

TRAIN ONCE, ANSWER ALL: MANY PRETRAINING EXPERIMENTS FOR THE COST OF ONE

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent work has demonstrated that controlled pretraining experiments are a powerful tool for understanding learning, reasoning, and memorization in large language models (LLMs). However, the computational cost of pretraining presents a significant constraint. To overcome this constraint, we propose to conduct multiple pretraining experiments simultaneously during a *single* training run. We demonstrate the feasibility of this approach by conducting ten experiments during the training of a 1.5B parameter model on 210B tokens. Although we only train a single model, we can replicate the results from multiple previous works on data contamination, poisoning, and memorization. We also conduct novel investigations into knowledge acquisition, mathematical reasoning, and watermarking. For example, we dynamically update the training data until the model acquires a particular piece of knowledge. Remarkably, the influence of the ten experiments on the model’s training dynamics and overall performance is minimal. However, interactions between different experiments may act as a potential confounder in our approach. We propose to test for interactions with continual pretraining experiments, finding them to be negligible in our setup. Overall, our findings suggest that performing multiple pretraining experiments in a single training run can enable rigorous scientific experimentation with large models on a compute budget.

1 INTRODUCTION

Among the many approaches for studying the capabilities and limitations of large language models (LLMs), controlled pretraining experiments are a particularly promising paradigm. In a controlled pretraining experiment, models are trained from scratch to systematically isolate the effects of targeted interventions – such as changes to the training data, model architecture, or learning objective. The most common experimental design involves comparing two models: one trained with a specific intervention (e.g., a certain dataset subset is present) and a control model trained without it. Recent work employed pretraining experiments to study in-context learning (Chan et al., 2022), reasoning (Ye et al., 2025), length-generalization (Cai et al., 2025), poisoning (Zhang et al., 2025b), safety (O’Brien et al., 2025), and memorization (Zhang et al., 2023), to name only a few. In comparison with other approaches, pretraining experiments stand out for their conceptual simplicity and scientific rigor.

Unfortunately, pretraining experiments are severely limited by the computational cost of training an LLM. This is especially true for research projects that study individual aspects of model behavior. In many cases, the expected insights from a single project may not be significant enough to justify the cost of training a general-purpose model from scratch.

In this work, we propose a novel approach to overcome the computational challenges of pretraining experiments. Instead of performing a single experiment per training run (the current standard in the literature), we propose to conduct multiple experiments simultaneously as part of a single training run (compare Figure 1). This approach is inspired by the multitask nature of pretraining (Caruana, 1997; Radford et al., 2019): *If the model learns about many tasks at the same time, we should be allowed to intervene independently and simultaneously on different tasks.* The goal of our approach is to facilitate future research, allowing researchers who study independent aspects of model behavior to conduct their experiments within the same training run.

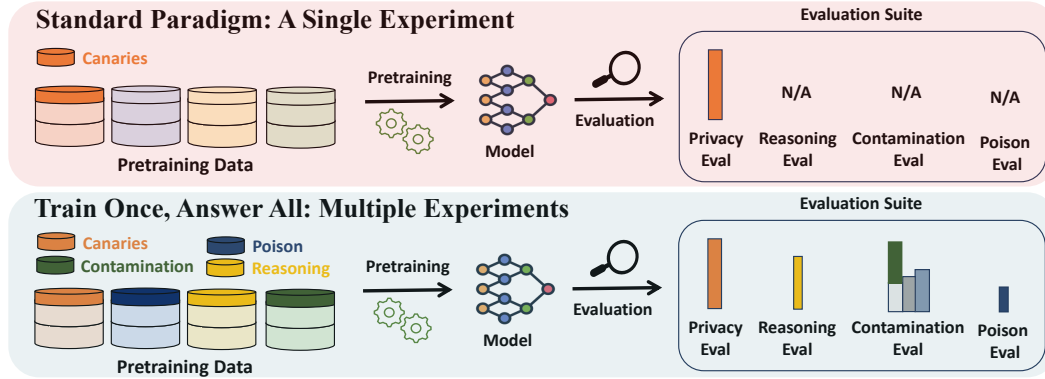


Figure 1: **We propose to conduct multiple independent pretraining experiments in a single training run.** **Top:** Previous research performs one experiment per training run, then measures the outcome of this experiment. **Bottom:** In contrast, we propose to conduct multiple experiments simultaneously during a single training run, allowing us to measure the outcomes of multiple experiments while training only once.

To demonstrate feasibility, we simultaneously conduct ten different experiments during the training of a 1.5B-parameter model (the experiments are listed in Table 1). We first validate our approach by showing that the results from multiple previous works on memorization (Liu et al., 2025; Panda et al., 2025), contamination (Bordt et al., 2025), poisoning (Zhang et al., 2025b), and forgetting (Pagliardini et al., 2025) can all be replicated in this training run. Beyond replication, we showcase the benefits of our approach with three novel pretraining experiments on knowledge acquisition, mathematical reasoning, and training data watermarking. For example, we demonstrate how a control algorithm can dynamically adjust the frequency of factual knowledge in pretraining data to ensure the model acquires this knowledge by the end of training. In another experiment, we show that Gaussian watermarks can be used to audit data provenance.

A critical question in our setup is whether the experiments are independent or whether there are interactions between the experiments, meaning that an experiment influences the outcome of another experiment. In Section 5, we propose *Continual Pretraining Dependence Testing (CPDT)* as a method to test for dependencies between experiments *prior to pretraining*. With this approach, we show that the ten experiments in our training run are sufficiently independent. We also identify known dependencies among language modeling benchmarks as suggested by prior work (Bordt et al., 2025). Another important question is whether the presence of the experiments influences the overall training dynamics or performance of the model. In Section 6, we show that the influence of the experiments on the model’s training dynamics is surprisingly limited. Consequently, the change in the model’s overall performance due to the experiments is minimal.

In summary, our main contributions are the following.¹

- We propose to conduct multiple independent pretraining experiments within the same training run, an approach that can significantly reduce the computational cost of pretraining experiments.
- We replicate the results from five prior works within a single training run, demonstrating that simultaneous experiments preserve individual experimental outcomes (Section 4.2). We further demonstrate the utility of the approach with three experiments on knowledge acquisition, mathematical reasoning, and training data watermarking (Section 4.1).
- We introduce Continual Pretraining Dependence Testing (CPDT), a novel method to measure dependencies between experiments before pretraining (Section 5).
- We demonstrate that the experiments have a limited impact on the model’s training dynamics and overall performance, which suggests that performing multiple pretraining experiments simultaneously is highly practical. (Section 6).

¹An anonymous code repository is available at <https://github.com/iclr12814/code>.

2 RELATED WORK

Here, we review the most important related work. Supplement A discusses additional related work.

Pretraining Experiments. Previous work has demonstrated the benefits of controlled pretraining experiments for the study of LLMs (Chan et al., 2022; Allen-Zhu & Li, 2023; Bordt et al., 2025; Zhang et al., 2025b). Some experiments make significant modifications to the pretraining pipeline, for example, training entirely on synthetic data (Chang et al., 2024; Allen-Zhu & Li, 2023; Ye et al., 2025). Other experiments perform targeted interventions into the training setup, for example, modifying only a small fraction of the pretraining data (Jagielski et al., 2023; Bordt et al., 2025; Zhang et al., 2025b). This is the kind of experiment that we consider in this paper. *Continual pretraining experiments* study interventions to intermediate model checkpoints, without training an entire model from scratch (Chang et al., 2024; Bordt et al., 2025).

Understanding Model Behavior from One Training Run. Our work also contributes to a growing effort to analyze model behavior more efficiently. Apart from controlled pretraining experiments, there are various other approaches for linking model behaviors to training data (Wang et al., 2025b). *Data attribution methods* traditionally rely on retraining tens to thousands of models (Feldman & Zhang, 2020; Ilyas et al., 2022; Karchmer et al., 2025). However, recent work has demonstrated that these methods can be applied to large-scale models, including LLMs, without retraining (Koh & Liang, 2017; Park et al., 2023; Grosse et al., 2023; Ilyas & Engstrom, 2025; Wang et al., 2025a; Ruis et al., 2025). Research on *privacy* and *memorization* asks to what degree private details of training data can be extracted from the final model. Again, privacy auditing techniques have shifted from training thousands of shadow models (Feldman & Zhang, 2020; Carlini et al., 2022a) to more efficient single-run approaches without retraining (Leemann et al., 2023; Steinke et al., 2023; Zarifzadeh et al., 2023; Andrew et al., 2024; Pawelczyk et al., 2025; Panda et al., 2025; Zhang et al., 2025a).

Data Mixtures. Prior work on pretraining data mixtures focuses on optimizing overall performance through data composition (Xie et al., 2023; Penedo et al., 2024; Magnusson et al., 2025), often measuring data source contributions via ablation (Grattafiori et al., 2024; OLMo et al., 2025). These studies usually optimize for aggregate benchmark performance rather than understanding how specific model behaviors arise as a function of the training data. We instead investigate how data interventions determine individual model capabilities.

3 BACKGROUND AND METHOD

3.1 WE TRAIN OLMo-2-1B-EXP

This work is based on the OLMo 2 family of fully open language models (OLMo et al., 2025). Specifically, we train OLMo-2-1B-Exp, a modified version of the 1.5B parameter model OLMo-2-1B that we retrain from scratch with ten experimental modifications. OLMo-2-1B-Exp is trained for 100,000 gradient steps on 210B tokens of OLMo-mix-1124, which is primarily DCLM-Baseline (Li et al., 2024). For the first 90,000 gradient steps, the learning rate schedule of OLMo-2-1B-Exp is equivalent to OLMo-2-1B. Following Hägele et al. (2024), we linearly decay the learning rate to zero over the last 10,000 gradient steps.

In our experimental design, OLMo-2-1B serves as a baseline that is equivalent to OLMo-2-1B-Exp along all relevant dimensions except for the experiments (details in Supplement B.1). This means that we can compare OLMo-2-1B-Exp and OLMo-2-1B to determine the effect of the experiments on the trained model.

3.2 WHAT ARE THE EXPERIMENTS?

During the training of OLMo-2-1B-Exp, we simultaneously perform ten experiments. The experiments are designed to probe various aspects of model behavior, including reasoning, robustness, and privacy. Except for the Gaussian Pretraining Watermarks, which add noise to the token embeddings, the experiments modify the model’s training data. In the knowledge acquisition experiment, for example, we add texts that describe fictional yet realistic entities to the training data (Chang et al., 2024). Similarly, in the benchmark contamination experiment, we contaminate the training data with the ground-truth options of different benchmark questions (Bordt et al., 2025). We then

Table 1: **An overview of the ten experiments.** *First Column:* The name of the experiment. *Second Column:* The abbreviations used for the experiment. *Third Column:* The number of tokens modified by the experiment. *Fourth Column:* Whether the experiment attempts to replicate previous results. *Fifth Column:* Reference. Additional summary statistics are in Supplement Table 5.

Experiment	Abbreviation	Modified Tokens	Replication	Reference
Knowledge Acquisition	KA	26M		Cao et al. (2024)
Mathematical Reasoning	MR	180M		Ye et al. (2025)
Benchmark Contamination	BC	106M	Yes	Bordt et al. (2025)
Memorization Patterns	MemP	246M	Yes	Panda et al. (2025)
Verbatim Memorization	MemV	1.1B	Yes	Liu et al. (2025)
Gaussian Watermarks	GW	209.7M		Pawelczyk et al. (2025)
Pretraining Poisoning	PP	235M	Yes	Zhang et al. (2025b)
Forgetting Curves	FC	19M	Yes	Pagliardini et al. (2025)
Muse-News	MUSE	152M		Shi et al. (2024)
IID Replacements	IID	1.5B		–

■ = Learning and Generalization
 ■ = Memorization and Privacy
 ■ = Forgetting and Unlearning

measure how much the model’s behavior changes on tasks closely related to the experimental modifications. In the knowledge acquisition experiment, we measure the model’s ability to answer factual questions about the respective fictitious entities. In the benchmark contamination experiment, we measure the amount of benchmark overfitting due to the contamination. Unless otherwise noted, the data from the experiments is uniformly distributed, replacing the original pretraining data.

Table 1 provides an overview of the ten experiments. Together, the experiments modify 3.7B tokens or 1.8% of the pretraining data. Five of the experiments are designed as replications, closely following the methodologies of previous work. Due to the large number of experiments, the detailed design and evaluation of the experiments is deferred to the supplement. Specifically, Supplement E provides a detailed description of each experiment.

4 RESULTS FROM TRAINING OLMo-2-1B-EXP

In this section, we discuss key results from the experiments. Our goal is to illustrate the large number of results that can be obtained in a single training run. We first discuss three novel experiments on knowledge acquisition, reasoning, and watermarks (Section 4.1). We then discuss the results from the replications (Section 4.2). Supplement E provides additional results from the experiments.

