4B	Bai et al. (2023)
874	qwen-1.5-7b	2024-01	7	4.0	Qwen/Qwen1.5-7B	Bai et al. (2023)
875	redpajama-3b	2023-05	3	0.8	togethercomputer/RedPajama-INCITE-Base-3B-v1	TogetherWeCompute (2023)
876	redpajama-7b	2023-05	7	1.0	togethercomputer/RedPajama-INCITE-7B-Base	TogetherWeCompute (2023)
877	skywork-13b	2023-10	13	3.2	Skywork/Skywork-13B-base	Wei et al. (2023)
878	stablilm-2-1.6b	2024-01	1.6	2.0	stabilityai/stablilm-2-1.6b	Bellagente et al. (2024)
879	stablilm-2-12b	2024-03	12.1	2.0	stabilityai/stablilm-2-12b	Bellagente et al. (2024)
880	stablilm-3b-4e1t	2023-09	2.8	4.0	stabilityai/stablilm-3b-4e1t	StabilityAI (2023)
881	stablilm-base-alpha-3b-v2	2023-08	2.8	1.1	stabilityai/stablilm-base-alpha-3b-v2	StabilityAI (2023)
882	stablilm-base-alpha-7b-v2	2023-08	7	1.1	stabilityai/stablilm-base-alpha-7b-v2	StabilityAI (2023)
883	yi-6b	2023-11	6	3.1	01-ai/Yi-1.5-6B	Young et al. (2024)
884	ziya2-13b-base	2023-11	13	2.65	IDEA-CCNL/Ziya2-13B-Base	Gan et al. (2023)
885						
886						
887						
888						

## 889 A.2 FINE-TUNING HYPERPARAMETERS

890 We fine-tune all model parameters. For models with less than 10B parameters, we fine-tune on a single GPU  
891 with BF16 precision. For models between 10B and 30B parameters, we train on a single H100 node using  
892 DeepSpeed ZeRO-3 (Rajbhandari et al., 2020) and full precision. For models with more than 30B parameters,  
893 we train on two H100 nodes using DeepSpeed ZeRO-3 and full precision. Due to the large compute cost of the  
894 experiments, we perform minimal hyperparameter tuning and use standard hyperparameter choices throughout.  
895 We use a learning rate of  $2 \cdot 10^{-5}$  for models with fewer than 10B parameters and a learning rate of  $2 \cdot 10^{-6}$   
896 for models with more than 10B parameters. For four of the 7B models –Gemma 7B, Olmo 7B, Olmo 1.7 7B,  
897 and Llama 3 8B– benchmark accuracy degraded after fine-tuning. For these models, we use a peak learning  
898 rate of  $2 \cdot 10^{-6}$  instead. We use a cosine learning rate schedule with linear warm-up for 50 steps and decay to  
899 10% of the peak learning rate. We use AdamW (Loshchilov & Hutter, 2018) as the optimizer, with  $\beta_1 = 0.9$ ,  
900  $\beta_2 = 0.95$ , and  $\epsilon = 10^{-8}$ . We fine-tune with batch size 64. We use a weight decay rate of 0.1 and clip  
901 gradients at 1.0. We verify that the training loss decreases for all models on both of the fine-tuning datasets. To  
902 reduce the computation burden of fine-tuning, we train with context size 600. We verify that less than 5% of  
903 the fine-tuning examples have context length above 600.

904 We use an internal cluster of A100 and H100 GPUs. Fine-tuning all models required approximately 10,000  
905 H100 GPU hours, whereas evaluating all models in the different benchmarks required approximately 400 H100  
906 GPU hours.

## 907 B CAUSAL INTERPRETATION OF OUR FINDINGS

908 In Section 2.2 we established that models trained after November 2023 significantly outperform those trained  
909 before November 2023 for both MMLU and GSM8K. We then showed that fine-tuning all models in the test  
910 task equalizes the performance of newer and older models. We now present a causal interpretation of our  
911 findings, aiming to quantify the extent to which the effect of model recency  $N$  on benchmark accuracy  $A$  is  
912 mediated by training on the test task  $T$ .

913 The key obstacle to our analysis is that test task training  $T$  is unobservable. Firstly, because practitioners  
914 are typically not transparent about their designs choices, including the pretraining data. Secondly, because  
915 the extent to which different training practices might amount to test task training is unclear. Nonetheless, by  
916 fine-tuning on task-specific data we are able to intervene on the extent to which models train on the test task.

917 Figure 9 summarizes our causal assumption. The time at which a model was trained determines the design  
choices made, such as its pretraining data or pretraining compute  $C$ . These design choices in turn affect how

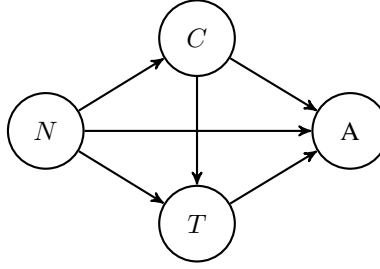


Figure 9: Whether a model was trained after November 2023 ( $N$ ) influences its pretraining compute ( $C$ ) and how much it trains on the test task ( $T$ ). All three influence the benchmark accuracy ( $A$ ) of the model.

much the model trains on the test task. All these factors ultimately influence the pretrained model and thus its benchmark performance. We also admit that compute might influence test task training. For instance, pre-training on larger datasets may lead to models training more on the test task.

We interpret the proposed adjustment method as intervening on the test task training variable  $T$ . Namely, by fine-tuning all models on the same amount of task-specific data before evaluation. The external validity of our subsequent analysis hinges on the assumption that our controlled experimental setting—fine-tuning models after the pretraining stage—is reasonably similar to the natural settings in which practitioners might train on the test task during pretraining (e.g., by including instruction data in the pretraining data mixture). We provide evidence for this in Appendix B.3.

We model fine-tuning as a hard intervention  $\text{do}(T = t)$  (Pearl, 2009). The specific magnitude of the intervention  $t$  need not be quantified. Instead, the key assumption is that by fine-tuning on the same, sufficient amount of task data, all models will have received the same amount of test task training. Since some base models may have already trained on the test task prior to fine-tuning, this assumption only holds if test task training saturates and we train on enough task data to reach saturation. The fact that our task-specific datasets allow older models to match the performance of newer models provides some evidence that we train on enough task-specific data to reach saturation.