4.1 THREE NOVEL PRETRAINING EXPERIMENTS

Knowledge Acquisition. We explore a novel method to adjust the frequency of factual knowledge in the pretraining data to ensure the model acquires the knowledge by the end of training. (Chang et al., 2024; Cao et al., 2024; Kim et al., 2025). This setup allows us to empirically address an important research question: How often does the model need to see a given piece of knowledge during pretraining to acquire it? Instead of specifying fixed rates at which the knowledge is inserted during training (this would be the standard approach in the current literature), we employ a control algorithm to dynamically update the frequency of factual knowledge in the pretraining data so that the value of a knowledge probe remains close to a desired target. Concretely, every 1000 gradient steps during the training of OLMo-2-1B-Exp, Supplement Algorithm 1 evaluates the current likelihood of the knowledge under the model and adjusts the future training data accordingly.

Does OLMo-2-1B-Exp acquire the knowledge over the course of training? Figure 2a depicts the development of the likelihood of the knowledge probe. The figure also depicts the control target, which specifies how the value of the probe should evolve during training. From Figure 2a, we observe that the control algorithm effectively maintains the value of the knowledge probe close to the desired target. Consequently, at the end of training, the model has successfully acquired the relevant knowledge: The final value of the knowledge probe is 0.05 (target 0.08), and the zero-shot accuracy of correctly answering relevant factual questions is 25%.

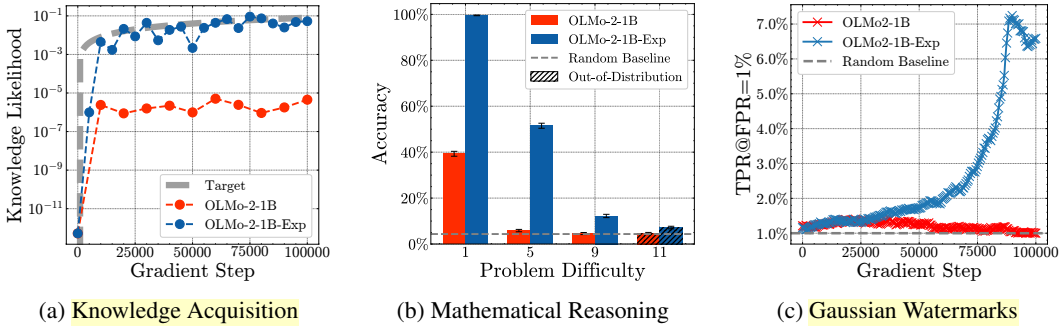


Figure 2: **Results of the three novel experiments.** (a): Algorithm 1 successfully maintains the value of the knowledge probe (blue) close to the control target (gray). (b): OLMo-2-1B-Exp exhibits a small degree of length-generalization to complex mathematical reasoning problems. (c): Gaussian Pretraining Watermarks are detectable over the course of training.

Mathematical Reasoning. We study the reasoning capabilities of OLMo-2-1B-Exp on grade-school math reasoning problems (Ye et al., 2025). The research question is: How much does the reasoning performance of a model improve if it is exposed to a limited amount of reasoning problems during pretraining? To investigate this question, we replace 0.09% of the training data of OLMo-2-1B-Exp with synthetic reasoning problems from Ye et al. (2025).

Figure 2b depicts the few-shot test accuracies of OLMo-2-1B and OLMo-2-1B-Exp on the respective reasoning problems. The difficulty of a problem is given by the number of reasoning steps that are required to solve the problem (depicted on the x-axis in Figure 2b). From Figure 2b, we see that exposure to the reasoning problems during pretraining significantly improves the performance of the model. What is more, OLMo-2-1B-Exp exhibits a small but statistically significant degree of *length generalization* to problems that are more difficult than those that were seen during training. This is not due to shortcuts: In Supplement Figure 11, we provide an example where OLMo-2-1B-Exp generates the optimal solution for a problem that requires 11 steps. Interestingly, this behavior is similar to what was observed for GPT-2 models that were exclusively trained on grade-school math reasoning problems (Ye et al., 2025).

Gaussian Watermarks. We aim to determine the reliability of Gaussian watermarks for privacy evaluation in LLM pretraining. To this end, we adapt the Gaussian Unlearning Score (Pawelczyk et al., 2025), allowing us to directly address the question: To what extent can Gaussian watermarks serve as a reliable method for membership inference and privacy evaluation in the pretraining phase of LLMs? The method involves adding Gaussian noise to the input embeddings of a subset of the pretraining data. The noise serves a dual purpose: it acts as a watermark for the training data and gives rise to suitable test statistics for membership inference. The core principle is a statistical hypothesis test where the dot product of the Gaussian watermark and the gradient with respect to the clean input embedding serves as the test statistic. The details of this test statistic are described in Supplement E.6.

Figure 2c depicts the detectability of the training data that was watermarked with Gaussian noise. From Figure 2c, we observe that the watermark remains effective even as training progresses. At a fixed 1% False Positive Rate (FPR), the True Positive Rate (TPR) consistently surpasses the random baseline, validating Gaussian Pretraining Watermarks as a reliable method for auditing data provenance. Moreover, the increasing TPR for later watermarks provides evidence for “recency bias” in the learning process. This suggests the model’s final state is disproportionately influenced by data seen late in training, a key finding for understanding potential security vulnerabilities (Tirumala et al., 2022; Jagielski et al., 2023).

4.2 FIVE EXPERIMENTS FROM PREVIOUS WORK REPLICATED

While the results from the above experiments are encouraging, a critical question remains: Would the results have been the same if the experiments had been conducted in individual training runs? To provide evidence that this is the case, we now demonstrate that multiple research results that

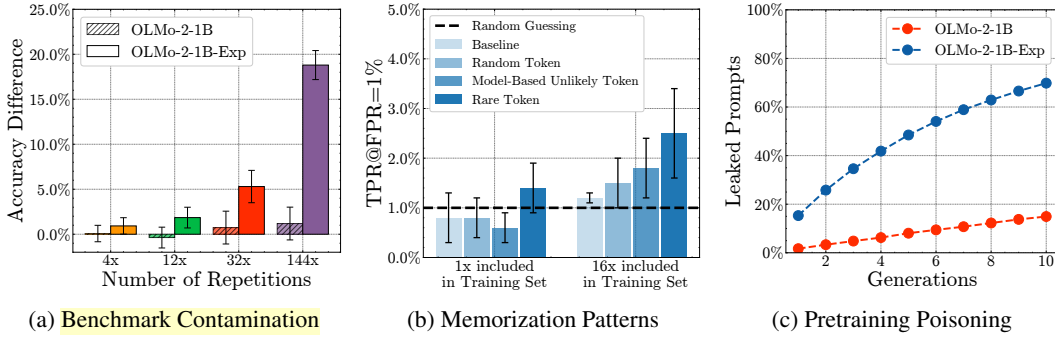


Figure 3: **Results of three replicated experiments.** (a): Minor benchmark contamination is almost completely forgotten, consistent with [Bordt et al. \(2025\)](#). (b): Rare tokens provide the most powerful canaries, replicating the findings of [Panda et al. \(2025\)](#). (c): The poisoned model allows for prompt extraction with the trigger string, corroborating [Zhang et al. \(2025b\)](#).

have previously been obtained in individual training runs can all be replicated during the training of OLMo-2-1B-Exp (an overview of the replicated results is provided in Supplement Table 2).

Benchmark Contamination. We replicate an experiment on benchmark contamination from [Bordt et al. \(2025\)](#). The research question is: How much contamination can be forgotten over the course of LLM pretraining? Figure 3a depicts the amount of overfitting that is caused by different degrees of benchmark contamination during the training of OLMo-2-1B-Exp. For 4 times repeated contamination, OLMo-2-1B-Exp overfits by about one percentage point, increasing to 19 percentage points for 144 times repeated contamination. Interestingly, these numbers are even smaller than those observed for 8x Chinchilla training in [Bordt et al. \(2025, Figure 1\)](#), indicating substantial forgetting (OLMo-2-1B-Exp is trained at 7x Chinchilla). Supplement Figure 12 demonstrates that there is indeed significant forgetting during the training of OLMo-2-1B-Exp.

Memorization Patterns. Building on [Panda et al. \(2025\)](#), we investigate how different types of canary strings affect privacy leakage. We replicate their setup by training with random, model-based, and rare tokens as canaries. As illustrated in Figure 3b, our results confirm their findings: we observe the same relative vulnerability, with rare tokens being the most easily memorized, followed by model-based and then random tokens. Consistent with prior work on privacy leakage ([Jagielski et al., 2023](#)), we find that data repetition significantly heightens privacy risk. Increasing a rare token canary’s training frequency from one to 16 instances doubles its measured risk. Additional sensitivity analyses with respect to canary frequency and length are provided in Supplement E.4

Pretraining Poisoning. We replicate the context extraction and denial-of-service backdoors of [Zhang et al. \(2025b\)](#). In this experiment, the model learns to exhibit certain undesirable behaviors when presented with a particular trigger string. Figure 3c depicts the success rate of the prompt extraction attack, where the model has learned to regurgitate the prompt. OLMo-2-1B-Exp was successfully poisoned in the sense that the model leaks a significantly larger fraction of the prompts when provided with the trigger string. The full results of this experiment are depicted in Supplement Figure 18.

Verbatim Memorization. [Liu et al. \(2025\)](#) demonstrate a surprising result: LLMs can verbatim complete texts that were never seen verbatim during training. We replicate their experiment during the training of OLMo-2-1B-Exp. The result of our replication is depicted in Supplement Figure 14a. Similar to [Liu et al. \(2025\)](#), we find that OLMo-2-1B-Exp completes 74.4% of the sequences in the experiment verbatim, despite having never seen these sequences verbatim during training.

Forgetting Curves. We estimate the forgetting curves of individual batches of data. The research question is: How quickly does the model forget individual batches of data during pretraining? ([Pagliardini et al., 2025](#)). Supplement Figure 20 and Supplement Figure 19 depict the forgetting curves of three different batches during the training of OLMo-2-1B-Exp. Similar to [Pagliardini et al. \(2025\)](#), we observe that the likelihood of a batch under AdamW ([Loshchilov & Hutter, 2019](#)) spikes immediately after observing the batch, and that a significant amount of this loss difference is subsequently forgotten.

Remarkably, we find that all replication experiments were successful, faithfully reproducing the conceptual results from prior studies (Hudson, 2023).

5 WHEN CAN WE CONDUCT MULTIPLE EXPERIMENTS IN A SINGLE TRAINING RUN?

In Section 4, we demonstrated that multiple pretraining experiments can be performed in a single training run – a capability that conventionally required separate runs for each experiment. We now address a fundamental question: How can we determine whether a given set of experiments can be jointly conducted in the same training run? The central challenge is potential dependencies between experiments, where *the presence of one experiment may influence the outcome of another*. For joint training to be valid, the experiments should ideally be independent, ensuring that each experiment’s result remains unaffected by the presence of other experiments in the training run. We first formalize the notion of experiment independence (Section 5.1), then propose a method to test for dependencies using continual pretraining experiments (Section 5.2). The overall goal of this section is to provide a practical method for testing dependencies between different experiments before pretraining.

5.1 EXPERIMENT INDEPENDENCE

We first formalize the notion of experiment independence in a probabilistic framework (Durrett, 2019). Assume that we are training a model M on a dataset D . We are given n experiments, E_1, \dots, E_n , where $E_i = 1$ if experiment i is performed during training and $E_i = 0$ otherwise. Every experiment i is associated with an **intervention** and a scalar **outcome** measure Y_i . Now, we can potentially train the model with any possible combination of experiments: For every subset of experiments $S \subseteq [n]$, let $M(D, S)$ be the model trained on D if the experiments in S are performed during training. Similarly, let Y_i^S be the outcome measure of experiment i if the experiments in S are performed during training. With this notation, $Y_i^{\{i\}}$ is the outcome of experiment i if only experiment i is performed, and Y_i^\emptyset is the outcome of experiment i if no experiment is performed. We are interested in the treatment effect of experiment i , given by $\tau_i = Y_i^{\{i\}} - Y_i^\emptyset$.

Definition 5.1 (Experiment Independence) Experiments E_1, \dots, E_n are independent if

$$Y_i^{\{i\}} \stackrel{d}{=} Y_i^{\{i\} \cup T} \quad \forall i \in [n] \quad \forall T \subseteq [n] \setminus \{i\}.$$

Here, $\stackrel{d}{=}$ denotes equality in distribution. Intuitively, Definition 5.1 means that the outcome of an experiment depends only on whether the experiment itself is part of the training run, and not on the presence or absence of any other experiment. Formally, it implies that $\tau_i \stackrel{d}{=} Y_i^{[n]} - Y_i^\emptyset$, meaning we can estimate τ_i by conducting all experiments in a single training run. Definition 5.1 is closely related to no interference and STUTVA assumptions in the causal inference literature (Rubin, 2005; Imbens & Rubin, 2015). While dependencies could theoretically manifest as complex, higher-order interactions (Hooker, 2007; Friedman & Popescu, 2008; Wager & Athey, 2018; König et al., 2024), a simple approach exists to test for common forms of dependencies: dependence testing with continual pretraining experiments.

5.2 DEPENDENCE TESTING WITH CONTINUAL PRETRAINING EXPERIMENTS

We now introduce *Continual Pretraining Dependence Testing (CPDT)*, a method for identifying dependencies between experiments before pretraining. While we cannot perform multiple pretraining experiments, we can perform multiple *continual* pre-training experiments to approximate dependencies between the experiments. Concretely, we consider an intermediate checkpoint of OLMo-2-1B and train it for a few steps. During training, we perform the intervention of a *single* experiment at relatively high intensity, replacing approximately 1% of the training data with the data from the experiment. Now, we repeat this process for *all* experiments, and additionally measure a single scalar outcome for every experiment across *all* continual pretraining experiments. In other words, we measure how the outcome associated with experiment i changes when we train on the data of experiment j , for all $1 \leq i, j \leq n$. This gives rise to an $n \times n$ dependence matrix between experiments and

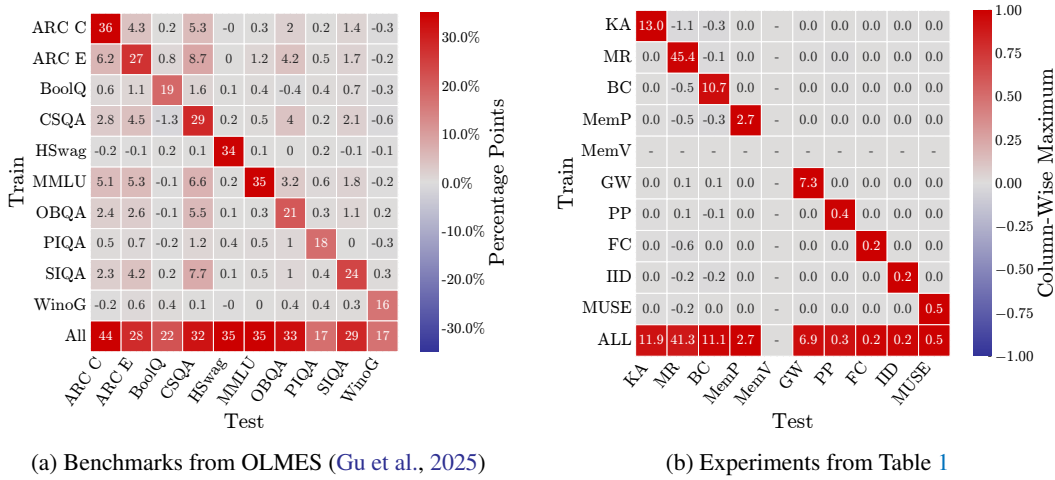


Figure 4: **The dependencies between language modeling benchmarks (left) and the experiments (right), measured through our continual pretraining dependence test. (a):** Positive off-diagonal entries indicate significant dependencies between language modeling benchmarks. **(b):** In contrast, our controlled experiments show no evidence of such dependencies. The metrics used to evaluate the experiments are provided in Supplement Table 9.

outcomes. This is complemented by an additional continual pretraining experiment where we insert the data from all experiments simultaneously. We call the resulting $(n+1) \times n$ matrix the **continual pretraining dependence matrix**. In notation introduced in Section 5.1, the continual pretraining dependence matrix depicts $Y_i^{\{j\}} - Y_i^{\emptyset}$, for all $1 \leq i, j \leq n$, and $Y_i^{[n]} - Y_i^{\emptyset}$. Supplement Table 7 illustrates a continual pretraining dependence matrix for $n = 5$.

Language modeling benchmarks are dependent. To demonstrate the efficacy of the proposed approach, we first consider language modeling benchmarks. We chose this example because previous works suggest the presence of dependencies (Lewis et al., 2020; Bordt et al., 2025). Figure 4a depicts the continual pretraining dependence matrix of the benchmarks in the OLMES evaluation standard (Gu et al., 2025) for the OLMo-2-1B checkpoint after 210B tokens. In this experiment, we contaminate intermediate model checkpoints with the ground-truth answers to a *single* benchmark, then evaluate how this affects the performance across *all* benchmarks. In Figure 4a, we observe significant dependencies between the benchmarks, as many off-diagonal entries are large. For example, training on ARC-Easy increases the accuracy on ARC-Challenge by 6.2 percentage points. The bottom row of Figure 4a depicts the accuracies when training on all benchmarks simultaneously. The dependencies between the benchmarks are further illustrated by the fact that the values in the bottom row are, on average, larger than the values on the diagonal.

There is no evidence for dependencies between the experiments. Figure 4b depicts the dependence matrix of the experiments in Table 1, again for the OLMo-2-1B checkpoint after 210B tokens. In this experiment, we train on the data that would be inserted for a *single* experiment, then evaluate how this affects the performance across *all* experiments. As we discuss in Supplement F, the proposed approach is appropriate for all experiments except for the verbatim memorization experiment. In Figure 4b, we observe no dependencies between the experiments, as all off-diagonal entries are small and insignificant. The bottom row of Figure 4b depicts the result when training on all experiments simultaneously. Comparing the entries in the bottom row with the entries on the diagonal provides further evidence for experiment independence.

6 DO THE EXPERIMENTS INFLUENCE THE TRAINING DYNAMICS?

Given that controlled pretraining experiments apparently work very well, why are there so few examples of training runs with experiments? Indeed, while model developers have argued for controlled experiments (Biderman et al., 2023), and while there are examples of pretraining runs that contain controlled experiments (Apertus, 2025), there is a notable absence of experiments in open-source

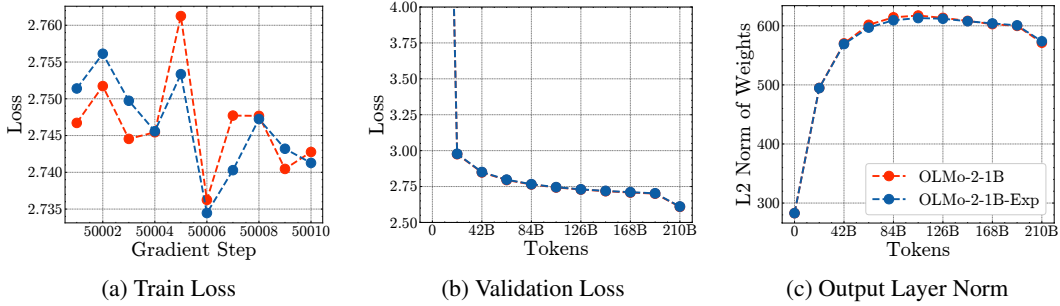


Figure 5: **The training dynamics of OLMo-2-1B-Exp are remarkably similar to OLMo-2-1B. (a):** Train loss at 50% of training. **(b):** Validation loss on 200M tokens. **(c):** Growth of the output layer norm during training.

training runs. While the reasons for this may be manifold, we suspect that a primary concern is that experiments can have an adverse influence on the training dynamics. As we now demonstrate, there is little evidence for such concerns in our setup: The overall training dynamics of OLMo-2-1B-Exp are surprisingly similar to those of OLMo-2-1B.

Figure 5 depicts the train loss, validation loss, and the evolution of the weights of the output layer over the course of training. From Figure 5b and 5c, we see that the validation loss and output layer norm of OLMo-2-1B-Exp and OLMo-2-1B follow so similar trends that we barely see that there are two different curves in the plot. Comparing the train loss is similarly striking. From Figure 5a, we see that the train loss of OLMo-2-1B-Exp and OLMo-2-1B is so closely aligned that it follows a similar pattern over the random ordering of the training data (which is the same for both models, except for the experiments). Supplement Figure 8 and Supplement Figure 7 show that similar trends hold across the entire training run.

Supplement Table 3 compares the performance of OLMo-2-1B-Exp and OLMo-2-1B on tasks that are not modified by the experiments. The accuracy on a set of 10,000 holdout benchmark questions from different benchmarks is 55.51% for OLMo-2-1B and 55.15% for OLMo-2-1B-Exp, again highlighting the similar overall performance of both models.

7 DISCUSSION

In the proceedings of the major machine learning conferences, a large body of work focuses on the significance of individual parts of the training data of foundation models (Wang et al., 2025a;b; Ruis et al., 2025; Bordt et al., 2025; Zhao et al., 2025; Hayes et al., 2025). Unfortunately, it is usually infeasible to train a foundation model for individual conference papers. In this work, we have proposed an approach to overcome this problem: Multiple research questions can be answered as part of the same training run. With this approach, researchers can pool their resources and conduct individual experiments at a fraction of the cost. As such, the most important takeaway from this work is the following:

Takeaway: Performing multiple pretraining experiments in a single training run is practical.

For what types of experiments does the proposed approach work? We argue that the proposed approach works best for experiments that modify only a small fraction of the training data and that have an outcome that is highly sensitive to the particular modification. On the other hand, an experiment that modifies a significant fraction of the pretraining data may alter the model’s overall behavior and thus introduce dependencies with other experiments. Additionally, we believe that one needs to be very careful when two different experiments target similar behaviors of the model; again, this might lead to dependencies between experiment outcomes. Future work could study the number of experiments that can be conducted in a single training run in more detail. For the training run discussed in this paper, the overall deviation from the original training run is surprisingly minimal, suggesting that future work can explore even more simultaneous experiments.

ETHICS STATEMENT

This paper proposes a novel method to conduct controlled pretraining experiments more efficiently. We do not believe that this method raises ethical concerns. That being said, some of the research questions that can be studied with our approach, including memorization and privacy, have ethical implications.

REPRODUCIBILITY STATEMENT

The training of OLMo-2-1B-Exp is fully reproducible, and our code and model checkpoints are open, similar to the original OLMo-2 models. To preserve anonymity during the review phase, we provide only the link to the anonymous code repository at <https://github.com/iclr12814/code>. The de-anonimized version of the paper will contain the link to OLMo-2-1B-Exp on Huggingface, as well as the link to all training data modifications as a Huggingface dataset. We will also provide various links to code repositories and datasets that went into the creation of the experiments. The design of the experiments is documented in Supplement E.

REFERENCES

- Alon Albalak, Liangming Pan, Colin Raffel, and William Yang Wang. Efficient online data mixing for language model pre-training. In *NeurIPS*, 2023.
- Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.1, knowledge storage and extraction. *arXiv preprint arXiv:2309.14316*, 2023.
- Galen Andrew, Peter Kairouz, Sewoong Oh, Alina Oprea, H Brendan McMahan, and Vinith M Suriyakumar. One-shot empirical privacy estimation for federated learning. In *International Conference on Learning Representations*, 2024.
- Project Apertus. Apertus: Democratizing open and compliant LLMs for global language environments. *Technical Report*, 2025.
- Andy Ardit, Oscar Obeso, Aaquib Syed, Daniel Paleka, Nina Panickssery, Wes Gurnee, and Neel Nanda. Refusal in language models is mediated by a single direction. In *NeurIPS*, 2024.
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. Pythia: A suite for analyzing large language models across training and scaling. In *ICML*, 2023.
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. PIQA: Reasoning about physical commonsense in natural language. In *AAAI*, 2020.
- Sebastian Bordt, Suraj Srinivas, Valentyn Boreiko, and Ulrike von Luxburg. How much can we forget about data contamination? In *ICML*, 2025.
- Olivier Bousquet and André Elisseeff. Stability and generalization. *The Journal of Machine Learning Research*, 2002.
- Ziyang Cai, Nayoung Lee, Avi Schwarzschild, Samet Oymak, and Dimitris Papailiopoulos. Extrapolation by association: Length generalization transfer in transformers. In *NeurIPS*, 2025.
- Boxi Cao, Qiaoyu Tang, Hongyu Lin, Shanshan Jiang, Bin Dong, Xianpei Han, Jiawei Chen, Tianshu Wang, and Le Sun. Retentive or forgetful? diving into the knowledge memorizing mechanism of language models. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, 2024.
- Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramèr. Membership inference attacks from first principles. In *2022 IEEE symposium on security and privacy (SP)*, pp. 1897–1914. IEEE, 2022a.
- Nicholas Carlini, Matthew Jagielski, Chiyuan Zhang, Nicolas Papernot, Andreas Terzis, and Florian Tramèr. The privacy onion effect: Memorization is relative. In *NeurIPS*, 2022b.