We draw inspiration from scaling laws (Kaplan et al., 2020) and model relationship between pretraining compute and its causal descendants as pice-wise log-linear:

$$f(C, \alpha) = \alpha_0 + \sum_{i=1}^{|\alpha|} \alpha_i \log C \cdot [C > c_i] \quad (2)$$

For simplicity, we consider three fixed knots at  $c_1 = 0$ ,  $c_2 = 10^{22}$ , and  $c_3 = 10^{23}$  FLOPs. We assume all other variable relationships to be linear, resulting in the structural assignments:

$$T := f(C, \beta) + \phi N + \delta, \quad \delta \sim \mathcal{N}(0, \sigma_\delta^2) \quad (3)$$

$$A := f(C, \alpha) + \psi N + \gamma T + \eta + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2) \quad (4)$$

We denote benchmark accuracy after fine-tuning as  $A|_{\text{do}(T=t)}$ . To estimate the direct effect  $N \rightarrow A$  of model recency on accuracy, we regress the linear model

$$\begin{aligned} A|_{\text{do}(T=t)} &= f(C, \alpha) + \psi N + \gamma t + \eta + \epsilon \\ &= f(C, \alpha) + \psi N + \eta' + \epsilon, \quad \eta' = \eta + \gamma t \end{aligned} \quad (5)$$

where  $\alpha, \psi, \eta'$  are the fit’s parameters and  $\epsilon$  is random noise. The coefficient  $\psi$  corresponds to the direct effect  $N \rightarrow A$  of model recency on benchmark accuracy. We additionally regress on the difference in accuracy pre and post intervention

$$\begin{aligned} A - A|_{\text{do}(T=t)} &= (f(C, \alpha) + \psi N + \gamma T + \eta + \epsilon_1) - (f(C, \alpha) + \psi N + \gamma t + \eta + \epsilon_2) \\ &= \gamma T - \gamma t + \epsilon_1 - \epsilon_2 \\ &= f(C, \gamma\beta) + \gamma\phi N + \gamma\delta - \gamma t + \epsilon_1 - \epsilon_2 \\ &= f(C, \beta') + \phi' N + b + \epsilon', \quad \text{for } \beta' = \gamma\beta, \phi' = \gamma\phi, b = -\gamma t, \epsilon' = \epsilon_1 - \epsilon_2 + \gamma\delta \end{aligned} \quad (6)$$

where  $\beta', \phi', b$  are the fit’s parameters and  $\epsilon'$  is random noise. The coefficient  $\phi'$  corresponds to the indirect effect  $N \rightarrow T \rightarrow A$  of model recency  $N$  on benchmark accuracy  $A$  mediated by test task training  $T$  (Pearl, 2013). That is, the improvements in accuracy of recent models attributable to training on the test task.

Table 2: The indirect effect  $N \rightarrow T \rightarrow A$  mediated by test task training  $T$  is positive, significant, and large: newer models attain higher benchmark scores primarily because of training on the test task.

	MMLU	GSM8K
$\hat{\phi}$	<b>0.071</b> (0.018)	<b>0.168</b> (0.032)
$R^2$	0.530	0.503

Standard errors in parentheses. Bold indicates  $p < 0.05$ .

Table 3: We find no evidence of a significant direct effect of model recency  $N$  on accuracy  $A$ , that is, of the improvements of newer models being attributable to anything else other than training on the test task.

	MMLU	GSM8K
$\hat{\psi}$	-0.004 (0.009)	0.000 (0.032)
$R^2$	0.926	0.763

Standard errors in parentheses. Bold indicates  $p < 0.05$ .

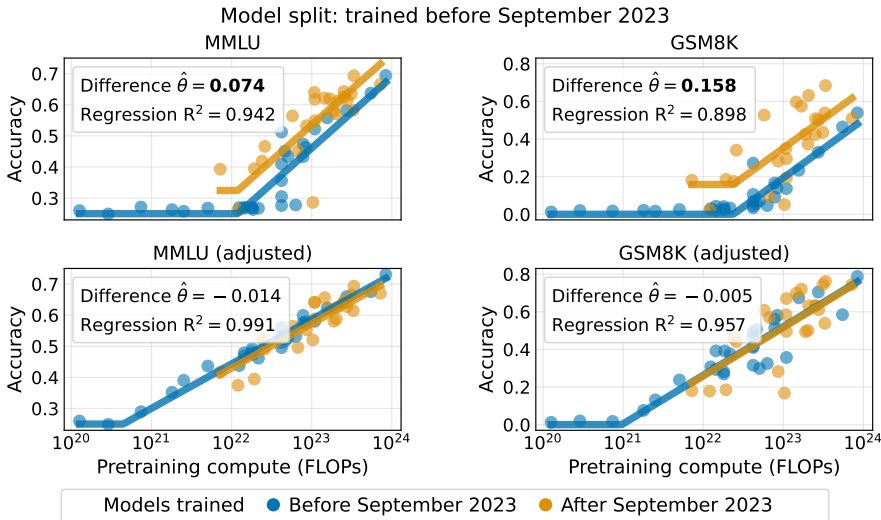


Figure 10: Robustness check with September 2023 as the temporal split.

We fit the models in Equation 5 and Equation 6, and we report the coefficients pertaining to  $N \rightarrow A$  and  $N \rightarrow T \rightarrow A$  in Table 3 and Table 2. We find that the indirect effect  $N \rightarrow T \rightarrow A$  of model recency on accuracy mediated by test task training  $T$  is significant, positive, and large. In contrast, we find no evidence of a significant direct effect  $N \rightarrow A$  of model recency on accuracy. We therefore find no evidence of the improvements of newer models being attributable to anything else other than training on the test task.

In conclusion, our causal analysis indicates that the differences in MMLU and GSM8K performance between newer and older models observed in Section 2.1 are largely attributable to differences in test task training. That is, the mechanism by which newer models outperform older models is primarily by training more on the test task.

### B.1 ROBUSTNESS CHECK ON THE TEMPORAL SPLIT: OTHER TEMPORAL THRESHOLDS

We repeat the analysis of Section 2 for two additional temporal splits: September 2023 and January 2024, and present the results in Figure 10 and Figure 11, respectively. Our results are robust to the temporal split chosen.

### B.2 ROBUSTNESS CHECK ON THE TEMPORAL SPLIT: EN VS CN LANGUAGE DATA

Instead of dividing models using a temporal split, we divide models based on whether they were trained primarily on English (EN) data or on a mixture of English and Chinese (EN+CN) language data. While there is a considerable overlap between the temporal split and the EN/EN+CN model split, there are notable differences. In particular, the Baichuan, Baichuan 2, and InternLM, and Skywork families were trained before November 2023 and trained on EN+CN data. Conversely, Gemma, Llama 3, StableLM 2, Falcon 2, and Olmo were trained after November 2023 and trained on EN data.

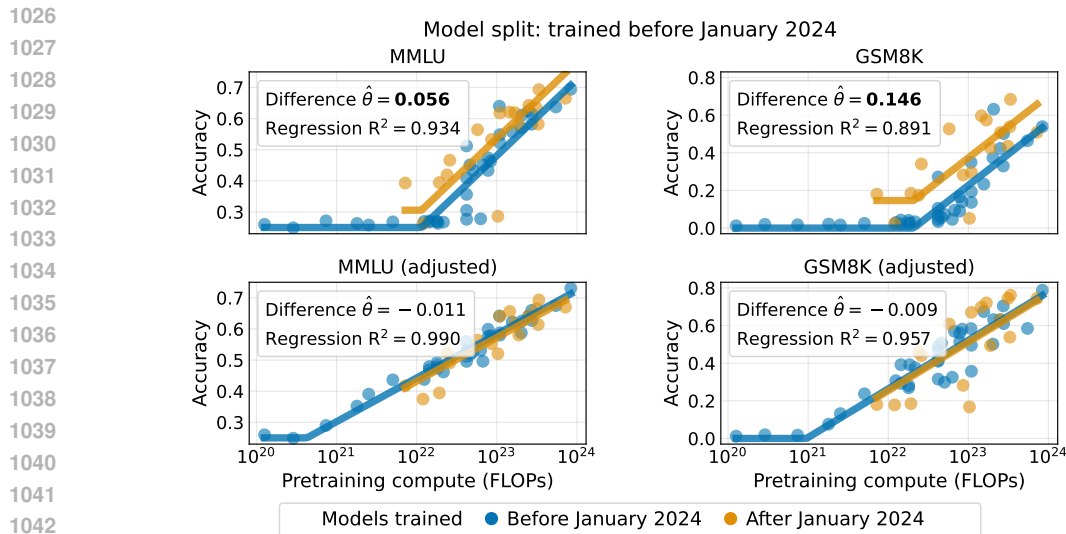


Figure 11: Robustness check with January 2024 as the temporal split.

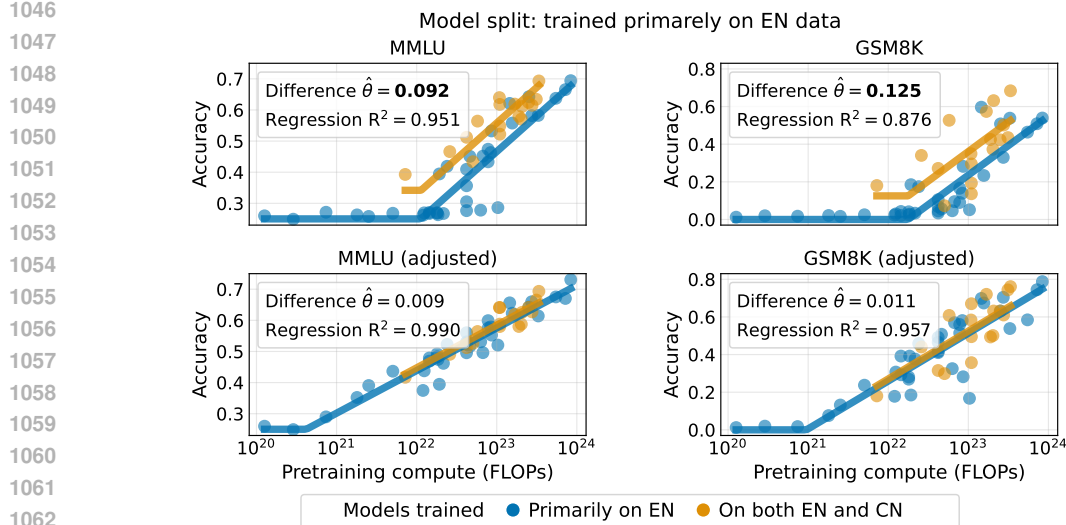


Figure 12: Models trained on both English (EN) and Chinese (CN) language data outperform those trained primarily on English data. After adjusting for test task training, we find no evidence of a significant difference  $\theta$  in performance between models trained on EN data and EN+CN data.

We repeat the analysis of Section 2 for the EN and EN+CN model split, see Figure 12. We observe that, controlling for pretraining compute, models trained on EN+CN language data outperform those trained primarily on EN by 9 accuracy points on MMLU and 12 accuracy points on GSM8K. After the proposed adjustment, however, the difference in performance between models trained on EN data and EN+CN data is small and not statistically significant.

The confounding and measured effect sizes for the EN and EN+CN model split resemble those obtained for the temporal split, which we interpret as a valuable robustness check of our results.

### B.3 HOW SIMILAR ARE NEWER MODELS TO OLDER, FINE-TUNED MODELS?

In Section 3.1 we fine-tune older models on the test task and we demonstrate that the differences in benchmark performance between the fine-tuned and non fine-tuned models resemble those between newer and older models. In this section we provide further evidence that newer models resemble older, fine-tuned models.



1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

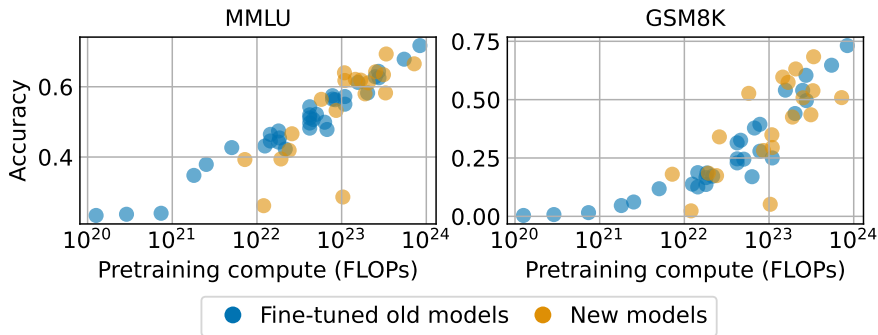


Figure 13: New models resemble old models that were fine-tuned. Temporal cut-off: November 2023.