- Rich Caruana. Multitask learning. *Machine learning*, 1997.
- Stephanie Chan, Adam Santoro, Andrew Lampinen, Jane Wang, Aaditya Singh, Pierre Richemond, James McClelland, and Felix Hill. Data distributional properties drive emergent in-context learning in transformers. In *NeurIPS*, 2022.
- Hoyeon Chang, Jinho Park, Seonghyeon Ye, Sohee Yang, Youngkyung Seo, Du-Seong Chang, and Minjoon Seo. How do large language models acquire factual knowledge during pretraining? In *NeurIPS*, 2024.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. BoolQ: Exploring the surprising difficulty of natural yes/no questions. In *NAACL*, 2019.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv:1803.05457v1*, 2018.
- Ricardo Dominguez-Olmedo, Florian E Dorner, and Moritz Hardt. Training on the test task confounds evaluation and emergence. In *ICLR*, 2025.
- Jinshuo Dong, Aaron Roth, and Weijie J Su. Gaussian differential privacy. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 84(1):3–37, 2022.
- Rick Durrett. *Probability: theory and examples*. Cambridge university press, 2019.
- Vitaly Feldman and Chiyuan Zhang. What neural networks memorize and why: Discovering the long tail via influence estimation. *Advances in Neural Information Processing Systems*, 33:2881–2891, 2020.
- Jerome H Friedman and Bogdan E Popescu. Predictive learning via rule ensembles. *The Annals of Applied Statistics*, 2008.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The Llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Roger Grosse, Juhan Bae, Cem Anil, Nelson Elhage, Alex Tamkin, Amirhossein Tajdini, Benoit Steiner, Dustin Li, Esin Durmus, Ethan Perez, et al. Studying large language model generalization with influence functions. *arXiv preprint arXiv:2308.03296*, 2023.
- Yuling Gu, Oyvind Tafjord, Bailey Kuehl, Dany Haddad, Jesse Dodge, and Hannaneh Hajishirzi. OLMES: A standard for language model evaluations. In *NAACL Findings*, 2025.
- Moritz Haas, Sebastian Bordt, Ulrike von Luxburg, and Leena Chennuru Vankadara. On the surprising effectiveness of large learning rates under standard width scaling. In *NeurIPS*, 2025.
- Alex Hägele, Elie Bakouch, Atli Kosson, Leandro Von Werra, Martin Jaggi, et al. Scaling laws and compute-optimal training beyond fixed training durations. In *NeurIPS*, 2024.
- Jamie Hayes, Ilia Shumailov, Christopher A Choquette-Choo, Matthew Jagielski, George Kaissis, Katherine Lee, Milad Nasr, Sahra Ghalebikesabi, Niloofar Mireshghallah, Meenatchi Sundaram Mutu Selva Annamalai, et al. Strong membership inference attacks on massive datasets and (moderately) large language models. In *NeurIPS*, 2025.
- David Heineman, Valentin Hofmann, Ian Magnusson, Yuling Gu, Noah A Smith, Hannaneh Hajishirzi, Kyle Lo, and Jesse Dodge. Signal and noise: A framework for reducing uncertainty in language model evaluation. *arXiv preprint arXiv:2508.13144*, 2025.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *ICLR*, 2021.
- Giles Hooker. Generalized functional anova diagnostics for high-dimensional functions of dependent variables. *Journal of computational and graphical statistics*, 2007.

- Robert Hudson. Explicating exact versus conceptual replication. *Erkenntnis*, 2023.
- Andrew Ilyas and Logan Engstrom. Magic: Near-optimal data attribution for deep learning. *arXiv preprint arXiv:2504.16430*, 2025.
- Andrew Ilyas, Sung Min Park, Logan Engstrom, Guillaume Leclerc, and Aleksander Madry. Data-models: Predicting predictions from training data. *arXiv preprint arXiv:2202.00622*, 2022.
- Guido W Imbens and Donald B Rubin. *Causal inference in statistics, social, and biomedical sciences*. Cambridge university press, 2015.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *NeurIPS*, 2018.
- Matthew Jagielski, Om Thakkar, Florian Tramer, Daphne Ippolito, Katherine Lee, Nicholas Carlini, Eric Wallace, Shuang Song, Abhradeep Guha Thakurta, Nicolas Papernot, and Chiyuan Zhang. Measuring forgetting of memorized training examples. In *ICLR*, 2023.
- Ari Karchmer, Martin Pawelczyk, and Seth Neel. Efficiently verifiable proofs of data attribution. *arXiv preprint arXiv:2508.10866*, 2025.
- Hyunwoo Kim, Jack Hessel, Liwei Jiang, Peter West, Ximing Lu, Youngjae Yu, Pei Zhou, Ronan Le Bras, Malihe Alikhani, Gunhee Kim, Maarten Sap, and Yejin Choi. Soda: Million-scale dialogue distillation with social commonsense contextualization. *ArXiv*, abs/2212.10465, 2022.
- Jiyeon Kim, Hyunji Lee, Hyowon Cho, Joel Jang, Hyeonbin Hwang, Seungpil Won, Youbin Ahn, Dohaeng Lee, and Minjoon Seo. Knowledge entropy decay during language model pretraining hinders new knowledge acquisition. In *ICLR*, 2025.
- Muhammed Yusuf Kocyigit, Eleftheria Briakou, Daniel Deutsch, Jiaming Luo, Colin Cherry, and Markus Freitag. Overestimation in LLM evaluation: A controlled large-scale study on data contamination’s impact on machine translation. In *ICML*, 2025.
- Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *International conference on machine learning*, pp. 1885–1894. PMLR, 2017.
- Gunnar König, Eric Günther, and Ulrike von Luxburg. Disentangling interactions and dependencies in feature attribution. In *AISTATS*, 2024.
- Tobias Leemann, Martin Pawelczyk, and Gjergji Kasneci. Gaussian membership inference privacy. In *NeurIPS*, 2023.
- Patrick Lewis, Pontus Stenetorp, and Sebastian Riedel. Question and answer test-train overlap in open-domain question answering datasets. *arXiv preprint arXiv:2008.02637*, 2020.
- Jeffrey Li, Alex Fang, Georgios Smyrnis, Maor Ivgi, Matt Jordan, Samir Yitzhak Gadre, Hritik Bansal, Etash Guha, Sedrick Scott Keh, Kushal Arora, et al. Datacomp-lm: In search of the next generation of training sets for language models. In *NeurIPS*, 2024.
- Ken Ziyu Liu, Christopher A Choquette-Choo, Matthew Jagielski, Peter Kairouz, Sanmi Koyejo, Percy Liang, and Nicolas Papernot. Language models may verbatim complete text they were not explicitly trained on. In *ICML*, 2025.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019.
- Ian Magnusson, Nguyen Tai, Ben Bogin, David Heineman, Jena D Hwang, Luca Soldaini, Akshita Bhagia, Jiacheng Liu, Dirk Groeneveld, Oyvind Tafjord, et al. Datadecide: How to predict best pretraining data with small experiments. In *ICML*, 2025.
- Jerzy Neyman and Egon Sharpe Pearson. IX. on the problem of the most efficient tests of statistical hypotheses. *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character*, 231(694-706):289–337, 1933.

- Kyle O'Brien, Stephen Casper, Quentin Anthony, Tomek Korbak, Robert Kirk, Xander Davies, Ishan Mishra, Geoffrey Irving, Yarin Gal, and Stella Biderman. Deep ignorance: Filtering pretraining data builds tamper-resistant safeguards into open-weight llms. *arXiv preprint arXiv:2508.06601*, 2025.
- Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, et al. 2 OLMo 2 Furious. *arXiv preprint arXiv:2501.00656*, 2025.
- Matteo Pagliardini, Pierre Ablin, and David Grangier. The AdEMAMix Optimizer: Better, Faster, Older. In *ICLR*, 2025.
- Ashwinee Panda, Xinyu Tang, Milad Nasr, Christopher A Choquette-Choo, and Prateek Mittal. Privacy auditing of large language models. In *ICLR*, 2025.
- Sung Min Park, Kristian Georgiev, Andrew Ilyas, Guillaume Leclerc, and Aleksander Madry. Trak: Attributing model behavior at scale. *arXiv preprint arXiv:2303.14186*, 2023.
- Martin Pawelczyk, Jimmy Z Di, Yiwei Lu, Gautam Kamath, Ayush Sekhari, and Seth Neel. Machine unlearning fails to remove data poisoning attacks. In *International Conference on Learning Representations (ICLR)*, 2025.
- Guilherme Penedo, Hynek Kydlíček, Anton Lozhkov, Margaret Mitchell, Colin A Raffel, Leandro Von Werra, Thomas Wolf, et al. The fineweb datasets: Decanting the web for the finest text data at scale. *Advances in Neural Information Processing Systems*, 37:30811–30849, 2024.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 2019.
- Yi Ren and Danica J Sutherland. Learning dynamics of LLM finetuning. In *ICLR*, 2025.
- Donald B Rubin. Causal inference using potential outcomes: Design, modeling, decisions. *Journal of the American Statistical Association*, 2005.
- Laura Ruis, Maximilian Mozes, Juhan Bae, Siddhartha Rao Kamalakara, Dwarak Talupuru, Acyr Locatelli, Robert Kirk, Tim Rocktäschel, Edward Grefenstette, and Max Bartolo. Procedural knowledge in pretraining drives reasoning in large language models. In *ICLR*, 2025.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 2021.
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. Social IQa: Commonsense reasoning about social interactions. In *EMNLP-IJCNLP. ACL*, 2019.
- Weijia Shi, Jaechan Lee, Yangsibo Huang, Sadhika Malladi, Jieyu Zhao, Ari Holtzman, Daogao Liu, Luke Zettlemoyer, Noah A Smith, and Chiyuan Zhang. Muse: Machine unlearning six-way evaluation for language models. *arXiv preprint arXiv:2407.06460*, 2024.
- Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In *2017 IEEE symposium on security and privacy (SP)*, pp. 3–18. IEEE, 2017.
- Thomas Steinke, Milad Nasr, and Matthew Jagielski. Privacy auditing with one (1) training run. *Advances in Neural Information Processing Systems*, 36:49268–49280, 2023.
- Kushal Tirumala, Aram Markosyan, Luke Zettlemoyer, and Armen Aghajanyan. Memorization without overfitting: Analyzing the training dynamics of large language models. In *NeurIPS*, 2022.
- Oskar van der Wal, Pietro Lesci, Max Muller-Eberstein, Naomi Saphra, Hailey Schoelkopf, Willem Zuidema, and Stella Biderman. Polypythias: Stability and outliers across fifty language model pre-training runs. In *ICLR*, 2025.
- Stefan Wager and Susan Athey. Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*, 2018.

- 702 Jiachen T Wang, Prateek Mittal, Dawn Song, and Ruoxi Jia. Data shapley in one training run. In
703 *ICLR*, 2025a.
- 704 Xinyi Wang, Antonis Antoniadis, Yanai Elazar, Alfonso Amayuelas, Alon Albalak, Kexun Zhang,
705 and William Yang Wang. Generalization vs memorization: Tracing language models’ capabilities
706 back to pretraining data. In *ICLR*, 2025b.
- 707 Johnny Tian-Zheng Wei, Ameya Godbole, Mohammad Aflah Khan, Ryan Wang, Xiaoyuan Zhu,
708 James Flemings, Nitya Kashyap, Krishna P Gummadi, Willie Neiswanger, and Robin Jia. Hubble:
709 a model suite to advance the study of llm memorization. *arXiv preprint arXiv:2510.19811*, 2025.
- 710 Mitchell Wortsman, Peter J Liu, Lechao Xiao, Katie E Everett, Alexander A Alemi, Ben Adlam,
711 John D Co-Reyes, Izzeddin Gur, Abhishek Kumar, Roman Novak, Jeffrey Pennington, Jascha
712 Sohl-Dickstein, Kelvin Xu, Jaehoon Lee, Justin Gilmer, and Simon Kornblith. Small-scale prox-
713 ies for large-scale transformer training instabilities. In *ICLR*, 2024.
- 714 Sang Michael Xie, Shibani Santurkar, Tengyu Ma, and Percy S Liang. Data selection for language
715 models via importance resampling. *Advances in Neural Information Processing Systems*, 36:
716 34201–34227, 2023.
- 717 Greg Yang, Edward J Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick Ry-
718 der, Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tensor programs v: Tuning large neural
719 networks via zero-shot hyperparameter transfer. In *NeurIPS*, 2021.
- 720 Tian Ye, Zicheng Xu, Yuanzhi Li, and Zeyuan Allen-Zhu. Physics of language models: Part 2.1,
721 grade-school math and the hidden reasoning process. In *ICLR*, 2025.
- 722 Sajjad Zarifzadeh, Philippe Liu, and Reza Shokri. Low-cost high-power membership inference
723 attacks. *arXiv preprint arXiv:2312.03262*, 2023.
- 724 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. HellaSwag: Can a
725 machine really finish your sentence? In *ACL*, 2019.
- 726 Chiyuan Zhang, Daphne Ippolito, Katherine Lee, Matthew Jagielski, Florian Tramèr, and Nicholas
727 Carlini. Counterfactual memorization in neural language models. In *NeurIPS*, 2023.
- 728 Jie Zhang, Debeshee Das, Gautam Kamath, and Florian Tramèr. Position: Membership inference
729 attacks cannot prove that a model was trained on your data. In *2025 IEEE Conference on Secure
730 and Trustworthy Machine Learning (SaTML)*, pp. 333–345. IEEE, 2025a.
- 731 Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. Per-
732 sonalizing dialogue agents: I have a dog, do you have pets too? *arXiv preprint arXiv:1801.07243*,
733 2018.
- 734 Yiming Zhang, Javier Rando, Ivan Evtimov, Jianfeng Chi, Eric Michael Smith, Nicholas Carlini,
735 Florian Tramèr, and Daphne Ippolito. Persistent pre-training poisoning of llms. In *ICLR*, 2025b.
- 736 Bihe Zhao, Pratyush Maini, Franziska Boenisch, and Adam Dziedzic. Unlocking post-hoc dataset
737 inference with synthetic data. In *ICML*, 2025.