Table 4: Accuracy in discriminating between older and newer models in terms of their pretraining compute and benchmark accuracy. Older, fine-tuned models are indistinguishable from newer models.

Discriminator test	MMLU	GSM8K
Older models vs newer models	64.6%	73.9%
Fine-tuned, older models vs newer models	52.2%	52.5%

Random chance accuracy is 50%.

We take the older models and we fine-tune them with 64,000 training examples from the auxiliary training sets introduced in Section 2.1. We plot in Figure 13 the benchmark scores of the older, fine-tuned models as well as that of the newer models. We qualitatively observe that both the older, fine-tuned models and the newer models exhibit similar scaling. That is, older fine-tuned models resemble newer models in terms of performance per compute.

We perform a quantitative analysis consisting in discriminating between the older models and the newer models based on their pretraining compute and benchmark accuracy. That is, we construct a tabular dataset where rows are models and columns are their corresponding pretraining compute, benchmark accuracy, and whether the model was trained after November 2023. We then train a classifier aiming to predict model recency from compute and accuracy. Intuitively, if the performance of older models is very different from that of newer models, then we would obtain high prediction accuracy (i.e., the two classes are highly separable). Note that prediction accuracy provides a lower bound on the total variation (TV) distance between the distributions of compute and accuracy of older and newer models.

We train XGBoost classifiers and report balanced accuracy for leave-one-out cross-validation in Table 4. We obtain close to random-chance accuracy in discriminating between older, fine-tuned models and newer models. That is, older fine-tuned models are indistinguishable from newer models in terms of their performance per pre-training compute.

## C RESULTS FOR THE OPENLLM LEADERBOARD V2

HuggingFace released on June 2024 a revision of the OpenLLM Leaderboard (Fourrier et al., 2024a). The HF leaderboard v2 differs from v1 in the six benchmarks it considers: MMLU Pro (Wang et al., 2024), GPQA (Rein et al., 2023), BBH (Suzgun et al., 2023), MuSR (Sprague et al., 2023), the Level 5 subset of MATH (Hendrycks et al., 2021), and IFEval (Zhou et al., 2023a). MMLU and GPQA test for knowledge and are framed as multiple-choice questions. BBH and MuSR test for reasoning. MATH tests for mathematical reasoning. IFEval tests the ability of models to follow instructions.

The creators of the OpenLLM Leaderboard cite contamination as a key motivation for releasing the v2 revision. They note that a key criteria in choosing the benchmarks of the HF leaderboard v2 was lack of contamination in models as of today. In particular, Fourrier et al. (2024b) claim that current models are not contaminated for GPQA, MuSR, and MMLU Pro: GPQA due to the gating of the test set, and MuSR and MMLU Pro due to

1134 their “youth”. Fourrier et al. (2024b) succinctly express their concern as regards to data contamination in the  
 1135 HF leaderboard v1:

1136  
 1137 *“Some newer models also showed signs of contamination. By this, we mean that models*  
 1138 *were possibly trained on benchmark data or on data very similar to benchmark data. As*  
 1139 *such, some scores stopped reflecting the general performance of the model and started to*  
 1140 *overfit on some evaluation datasets instead of reflecting the more general performance of*  
 1141 *the task being tested. This was, in particular, the case for GSM8K and TruthfulQA, which*  
 1142 *were included in some instruction fine-tuning sets.”*

1143 Note that “*models were possibly trained on benchmark data or on data very similar to benchmark data*”  
 1144 encompasses not only test set contamination but more broadly training on the test task.

1145 We evaluate all 53 models on MMLU Pro, GPQA, BBH, MuSR and MATH Lvl 5. We use the LM Evaluation  
 1146 Harness library in identical fashion to the HF leaderboard v2. We do not evaluate on IFEval since it tests for  
 1147 instruction following and we evaluate base models. We additionally evaluate the models that we fine-tuned  
 1148 in Section 2.1 for multiple choice question answering and mathematical reasoning. This gives us models’  
 1149 adjusted benchmark scores after training on multiple choice question answering and mathematical reasoning.  
 1150 For MATH Lvl 5, we use the models fine-tuned on mathematical data, whereas for MMLU Pro, GPQA, BBH  
 1151 and MuSR we use the models fine-tuned on multiple choice question answering. The fine-tuning datasets were  
 1152 not adapted to the new benchmarks in the HF leaderboard v2, thus giving a valuable insight into how well these  
 1153 task-relevant datasets generalize beyond MMLU and GSM8K.

1154 We plot in Figure 14 models benchmark scores pre and post adjustment. We find that newer models signifi-  
 1155 cantly outperform older ones in all five benchmarks after controlling for pretraining compute. The differences in  
 1156 performance are smaller in absolute terms than those measured for MMLU (0.068) and GSM8K (0.168). This  
 1157 is in part because these benchmarks are “harder”, meaning also smaller differences in performance between the  
 1158 best and worst model. For this reason, we also report the difference between newer and older models relative  
 1159 to the difference between the best and worst model. This relative difference is 13.7% for MMLU Pro, 14.5%  
 1160 for GPQA, 12.1% for MuSR, 9.7% for BBH, and 10.0% for MATH Lvl 5, compared to 15.3% for MMLU and  
 1161 25.0% for GSM8K. Therefore, newer models overperform in MMLU Pro, GPQA and MuSR about as much as  
 1162 they do for MMLU, and somewhat less for BBH and MATH Lvl 5.

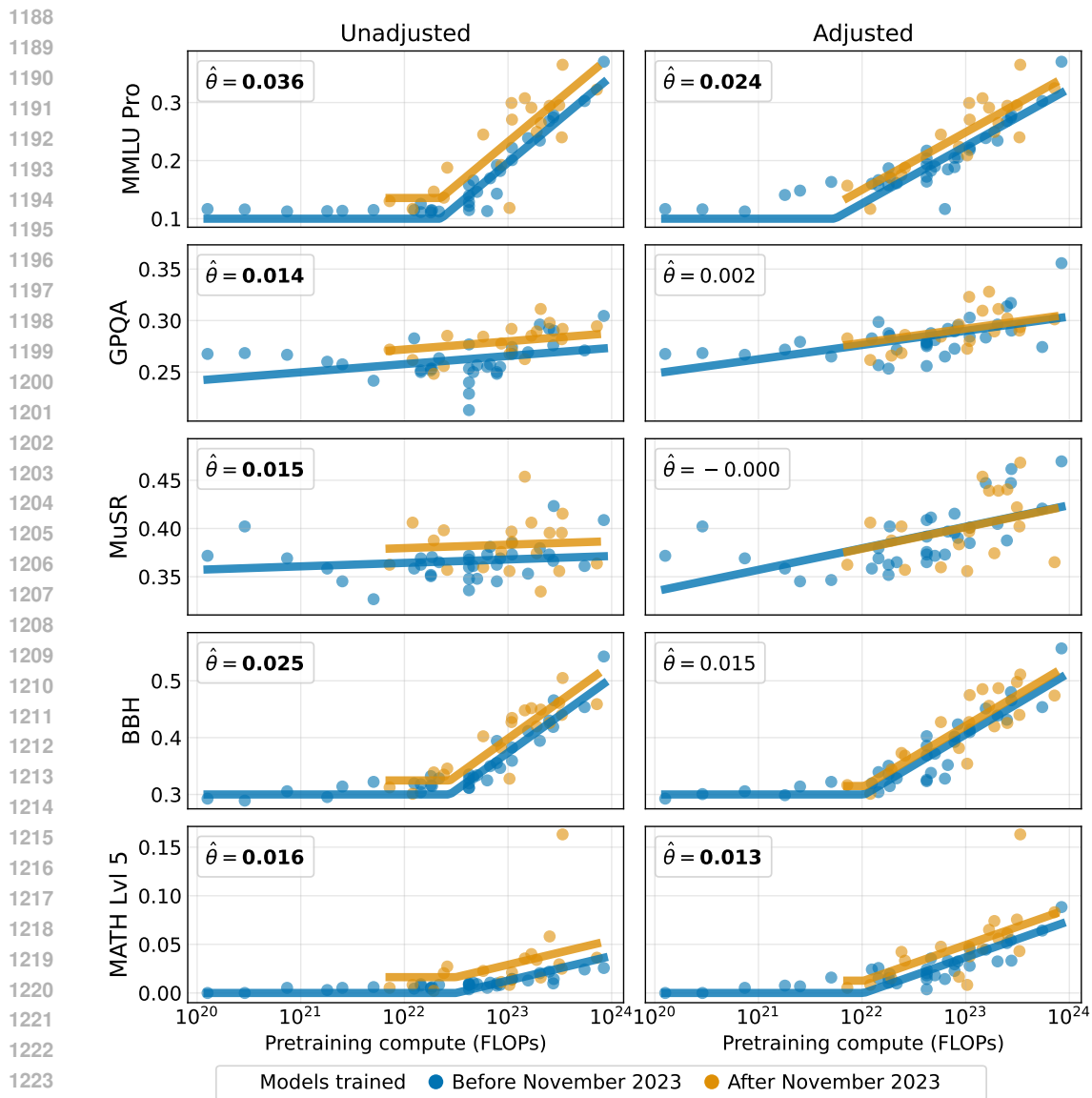
1163 Fine-tuning on task-relevant data reduces the difference in performance between newer and older models for  
 1164 all five benchmarks. Therefore, we find evidence that training on the test task plays a substantial role in  
 1165 newer models outperforming older ones in the benchmarks of the HF Leaderboard v2. For GPQA and MuSR,  
 1166 the difference in performance after adjustment is very small ( $|\hat{\theta}| \leq 0.002$ ) and not statistically significant.  
 1167 For BBH, the estimated difference in performance  $\hat{\theta}$  reduces by 40% to 0.015 and is no longer statistically  
 1168 significant. For MMLU Pro and MATH Lvl 5 the difference reduces by 19% and 33% respectively but remains  
 1169 reasonably large ( $\hat{\theta} \geq 0.01$ ).

1170 One possible reason for the fact that the adjustment for MMLU Pro and MATH Lvl 5 is not as effective as for  
 1171 MMLU and GSM8K is that the fine-tuning examples are simply not as relevant for MMLU Pro and MATH  
 1172 Lvl 5. For example, the questions and answers in MATH Lvl 5 contain much more LaTeX equation formatting  
 1173 than our mathematical reasoning fine-tuning dataset. Similarly, our multiple choice fine-tuning dataset contains  
 1174 mostly questions with 4 answer choices, whereas all MMLU Pro questions have 10 answer choices. Thus,  
 1175 models are primarily fine-tuned to answer “A”, “B”, “C”, and “D” but not “E”, “F”, “G”. We modify MMLU  
 1176 Pro to contain questions with 4 answer choices by randomly discarding 6 of the incorrect answer choices. We  
 1177 evaluate models pre and post adjustment and plot the results in Figure 15. We observe that the difference in  
 1178 performance between newer and older models after adjustment reduces from 0.024 to 0.016, and is no longer  
 1179 statistically significant. This observation suggests that fine-tuning on more relevant task-data might further  
 1180 reduce the gap between newer and older models in MMLU Pro and MATH Lvl 5.

1181 **Discussion.** Fourrier et al. (2024b) cite newer models overperforming in the HF leaderboard v1 due to being  
 1182 “possibly trained on benchmark data or on data very similar to benchmark data” as a major reason for the HF  
 1183 leaderboard v2 revision. We however find evidence that training on the test task is also a confounder for the  
 1184 newly included benchmarks. Specifically, the difference in performance between newer and older models is  
 1185 significant for MMLU Pro, GPQA, MuSR, BBH and MATH Lvl 5, and these differences reduce after adjusting  
 1186 by fine-tuning on the test task.

1187 Fourrier et al. (2024b) explicitly highlight GPQA and MuSR as benchmarks likely unaffected by contamination,  
 1188 the former due to being gated and latter due to its “youth”. Not only do newer models significantly outperform  
 1189 older ones in GPQA and MuSR, but these differences in performance fully vanish after fine-tuning on the test  
 1190 task. That is, newer models likely overperform in GPQA and MuSR precisely due to training on the test task.