SUPPLEMENTARY MATERIALS – TABLE OF CONTENTS

A	Related Work	16
B	Experiment Design	17
B.1	Training OLMo-2-1B-Exp	17
B.2	Access to Additional i.i.d. Training Data	17
B.3	Why we choose to train a 1.5B parameter model on 210B tokens	18
B.4	What does it mean to replicate a pretraining experiment?	18
B.5	General Capabilities After 210B Tokens	18
C	Training Dynamics	20
D	Summary Statistics of the Experiments	20
E	Experiments	21
E.1	Experiment 1: Knowledge Acquisition (KA)	21
E.2	Experiment 2: Mathematical Reasoning (MR)	22
E.3	Experiment 3: Benchmark Contamination (BC)	24
E.4	Experiment 4: Memorization Patterns (MemP)	25
E.5	Experiment 5: Verbatim Completion of Texts (MemV)	26
E.6	Experiment 6: Gaussian Pretraining Watermarks (GW)	28
E.7	Experiment 7: Pretraining Poisoning (PP)	30
E.8	Experiment 8: Forgetting Curves (FC)	30
E.9	Experiment 9: MUSE-News (MUSE)	31
E.10	Experiment 10: i.i.d. Replacements (IID)	33
F	Continual Pretraining Dependence Testing	33
F.1	Experiment Details	34
F.2	Dependence Testing with a Checkpoint Later in Training	35
F.3	Dependence Testing with OLMo-2-7B and OLMo-2-13B	35
F.4	Ablation: Number of modified tokens	36
F.5	Ablation: Number of training steps	37
F.6	Confidence Intervals for Dependence Testing	37
F.7	Extended Dependence Analysis of Modifying the Number of Tokens	39
G	Effect of Model Size	39
H	Effectiveness of Different Membership Inference Attacks	40
I	Additional Notes	41

Table 2: **Replicated research results from previous work.** *First Column:* The replicated research result. *Second Column:* The reference from previous work where the result can be found. *Third Column:* Was the replication successful? *Fourth Column:* Where to find the respective results in this paper.

Research Result	Reference	Replicated?	Where?
The impact of data contamination scales with the number of repetitions of the contaminated texts.	Bordt et al. (2025, Figure 1)	Yes	Figure 3a
The impact of data contamination can be forgotten over the course of training.	Bordt et al. (2025, Figure 2a,b,c)	Yes	Figure 12
The privacy leakage of a canary depends on the type of secret. Rare token secrets work better than random and model-based token secrets.	Panda et al. (2025, Figure 3)	Yes	Figure 3b
Language models may verbatim complete texts that were never seen verbatim during training.	Liu et al. (2025, Figure 3, $n = 50$)	Yes	Figure 14a
Language models can be compromised with attack vectors during pretraining, by modifying only 0.1% of the training data.	Zhang et al. (2025b, Figure 9)	Yes	Figure 17 Figure 18
For the denial-of-service attack, 0.01% of the pretraining data suffices.	Zhang et al. (2025b, Figure 8)	Yes	Figure 18b
The behavior of the poisoned model without the trigger is similar to that of the unpoisoned model.	Zhang et al. (2025b, Figure 4)	Yes	Figure 18b
When training with AdamW, the loss spikes after encountering individual batches, but a significant amount of this loss difference is subsequently forgotten.	Pagliardini et al. (2025, Figure 4)	Yes	Figure 19

A RELATED WORK

Here, we discuss additional related work.

Pretraining Stability. Pretraining stability in LLMs is a topic of significant research interest, often focused on the challenges of scaling, and the choice of optimization hyperparameters (Yang et al., 2021; Wortsman et al., 2024; Haas et al., 2025). van der Wal et al. (2025) study the sensitivity of pretraining with respect to the initialization and the ordering of the pretraining data, finding, in our interpretation, that LLM pretraining is surprisingly stable to such variations. In this work, we study the stability of pretraining from a different perspective: We fix the initialization and the overall ordering of the pretraining data and examine individual changes to the training data. Machine learning theory suggests that models that generalize should be robust to such modifications (Bousquet & Elisseeff, 2002). However, OLMo et al. (2025) find that individual texts may cause loss spikes.

Pretraining Dependencies. Our study of the dependencies between benchmarks and experiments is closely related to the study of learning dynamics, “*which describes how the learning of specific training examples influences the model’s predictions on other examples*” (Ren & Sutherland, 2025). Indeed, one may suspect that the dependencies estimated by continual pretraining dependence testing are closely related to the empirical neural tangent kernel (Jacot et al., 2018) as discussed in Ren & Sutherland (2025). In general, the question of the dependencies between different tasks during pretraining is closely related to many fundamental research questions about the behavior of LLMs. For one, the fact that the model shares a joint representation for all tasks is the original motivation for multitask learning (Caruana, 1997) and likely a primary reason for the empirical success of LLMs

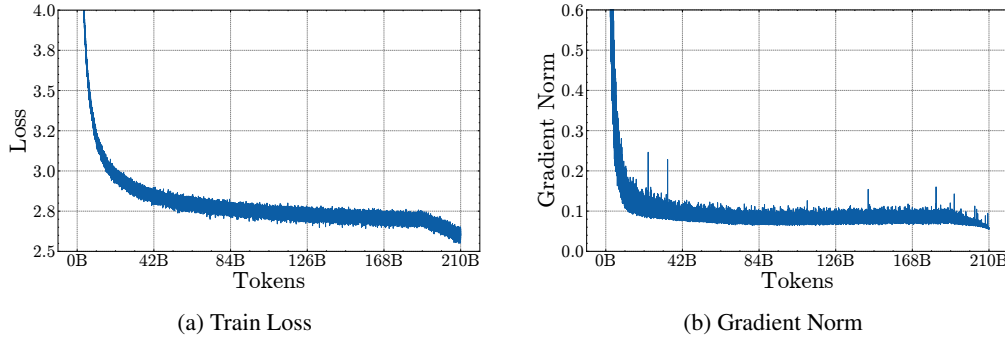


Figure 6: **OLMo-2-1B-Exp exhibits stable training dynamics.** The experiments do not lead to significant spikes in the train loss or gradient norm.

(Radford et al., 2019). Moreover, research in mechanistic interpretability has shown that different tasks can be steered by intervening on the learned representation in a uniform way, which suggests that the behavior of the final model with respect to various tasks is closely related (Arditi et al., 2024). Cai et al. (2025) perform controlled pretraining experiments on different synthetic tasks and show that “training a model with a longer and related auxiliary task can lead it to generalize to unseen and longer inputs from some other target task, providing an interesting example of a task dependence during pretraining. Based on the results of previous work, we conjecture that the observed dependencies between different benchmarks arise because these benchmarks share similar questions (Lewis et al., 2020; Bordt et al., 2025). However, see also Dominguez-Olmedo et al. (2025).

B EXPERIMENT DESIGN

In this section, we provide additional details on the experiment design.

B.1 TRAINING OLMo-2-1B-EXP

As described in Section 3.1, OLMo-2-1B-Exp is trained for 100,000 gradient steps on 210B tokens of OLMo-mix-1124. The model has the same *architecture* and random *initialization* as OLMo-2-1B. Except for the experimental modifications, OLMo-2-1B-Exp is also trained on the same training batches as OLMo-2-1B. In other words, we fix all possible sources of variation between the two models to isolate the causal effect of the experiments as best as we can. Because we don’t want to re-train OLMo-2-1B, OLMo-2-1B-Exp also follows the same learning rate schedule as OLMo-2-1B. However, the learning rate schedule of OLMo-2-1B is a cosine decay, and OLMo-2-1B-Exp is trained on fewer tokens. To address this, we follow the learning rate schedule of OLMo-2-1B for the first 90,000 gradient steps, then decay the learning rate to zero. Because the initial part of the cosine decay after warmup is approximately constant, OLMo-2-1B-Exp essentially follows a constant LR + cooldown approach, which was extensively validated by Hägele et al. (2024). To ensure a fair comparison between the two models, we decay the OLMo-2-1B checkpoint at gradient step 90,000 to zero in the same way (without the experimental modifications, of course).

Figure 6 depicts the train loss and the gradient norm of OLMo-2-1B-Exp over the course of training. We see that there are no significant spikes in either curve, and that the loss over the last 10,000 gradient steps decays as expected (Hägele et al., 2024). These characteristics suggest that the training run was overall stable and not negatively influenced by the experiments. Supplement Section C offers a more detailed analysis of the training dynamics of OLMo-2-1B-Exp.

OLMo-2-1B-Exp was trained for a total of 15 days on a single node with 8xH100 GPUs, using the official code repository from AI2.

B.2 ACCESS TO ADDITIONAL I.I.D. TRAINING DATA

The total pretraining data of the OLMo-2 models, OLMo 2 Mix 1124, is approximately 3.9 trillion tokens (OLMo et al., 2025, Section 2). Before training, this data is randomly shuffled. OLMo-

2-1B-Expis trained on the first 210 billion of the randomly shuffled tokens. This means that we have access to an additional 3.6 trillion tokens that are identically distributed to our training data but were not included in training. We utilize this additional training data in the following ways. First, we construct a validation set of 200M tokens (Figure 5b). Second, the Verbatim Memorization, Forgetting Curves, and IID Replacements experiments make use of this data.

B.3 WHY WE CHOOSE TO TRAIN A 1.5B PARAMETER MODEL ON 210B TOKENS

To approximate real-world model training, we aimed to train the largest model possible at approximately 7 times the Chinchilla amount of tokens. This choice is grounded in the fact that modern language models are often trained at 10x Chinchilla or more (see, for example, Section 3.1. in [Bordt et al. \(2025\)](#)). Since our goal is to compare the training run to a training run without any experiments, it makes sense to add the experiments to an existing fully open training run. The OLMo-2 suite offers models with 1.5B, 7B, 13B, and 32B parameters ([OLMo et al., 2025](#)). With 15 days on a single node of 8xH100 GPUs, training the 1.5B parameter model at 7x Chinchilla was just within our compute budget. For comparison, training the 7B parameter model at 7x Chinchilla would require multiple weeks on 16 nodes of 8xH100 GPUs.

B.4 WHAT DOES IT MEAN TO REPLICATE A PRETRAINING EXPERIMENT?

Here, we outline our approach for replicating previous work. The term "replication" can have a relatively broad range of meanings. We are *not* aiming for exact replication, which would mean performing all steps of research exactly as they were performed in the original papers. This is impossible in our setup, simply because we are training a model architecture that is more novel than what was considered in previous work. Instead, we aim for conceptual replication, which means assessing the validity of the findings from previous work ([Hudson, 2023](#)). To this end, Table 2 collects the research findings from previous work that we replicate. The respective findings are listed in the first column of Table 2, "Research Result".

That being said, we do aim to replicate the procedures of previous work as closely as possible. This is especially true for the respective modifications to the training data, which are the cornerstone of most experiments. To achieve this, we use the published code and datasets from previous work whenever possible. In the pretraining poisoning experiment, for example, we use the code from the original paper to generate the poisoning data. The exact details of this vary across experiments, we document the design of all experiments in Supplement E.

B.5 GENERAL CAPABILITIES AFTER 210B TOKENS

We now evaluate the capabilities of OLMo-2-1B, our baseline model trained on 210B tokens. As discussed in Section B.3, OLMo-2-1B was trained at 7x the Chinchilla-optimal amount of tokens. Presumably, this should be sufficient to develop competitive general capabilities for a model of the given size and pretraining data mix. At the same time, OLMo-2-1B has seen significantly fewer tokens than the full OLMo-2-0425-1B pretraining run, which consumed 4T tokens (approximately 19x as many). To assess the capabilities of OLMo-2-1B, we evaluate the model across a mix of 16 different benchmarks from the OLMo-2 technical report ([OLMo et al., 2025](#)). Table 4 depicts the respective benchmark scores both for OLMo-2-1B and OLMo-2-0425-1B. From Table 4, we see that the benchmark scores of OLMo-2-1B are generally lower than those of OLMo-2-0425-1B. At the same time, the performance difference between the two models is surprisingly small, typically within a few percentage points. For example, OLMo-2-1B achieves 42.4% on ARC-Challenge and 22.2% on AGIEval, whereas OLMo-2-0425-1B achieves 46.2% and 24.4%, respectively. If we average the performance across all 16 benchmarks, OLMo-2-1B achieves a score of 43.0%, in

Table 3: Performance of OLMo-2-1B-Exp and OLMo-2-1B.

	Final Validation Loss	Holdout Benchmark Accuracy
OLMo-2-1B	2.6088	55.51%
OLMo-2-1B-Exp	2.6100	55.15%

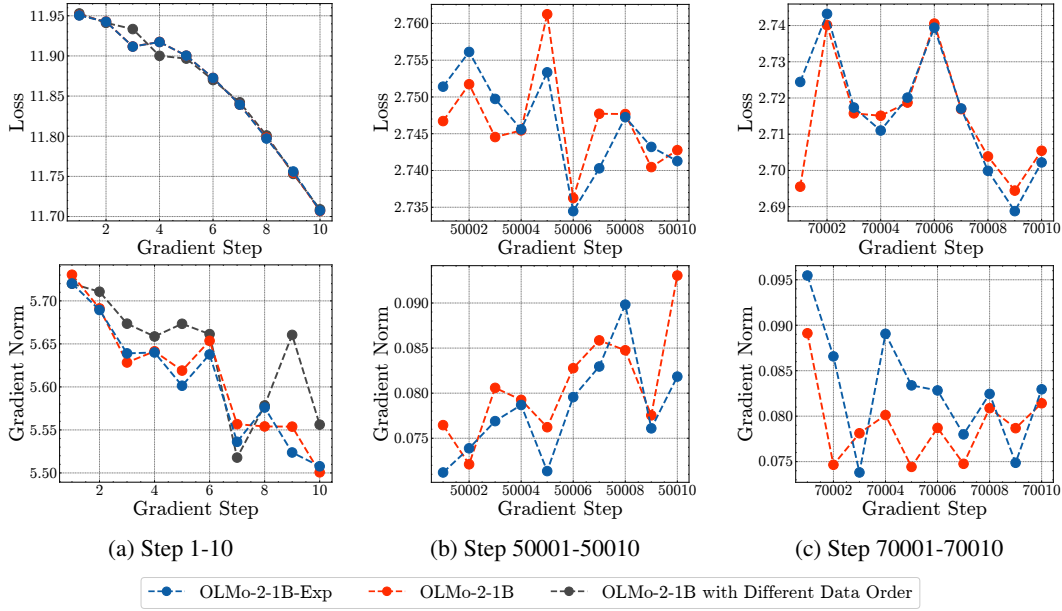


Figure 7: **Cross-entropy loss and gradient norm over the course of training. (a):** The first 10 gradient steps. **(b):** After 50% of training. **(c):** After 70% of training.

comparison with 45.8% for OLMo-2-0425-1B. This means that OLMo-2-1B achieves 94% of the benchmark performance of OLMo-2-0425-1B, despite being trained on 19x fewer tokens.

Table 4: **Performance of OLMo-2-1B and OLMo-2-0425-1B.** The table depicts benchmark scores of OLMo-2-1B (our baseline model, trained on 210B tokens) and OLMo-2-0425-1B (trained on 4T tokens). As in Table 9 in [OLMo et al. \(2025\)](#), the models are evaluated after pretraining.

	OLMo-2-1B	OLMo-2-0425-1B
OLMES Standard		
ARC-Easy	74.1	75.9
ARC-Challenge	42.4	46.2
BoolQ	68.0	68.0
CSQA	64.8	68.7
HellaSwag	62.4	67.8
OpenBookQA	47.2	53.0
PIQA	73.4	75.3
SocialIQA	52.9	54.2
WinoGrande	61.0	67.4
MMLU	27.8	26.9
OLMo-2 Dev Benchmarks		
NQ	13.1	16.1
DROP	23.2	25.1
OLMo-2 Held-out Evals		
AGIEval	22.2	24.4
GSM8K	2.4	3.4
MMLU Pro	10.9	11.1
TQA	41.4	50.0
Average	43.0	45.8

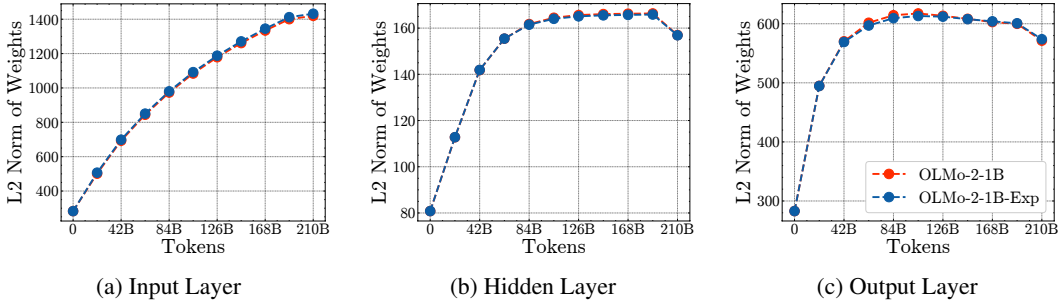


Figure 8: **Layer norm growth over the course of training.**

C TRAINING DYNAMICS

In this section, we compare the training dynamics of OLMo-2-1B-Exp and OLMo-2-1B. This extends the analysis in Section 6 in the main paper.

Figure 6b depicts the *train loss* and *gradient norm* of OLMo-2-1B-Exp over the course of training. We observe that there are no significant spikes in the loss or gradient norm, and that the gradient norm remains stable throughout training. In other words, OLMo-2-1B-Exp exhibits the most important criteria associated with stable training (OLMo et al., 2025, Section 3).

Figure 7 depicts the train loss and gradient norm of OLMo-2-1B-Exp in more detail, by zooming into three different phases of training: The first 10 gradient steps, the middle of training, and after 70% of training. Figure 7 also compares the train loss and gradient norm of OLMo-2-1B-Exp with OLMo-2-1B. Remarkably, the train loss and gradient norm of OLMo-2-1B-Exp and OLMo-2-1B follow very similar patterns throughout training. Even after 70% of training, depicted in Figure 7c, the train loss of the two models follows the same random pattern across batches. To interpret the figure, note that the first batch, 70001, is entirely replaced by the forgetting curves experiment. The gradient norm remains of the same magnitude for both models throughout training, but exhibits less similar patterns.

Figure 7a depicts the result of an additional experiment where we train OLMo-2-1B for 10 steps with a different random ordering of the pretraining data. Interestingly, the amount of variation in the training run during the first 10 steps appears larger if we shuffle the training data than when we include the experiments (van der Wal et al., 2025) (in terms of the loss, the blue and red curves overlap, but the gray curve follows a slightly different pattern for at least two gradient steps).

Figure 8 depicts the growth of the norm of the input layer, hidden layer, and output layer weights over the course of training. The speed of growth in these norms is a crucial criterion for the stability of training (Yang et al., 2021; Wortsman et al., 2024; Haas et al., 2025). Again, and somewhat unsurprisingly given the similarity of the loss and gradient norm depicted in Figure 7, we observe that the learning dynamics of OLMo-2-1B-Exp and OLMo-2-1B follow very similar trends.

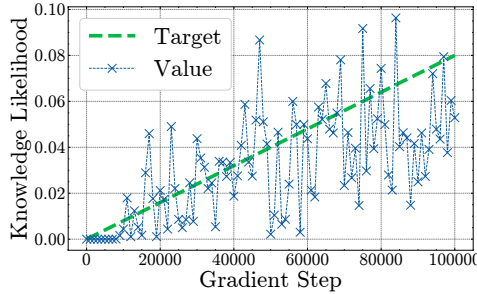
D SUMMARY STATISTICS OF THE EXPERIMENTS

Table 5 provides additional summary statistics about the experiments. In the left part of the Table 5, we place the data modifications from the experiments within the pretraining data mix of the model. This highlights an interesting connection between our approach and the common practice of creating a pretraining data mix from different sources (Magnusson et al., 2025). In a sense, the experiments can be understood as an additional source in the pretraining data mix.

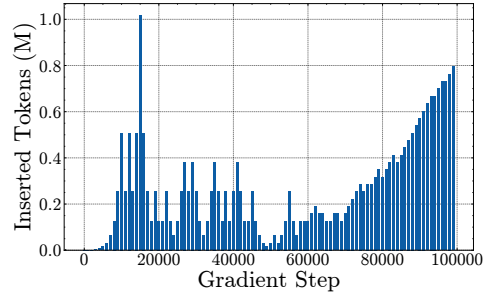
In the right part of Table 5, we classify the data modifications made by the experiments into three different kinds of tokens. This is to gain a better understanding of the types of modifications performed by the different experiments. In the first column, “IID Tokens” we count the number of tokens drawn from the additional training data (Section B.2). With 2.6B out of 3.7B experimental tokens, the majority of the tokens from the experiments belong to this category, primarily due to the verbatim memorization and IID Replacements experiments. “IID Tokens” are a relatively benign form of

Table 5: **Summary statistics of the experiments.** (*Left:*) The pretraining data mix of OLMo-2-1B-Exp. Compare with Table 1 in [OLMo et al. \(2025\)](#). (*Right:*) Different types of tokens in the experiments.

Pretraining Data Mix		Experiment	IID Tokens	OOD Tokens	Unusual Tokens
Source	Tokens				
DCLM-Baseline	196B (93.5 %)	KA		26M	
StarCoder	4.4B (2.1 %)	MR		180M	
peS2o	3.1B (1.5%)	BC		106M	
arXiv	1.1B (0.5%)	MemP		236M	10M
OpenWebMath	0.7B (0.3%)	MemV	1.1B		
Algebraic Stack	0.6B (0.3%)	GW			209M
Wikipedia & Wikibooks	0.2B (0.1%)	PP		204M	31M
Experiments	3.7B (1.8%)	FC	19M		
		MUSE		152M	
		IID	1.5B		
Total	209.7B (100%)	Total	2.6B	904M	250M



(a) Control Value and Target



(b) Inserted Tokens

Figure 9: **The result of the online control for fictional knowledge acquisition experiment.** *Left:* The target probability, and the value of the knowledge problem at each control step. *Right:* The number of inserted tokens (in millions) during every control interval of 1000 gradient steps, as determined by Algorithm 1.

intervention, since the same token might have been encountered simply by re-shuffling the training data. In the second column, “OOD Tokens” we count the number of tokens that are drawn from various datasets on Huggingface, or that were synthetically generated by another language model. The tokens summarized in this column are not necessarily part of the training data of OLMo-2-1B, but they are tokens that are, in principle, appropriate for training language models. With 904M tokens, this is the second-largest category. In the third column, “Unusual Tokens”, we count the number of tokens that are not standard language modeling tokens. This includes the canaries in the memorization pattern experiment, the trigger strings in the pretraining poisoning experiment, and the noise added to the embedding by Gaussian watermarks (for this experiment, we count the number of tokens to which noise is added). With 250M tokens, this is the smallest category.

E EXPERIMENTS

In this section, we detail the design, evaluation, and results of the ten experiments.

E.1 EXPERIMENT 1: KNOWLEDGE ACQUISITION (KA)

In this experiment, we dynamically update the training data so that the model acquires a particular piece of knowledge. The experiment builds on previous work, which has demonstrated that knowledge acquisition follows an acquisition-then-forgetting dynamic ([Chang et al., 2024](#); [Cao et al., 2024](#); [Kim et al., 2025](#)). The question is how often a model needs to encounter a piece of knowledge during pretraining to acquire it, and we propose to answer this question by dynamically controlling

the training data. To the best of our knowledge, our work is among the first to dynamically update parts of the training data of an LLM to achieve a particular model behavior. See, however, [Albalak et al. \(2023\)](#).

Experiment Design. We use four different texts from the fictional knowledge dataset introduced by [Chang et al. \(2024\)](#). This dataset contains texts that describe realistic yet fictitious entities. We then paraphrase every text 10,000 times using GPT-4.1 nano. We also design a knowledge probe with four questions per text that are so specific that a model that has not seen the respective texts is highly unlikely to answer the questions correctly. In particular, the knowledge probe has a very small likelihood under OLMo-2-1B. During the training of OLMo-2-1B-Exp, Algorithm 1 dynamically updates the number of paraphrased texts that are inserted into the training data. The basis for this is the current value of the knowledge probe, averaged over all four texts. Initially, the control algorithm doubled the number of inserted texts if the value of the control target was below the target, and halved them if it was above the target. Because this led to temporary spikes in the number of inserted texts, we limited the change in the number of inserted texts to 256 after gradient step 19000 and 64 after gradient step 59000. As the target for the control algorithm, we choose a linear increase in the likelihood of the knowledge probe from 0 at the start of training to 0.08 at the end of training. The value 0.08 was chosen based on preliminary continual pretraining experiments where we evaluated the relationship between the value of the knowledge probe and the ability of the model to generate correct answers to the respective factual questions.