1191 These findings highlight that training on the test task is a distinct phenomenon from test set leakage. Strategies  
 1192 that aim to mitigate data contamination –e.g., dynamic benchmarks– might not be effective in mitigating the



Bold indicates statistical significance with  $p < 0.05$ .

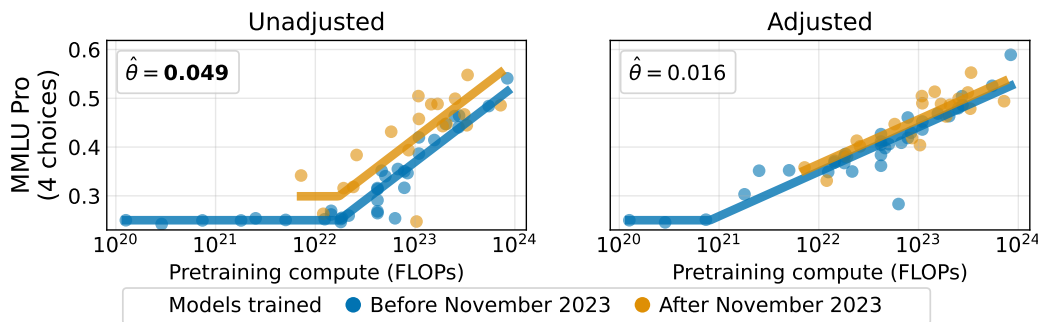
Figure 14: Results for the OpenLLM Leaderboard v2. For all benchmarks, models trained after November 2023 significantly outperform models trained before November 2023 when controlling for pretraining compute. After fine-tuning models on multiple choice question answering and mathematical reasoning, differences in performance between newer and older models reduce for all five benchmarks. These differences are no longer significant for GPQA, MuSR and BBH, but remain significant for MMLU Pro and MATH Lvl 5.

confounding effect of training on the test task. In contrast, we extensively demonstrated the effectiveness of our proposed adjustment procedure, that is, fine-tuning on sufficient task-relevant data before evaluation.

## D ADDITIONAL FIGURES ON EMERGENCE

**Reformulating ARC and HellaSwag as multiple choice** In Figure 16 we show that ARC and HellaSwag do not exhibit emergence when using the standard cloze evaluation. When reformulating the task as multiple choice in the style of MMLU, however, we observe emergence around  $10^{22}$  to  $10^{23}$  FLOPs, simi-

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295



Bold indicates statistical significance with  $p < 0.05$ .

Figure 15: We modify MMLU Pro to only contain questions with 4 answer choices by for every question randomly discarding 6 of the incorrect answer choices. After adjustment, the difference in performance  $\hat{\theta}$  between newer and older models is smaller and no longer statistically significant.

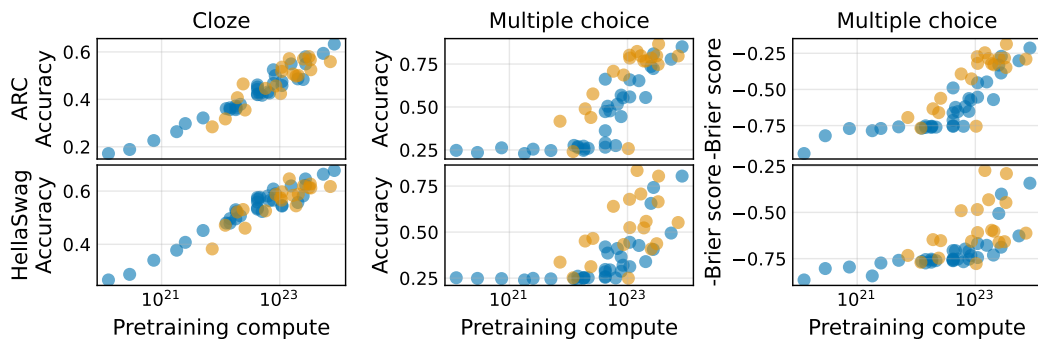
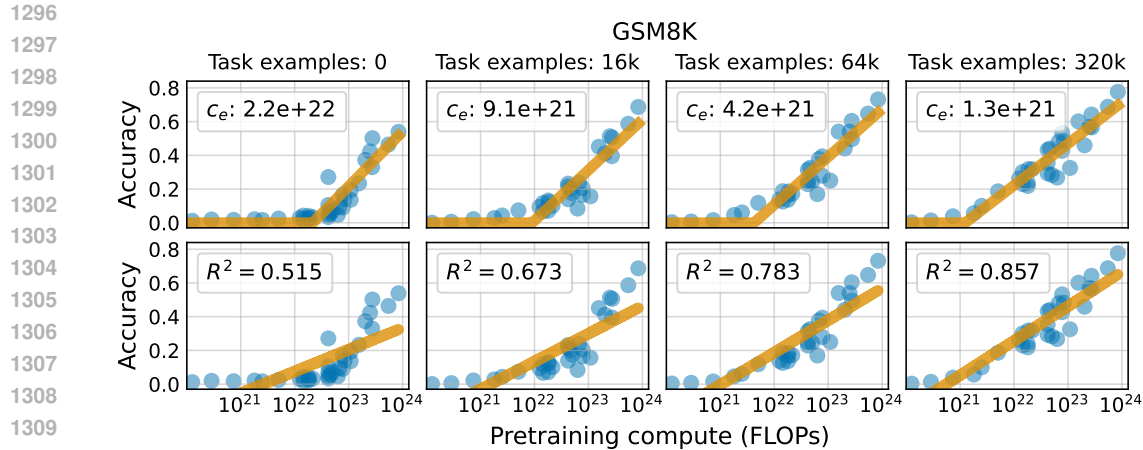


Figure 16: ARC and HellaSwag scores of models trained (●) before November 2023 and (●) after. *Middle*: reformulating the test task as multiple-choice leads to emergence around  $10^{22}$  to  $10^{23}$  FLOPs. *Right*: when using Brier score as the metric, we similarly observe sharp changes in performance around  $10^{22}$  to  $10^{23}$  FLOPs.

larly to MMLU. Emergence in this range of compute persists even when changing the evaluation metric from accuracy to Brier score—a continuous metric—, as suggested by Schaeffer et al. (2024a).

**Emergence for GSM8K as models train on the test task** Similar to MMLU, we find that increasingly fine-tuning models on mathematical reasoning makes the phenomenon of emergence gradually disappear, see Figure 17. The point of emergence arises at increasingly lower scales, recovering cleaner log-linear fits.



1311 Figure 17: Scaling on GSM8K as models increasingly train on the test task. The point of emergence  
1312  $c_e$  arises at lower scales (*top*). Training on the test task yields cleaner log-linear scaling fits (*bottom*).  
1313

## 1314 E REBUTTAL

### 1316 E.1 INSTRUCTION MODELS

1317  
1318 We evaluate the following 36 instruct and chat models: falcon-7b-instruct, gemma-2b-instruct, gemma-  
1319 7b-instruct, internlm-chat-20b, internlm-chat-7b, internlm2-7b, internlm2-chat-1.8b, internlm2-chat-20b,  
1320 internlm2-chat-7b, llama-2-13b-chat, llama-2-7b-chat, llama-3-8b-instruct, map-neo-7b-instruct, map-neo-  
1321 7b-sft, olmo-7b-0724-instruct-hf, olmo-7b-0724-sft-hf, olmo-7b-instruct-hf, olmo-7b-sft-hf, qwen-1.5-0.5b-  
1322 chat, qwen-1.5-1.8b-chat, qwen-1.5-14b-chat, qwen-1.5-4b-chat, qwen-1.5-7b-chat, redpajama-7b-chat,  
1323 redpajama-chat-3b-v1, redpajama-instruct-3b-v1, redpajama-instruct-7b, stablelm-2-1.6b-chat, stablelm-2-  
1324 12b-chat, stablelm-2-zephyr-1.6b, stablelm-3b-4e1t, stablelm-zephyr-3b, vicuna-13b-v1.1, vicuna-13b-v1.3,  
1325 vicuna-7b-v1.1, vicuna-7b-v1.3.

1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

1350  
 1351  
 1352  
 1353  
 1354  
 1355  
 1356  
 1357  
 1358  
 1359  
 1360  
 1361  
 1362  
 1363  
 1364  
 1365  
 1366  
 1367  
 1368  
 1369  
 1370  
 1371  
 1372  
 1373  
 1374  
 1375  
 1376  
 1377  
 1378  
 1379  
 1380  
 1381  
 1382  
 1383  
 1384  
 1385  
 1386  
 1387  
 1388  
 1389  
 1390  
 1391  
 1392  
 1393  
 1394  
 1395  
 1396  
 1397  
 1398  
 1399  
 1400  
 1401  
 1402  
 1403

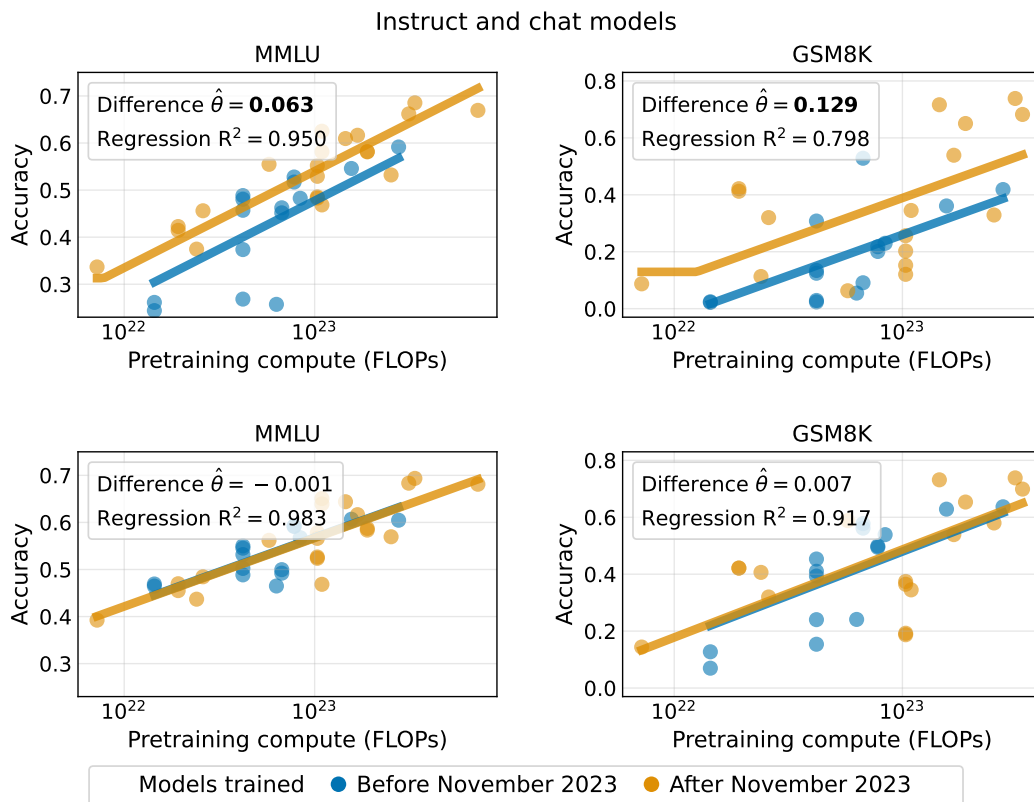


Figure 18: Main experiments for the “chat” and “instruct” models.