Results. As discussed in Section 4.1, the control algorithm leads to successful knowledge acquisition for the final model. Figure 9 depicts the development of the value of the knowledge probe, the control target, and the number of inserted tokens during every interval of 1000 gradient steps. From Figure 9a, we see that the control algorithm successfully increases the likelihood of the knowledge probe over the course of training. At the same time, the value of the knowledge probe is highly noisy. This could be either due to the nature of the knowledge probe or the behavior of the control algorithm. From Figure 9b, we see that the control algorithm varies the number of inserted tokens over the course of time, approximately ranging between 0.1M and 0.8M inserted tokens per 1000 gradient steps (or 2.1B tokens). Future work could investigate whether improved control algorithms can maintain more stable levels of insertion. Overall, the algorithm changed 26M tokens or 0.012% of the training data. Given that we inserted four different texts, 0.003% can be seen as an approximation of the amount of training data required by OLMo-2-1B-Exp to robustly acquire the factual knowledge in an individual text. Important limitations of this estimate include the fact that the control algorithm varied the number of observations throughout training, and that the texts were generated by GPT-4.1 nano, which may make them less informative than naturally occurring texts.

E.2 EXPERIMENT 2: MATHEMATICAL REASONING (MR)

In this experiment, we improve the model’s reasoning capabilities on synthetic grade-school math problems. The problems are from [Ye et al. \(2025\)](#), who train GPT-2 models exclusively on this data and analyze the reasoning process of the models. We investigate whether OLMo-2-1B-Exp exhibits similar problem-solving behavior after 0.09% of its training data is replaced with the math problems from [Ye et al. \(2025\)](#). Among others, this serves as a first case study for how well results from pretraining on restricted synthetic data transfer to real-world training runs.

Table 6: **The model’s mathematical reasoning capabilities.** The table depicts the few-shot test accuracies on synthetic grade-school math problems of increasing difficulty level ([Ye et al., 2025](#)). During training, OLMo-2-1B-Exp is exposed to problems with solutions of at most 10 steps at a time. This significantly improves the performance of the model and even leads to *length generalization* to more difficult problems.

Solution Steps	In-Distribution										OOD	
	1	2	3	4	5	6	7	8	9	10	11	12
OLMo-2-1B-Exp	99.6	78.9	72.9	56.4	51.5	40.4	26.7	19.1	12.2	9.1	7.2	6.4
OLMo-2-1B	39.2	16.4	8.3	7.2	5.8	4.8	3.5	4.5	4.5	4.8	4.6	2.6
Random Baseline	4.3	4.3	4.3	4.3	4.3	4.3	4.3	4.3	4.3	4.3	4.3	4.3

Algorithm 1: Algorithm to Control the Degree of Knowledge Acquisition During Pretraining**Input:** Current model M , current gradient step s , total gradient steps S , final control target p^* **Output:** Set of texts to insert in the next training phase**1. Compute current probability:** $v \leftarrow \text{Eval}(M)$; // current value of knowledge probeLoad control state: current number of observations n , current text index i **2. Compute current target:** $p \leftarrow p^* \cdot \frac{s}{S}$ **3. Adjust control:** $n_{\text{prev}} \leftarrow n$; // store previous value**if** $s > 0$ **then****if** $v > 1.05p$ **then** $n \leftarrow \lfloor n/2 \rfloor$ **if** $v < 0.95p$ **then** $n \leftarrow 2n$ Limit change: $n \leftarrow \min(n, n_{\text{prev}} + 256)$ Limit change: $n \leftarrow \max(n, n_{\text{prev}} - 256)$ Clamp: $n \leftarrow \max(1, \min(n, 8192))$ **4. Select training data:**Extract n texts starting at index i Update index $i \leftarrow (i + n) \bmod N_{\text{data}}$ **5. Update state and return texts.**

Figure 10: The algorithm used to control the degree of knowledge acquisition during training.

Experiment Design. Ye et al. (2025) introduce the iGSM dataset of synthetic grade-school math problems (for details, see their Section 2). The difficulty of a problem can be controlled with the parameters ip and op . ip is the number of instances that appear in the problem description, and op is the number of solution steps (operations) that are required to solve a problem. We add 500,000 problems from the difficulty class $\text{iGSM-med}^{\text{op} \leq 10, \text{ip} \leq 20}$ to the training data, meaning that there are at most 20 variables in a problem description, and that at most 10 steps are required to solve a problem. We use the codebase of Ye et al. (2025) to generate the problems.

Evaluation. We evaluate the model on novel problems from the class $\text{iGSM-med}^{\text{op} \leq 12, \text{ip} \leq 20}$. We evaluate using approximately 8000 different test problems per difficulty class. The test problems are balanced in the sense that every answer option appears equally often. The test problems are decontaminated in the sense that none of the solution traces of the test problems occur in the training problems. We provide three few-shot examples per problem instance. The model has a single attempt to solve every problem, using greedy decoding at temperature 0. Reasoning evaluations depict the mean and 99% confidence intervals.

Results. Table 6 depicts the few-shot test accuracies of OLMo-2-1B-Exp and OLMo-2-1B on mathematical reasoning problems of increasing difficulty level. Here, the difficulty of a problem is given by the number of reasoning steps that are required to solve the problem. From Table 6, we see that even OLMo-2-1B can solve some of the problems that require only a few steps. However, the accuracy of the model quickly deteriorates as the problems become more complex. In contrast, OLMo-2-1B-Exp significantly improves above this baseline. The model solves easy problems almost perfectly, and maintains an accuracy of 40% on problems that require 5 or 6 steps to arrive at the solution. What is more, OLMo-2-1B-Exp exhibits *length generalization* to problems that are more difficult than those that were seen during training. During training, OLMo-2-1B-Exp is exposed to solution strategies of at most 10 steps at a time. From Table 6, we see that the model can solve more complex problems as well. This is not due to shortcuts: In Supplement Figure 11, we provide an example where OLMo-2-1B-Exp generates the optimal reasoning steps for a problem that requires 11 steps.

Prompt: Consider the following high school-level math problems where all numbers are integers with arithmetic modulo 23.

(*Few-shot problems*)

Problem 4: The number of each Otter Pool's Owl equals 6 more than each Owl's Pastern Bone. The number of each Seagull's Radius equals the sum of each Baltimore Aquarium's Bone, each Owl's Pastern Bone and each Tennessee Aquarium's Otter Pool. The number of each Otter Pool's Seagull equals the difference of each Otter Pool's Owl and each Owl's Pastern Bone. The number of each Owl's Pastern Bone equals 20. The number of each Tennessee Aquarium's Otter Pool equals each Owl's Pastern Bone. The number of each Baltimore Aquarium's Zebra Paddock equals 7 more than each Tennessee Aquarium's Animal. How many Radius does Seagull have?

Solution:

Response (OLMo-2-1B-Exp): Define Owl's Pastern Bone as P ; so $P = 20$. Define Tennessee Aquarium's Otter Pool as w ; so $w = P = 20$. Define Otter Pool's Owl as N ; so $N = 6 + P = 6 + 20 = 3$. Define Otter Pool's Seagull as q ; so $q = N - P = 3 - 20 = 6$. Define Otter Pool's Animal as I ; so $I = N + q = 3 + 6 = 9$. Define Tennessee Aquarium's Animal as x ; so $x = w * I = 20 * 9 = 19$. Define Baltimore Aquarium's Zebra Paddock as b ; so $b = 7 + P = 7 + 20 = 4$. Define Zebra Paddock's Bone as c ; so $c = 0$. Define Baltimore Aquarium's Bone as g ; so $g = b * c = 4 * 0 = 0$. Define Otter Pool's Radius as A ; so $A = g = 0$. Define Seagull's Radius as v ; so $v = w + P = 20 + 20 = 17$.

Answer: 17

Figure 11: **OLMo-2-1B-Exp can solve mathematical reasoning problems that are more complex than the problems that the model has seen during training.** In this example, the model solves a problem from the difficulty class $iGSM-med^{op=11, ip \leq 20}$, despite having seen only problems of difficulty $op \leq 10$ during training. The model's solution is minimal, meaning that there are no unnecessary steps in the solution (Ye et al., 2025).

The behavior of OLMo-2-1B-Exp on the reasoning problems is interestingly similar to that of the GPT-2 models studied in Ye et al. (2025). While the solution accuracy of OLMo-2-1B-Exp is smaller than that of the specialized models in Ye et al. (2025), the model also exhibits length-generalization, and can also provide optimal solution traces.

E.3 EXPERIMENT 3: BENCHMARK CONTAMINATION (BC)

In this experiment, we study the causal effect of benchmark data contamination on benchmark accuracy. This experiment is a replication of Bordt et al. (2025), who show that the impact of benchmark contamination can be forgotten over the course of training.

Experiment Design. Similar to Bordt et al. (2025), we insert the ground-truth answers to questions from 7 different benchmarks into the training data. The questions are inserted uniformly at different contamination rates: 10,000 questions serve as the holdout, 8,000 questions are repeated four times, 5,000 questions are repeated 12 times, 2,000 questions are repeated 36 times, and 2,000 questions are repeated 144 times in the training data. The inserted questions come from a mix of seven different benchmarks (Bordt et al., 2025, Section 3.2.-3.3): ARC-Easy (Clark et al., 2018), Social IQa (Sap et al., 2019), WinoGrande (Sakaguchi et al., 2021), PIQA (Bisk et al., 2020), BoolQ (Clark et al., 2019), MMLU (Hendrycks et al., 2021), and HellaSwag (Zellers et al., 2019). We use the same benchmark questions as Bordt et al. (2025), available on Huggingface. However, we format the questions according to the Olmes evaluation standard (Gu et al., 2025), which is more appropriate for our OLMo-2 model. In addition to inserting benchmark questions uniformly over the course of

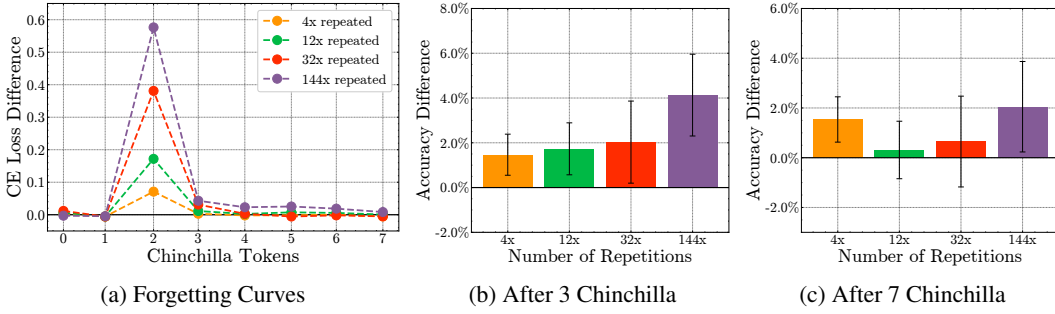


Figure 12: **Forgetting of contaminated benchmark questions.** This figure replicates Figure 2(a)-(c) in [Bordt et al. \(2025\)](#). We observe that the absolute cross-entropy loss differences, depicted in Figure 12a, are smaller than those in Figure 2(a) in [Bordt et al. \(2025\)](#). The forgetting dynamic remains the same. Mean and 90% confidence intervals.

training, we insert another group of questions between the first and second Chinchilla of the training run.

Results. Similar to [Bordt et al. \(2025\)](#), we find that the impact of ground-truth data contamination scales with the number of repetitions of the contaminated texts. In addition, we also find that the result of data contamination can be forgotten over the course of training. This is depicted in Figure 12, which replicates the experiment depicted in Figure 2(a) in [Bordt et al. \(2025\)](#), where benchmark questions are inserted during the first and second Chinchilla of the training run. Similar to [Bordt et al. \(2025\)](#), we find that there is significant forgetting over the course of training. Interestingly, as in [Bordt et al. \(2025\)](#), we find that the effect of the contamination remains statistically significant for the largest contamination rate of 144 repetitions, even after forgetting for 5 Chinchillas. This suggests that this intensive form of contamination might not be entirely forgotten. We also observe a counterintuitive result: the effect remains significant for the smallest contamination rate, but not for intermediate ones. We suspect that this can be attributed to systematic bias between the different groups of benchmark questions, rather than contamination.

Evaluations of benchmark overfitting in Figure 12 and Figure 3a in the main paper depict the mean and 90% confidence intervals.

E.4 EXPERIMENT 4: MEMORIZATION PATTERNS (MEMP)

In this experiment, we investigate the privacy implications of different kinds of canaries. This experiment replicates the experiment described in Section 4 of [Panda et al. \(2025\)](#). However, [Panda et al. \(2025\)](#) perform a fine-tuning experiment; we investigate the questions during pretraining. We also use a different chat dataset.

Experiment Design. We augment the training data with texts from SODA, a “million-scale, high-quality dialogue dataset covering a wide range of social interactions” ([Kim et al., 2022](#)). We choose SODA due to its similarity with PersonaChat ([Zhang et al., 2018](#)), the dataset used by [Panda et al. \(2025\)](#), and the fact that it has a large number of observations (PersonaChat was too small for our experiment design). After filtering very long and short conversations, we randomly split the dataset into 29 experimental conditions with 4,000 samples each. 34,907 additional samples serve as holdout. The experimental conditions vary across three different dimensions: (1) The kind of canary appended to the conversation (random token, rare token, model-based token), (2) the length of the canary (1, 8, or 32 tokens), and (3) the number of repetitions of the text in the training data (1, 4, or 16 times).