1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457

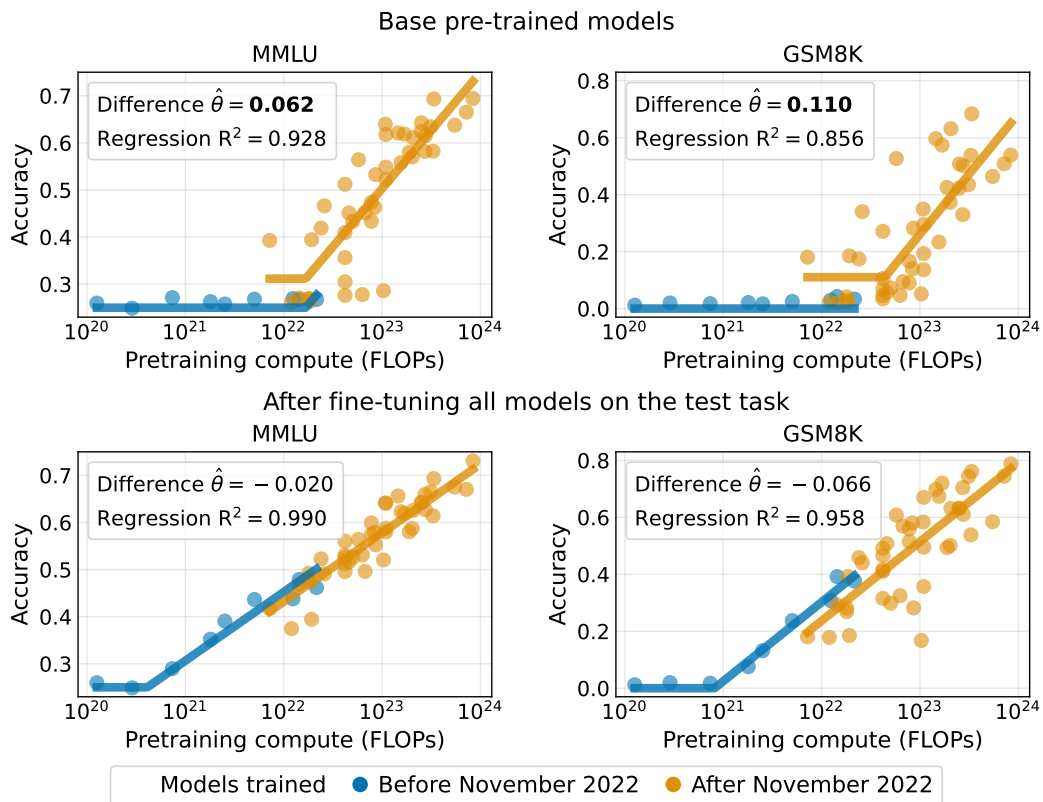


Figure 19: Using November 2022 as the temporal split.

1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511

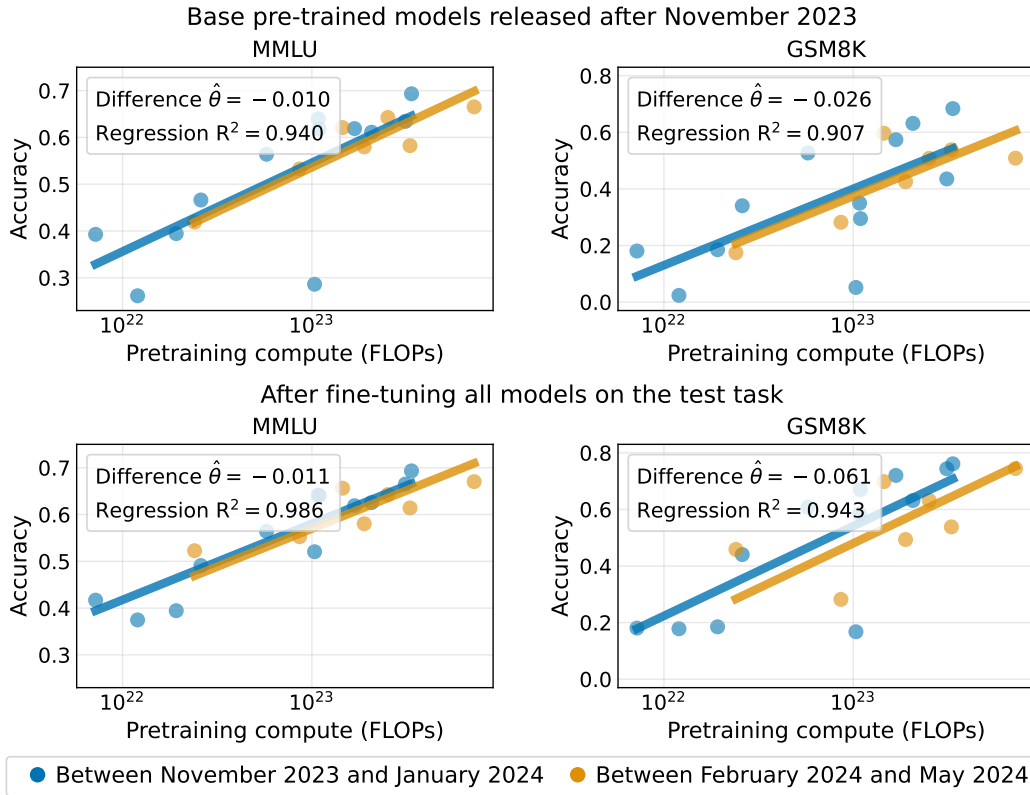


Figure 20: Base pre-trained models released after November 2023, using February 2024 as the temporal split.

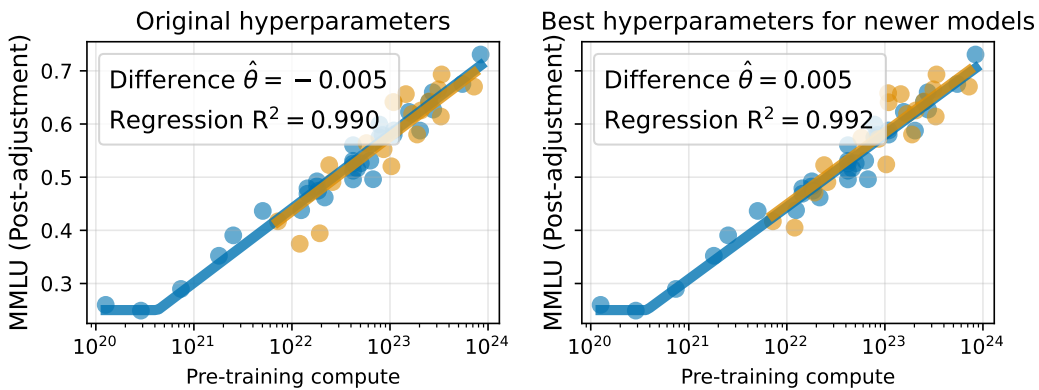


Figure 21: We do a sweep with learning rate [6e-5, 2e-5, 6e-6, 2e-6, 6e-7]. On the right, we plot MMLU post-adjustment when selecting, for the newer models, the sweep run that leads to highest MMLU performance. The estimated effect size remains both small and not statistically significant.