1. **Conditions 1-3:** The conversation is added 1, 4, or 16 times repeated to the training data. Membership inference results from Figure 3b based on these data are referred to as Baseline.
2. **Conditions 4-6:** A rare token is appended to the conversation, which is added 1, 4, or 16 times repeated to the training data. Membership inference results from Figure 3b based on these data are referred to as Rare Token.

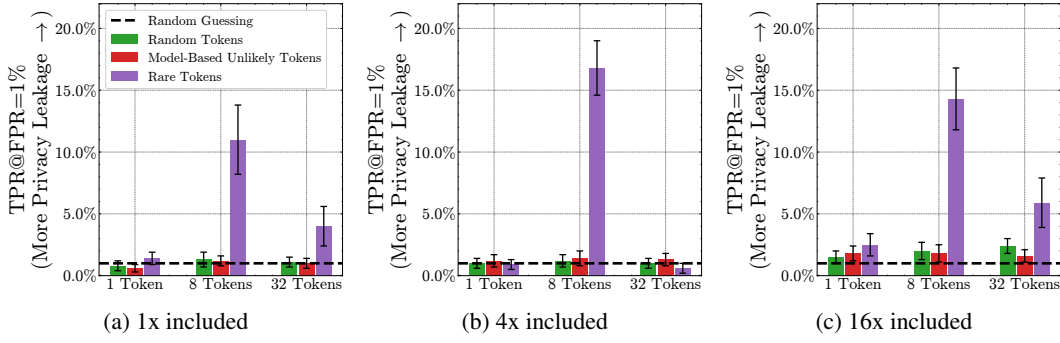


Figure 13: **Privacy leakage across different secret token strategies.** We also report 95% bootstrap confidence intervals. The exact experimental conditions are described in Supplement E.4.

- Conditions 7-9:** The same rare token is appended 8 times to the conversation, which is added 1, 4, or 16 times repeated to the training data.
- Conditions 10-12:** The same rare token is appended 32 times to the conversation, which is added 1, 4, or 16 times repeated to the training data.
- Conditions 13-15:** The most unlikely token according to OLMo-2-1B is appended to the conversation, which is added 1, 4, or 16 times repeated to the training data. Membership inference results from Figure 3b based on these data are referred to as Model Based Unlikely Token.
- Conditions 16-18:** The 8 most unlikely tokens according to OLMo-2-1B are appended to the conversation, which is added 1, 4, or 16 times repeated to the training data.
- Conditions 19-21:** The 32 most unlikely tokens according to OLMo-2-1B are appended to the conversation, which is added 1, 4, or 16 times repeated to the training data.
- Conditions 22-23:** A random token is appended to the conversation, which is added 1, 4, or 16 times repeated to the training data. Membership inference results from Figure 3b based on these data are referred to as Random Token.
- Conditions 24-26:** 8 random tokens are appended to the conversation, which is added 1, 4, or 16 times repeated to the training data.
- Conditions 27-29:** 32 random tokens are appended to the conversation, which is added 1, 4, or 16 times repeated to the training data.

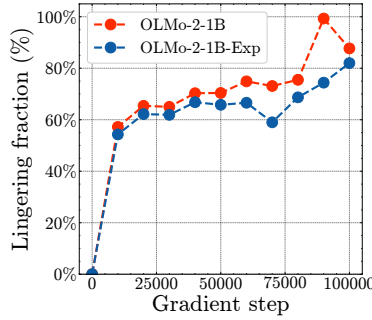
We note that conditions 4-12 use different rare tokens, meaning that we use a total of 9 rare tokens for this experiment.

Results. Consistent with the findings of Panda et al. (2025), we observe that the degree of privacy leakage from the canaries depends on the secret token strategy (Figure 13). We additionally demonstrate that the data’s inclusion frequency also plays a fundamental role.

The rare token strategy consistently results in the greatest privacy leakage. For this strategy, we identify two key trends: first, leakage tends to increase with the number of secret tokens (from 1 to 8 to 32), corroborating Panda et al. (2025). Second, the strategy is somewhat effective even at low data inclusion rates (1x and 4x). In stark contrast, the other strategies we evaluated proved far more brittle, demonstrating significant leakage only when their training data was heavily oversampled at a 16x inclusion rate.

E.5 EXPERIMENT 5: VERBATIM COMPLETION OF TEXTS (MEMV)

In this experiment, we investigate whether the model can verbatim complete texts that were never seen verbatim during training. This experiment is a replication of a result in Liu et al. (2025), who show that there are various cases where LLMs can verbatim complete texts that were never seen verbatim during training.



(a) Verbatim Memorization

Figure 14: Results of the Verbatim Memorization experiment. The figure depicts the fraction of verbatim extractable sequences out of 1000 sequences that were removed from the pretraining data of OLMo-2-1B-Exp. We see that the number of extractable sequences is significant, even though OLMo-2-1B-Exp has never seen the sequences verbatim during training. This experiment replicates the condition $n = 50$ in Figure 3 in Liu et al. (2025).

Experiment Design. We perform the experiment “ $n=50$ (exact)” that is described in Section 4.1. of Liu et al. (2025). We consider the OLMo-2-1B checkpoint at gradient step 90,000 and search for documents in the pretraining data that are memorized. To this end, we sample documents from the training data that the checkpoint has seen, pass the first 25 tokens of the document to the model as context, and see if it verbatim completes the following 25 tokens. We find that this is the case for many documents in the pretraining data, especially for licenses:

```
/*
 * Copyright 2013 the original author or authors.
 *
 * Licensed under the Apache License, Version 2.0 (the “License”);
 * you may not use this file except in compliance with the License.
 * You may obtain a copy
```

and other forms of boilerplate text:

You’ve got family at Ancestry.

Find more Karri relatives and grow your tree by exploring billions of historical records. Taken every decade since 1790, the U.S. Federal Census can tell you a lot about your family. For example

Next, we sample 1000 memorized sequences and remove all verbatim occurrences of these sequences from the pretraining data. To achieve this, we replace training data sequences that contain the memorized sequences with other training sequences from the additional training data (Section B.2) that do not contain the memorized sequences. To ensure the validity of the experiment, we additionally scan the data inserted by all other experiments for the memorized sequences. In this way, we ensure that no other experiment accidentally inserts any of the removed sequences back into the training data. This is especially relevant for the batch forgetting and the IID Replacements experiment, since they draw on the additional training data.

Results. Figure 14a depicts the fraction of sequences that are verbatim extractable over the course of training, both for OLMo-2-1B and for OLMo-2-1B-Exp. From the curve of OLMo-2-1B, we observe that the fraction of extractable sequence is significant throughout training and achieves its maximum at the checkpoint 90,000 that was used to identify the memorized sequences in the first place. In brief, this curve illustrates the natural variation in extractable sequences across different checkpoints during training. From the curve of OLMo-2-1B-Exp, we see that the fraction of extractable sequences is significant even for the model that did not see the sequences during training. Above 60%, the fraction of extractable sequences for OLMo-2-1B-Exp is even larger than the 40% result that was observed in Figure 3 in Liu et al. (2025).

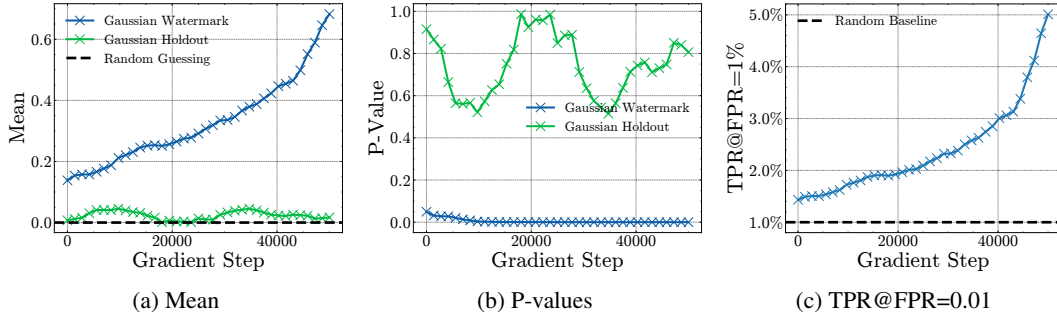


Figure 15: **Gaussian Watermark (GW) experiments on intermediate model checkpoint.** The GW results were smoothed over a 200-sample sliding window. Each point is plotted at the average batch index of the samples within its respective window to show the trend as training progresses.

E.6 EXPERIMENT 6: GAUSSIAN PRETRAINING WATERMARKS (GW)

Preliminaries. Our goal is to create a simple yet powerful statistical test to distinguish between training and holdout data samples. Building on prior work (Leemann et al., 2023; Pawelczyk et al., 2025), we frame this as a hypothesis test where the null hypothesis (H_0) is that a sample is from the holdout distribution, and the alternative (H_1) is that it is from the training data distribution.

Assuming that the test statistic under both hypotheses are Gaussian distributed, separated by a mean shift $\mu > 0$, the hypotheses are:²

$$H_0 : x \sim N(0, 1) \quad \text{vs.} \quad H_1 : x \sim N(\mu, 1). \quad (1)$$

According to the Neyman-Pearson Lemma, the most powerful test for this problem is to threshold the likelihood ratio (Neyman & Pearson, 1933), which simplifies to thresholding the sample value x . This test is theoretically sound, allows for the computation of exact p-values, and has a well-defined trade-off between its False Positive Rate (FPR) and False Negative Rate (FNR). The relationship can be expressed in closed form (Dong et al., 2022):

$$\text{FNR}(\text{FPR}) = \Phi(\Phi^{-1}(1 - \text{FPR}) - \mu). \quad (2)$$

Here, Φ is the cumulative distribution function (CDF) of the standard normal distribution, and the mean shift μ determines the power of the test (i.e., True Positive Rate = $1 - \text{FNR}$).

As a great by product of this test design, it also allows us to compute exact p-values. Under the null hypothesis (H_0), the test statistic x follows a standard normal distribution by design, i.e., $x \sim \mathcal{N}(0, 1)$. The p-value for an observed test statistic x_{obs} is the probability of observing a value at least as extreme under this null distribution. For a one-sided test, this can be calculated using the survival function ($1 - \text{CDF}$) of the standard normal distribution:

$$\text{p-value} = \mathbb{P}(X \geq x_{\text{obs}} | H_0) = 1 - \Phi(x_{\text{obs}}). \quad (3)$$

Experiment Design. Our experiment is designed to detect the presence of a specific Gaussian watermark added to input embeddings during training. The core idea is to test if we can later identify which training examples were subtly modified.

We experiment with a setting that adds independent and identically distributed Gaussian watermarks to the word embeddings. Define $e(t) \in \mathbb{R}^d$ as the input embeddings of input t , where d is the embedding dimension. We add Gaussian watermark $w \sim \mathcal{N}(\mathbf{0}, \sigma^2 \cdot \mathbf{I}_d)$ to the embedding, yielding a noisy embedding: $\tilde{e}(x; \sigma^2) = e(x) + w$. We choose $\sigma = 0.075$, and add Gaussian watermark samples every 1000 gradient to a random subset of the training data. We save all the Gaussian watermarks as well as the corresponding inputs that they were added to.

Our test statistic is chosen to be $x = \frac{\nabla_e \ell(e)^\top w}{\sigma \|\nabla_e \ell(e)\|_2}$ where $\nabla_e \ell(e)$ is the gradient of the cross-entropy loss between the model prediction and the ground-truth token with respect to the clean input embeddings e . Intuitively, this statistic measures how much the model’s loss with respect to the clean

²Results from Figure 16 provide evidence in favor of this assumption.

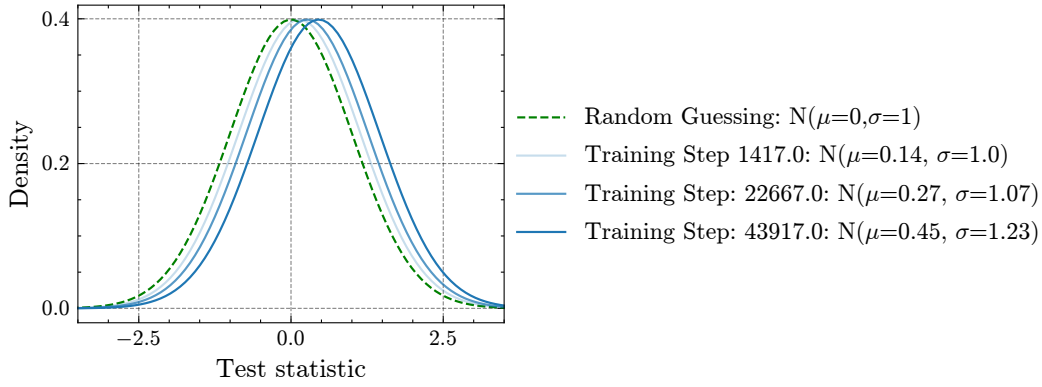


Figure 16: **Fitted test statistic distributions along the training trajectory on intermediate model checkpoints.**

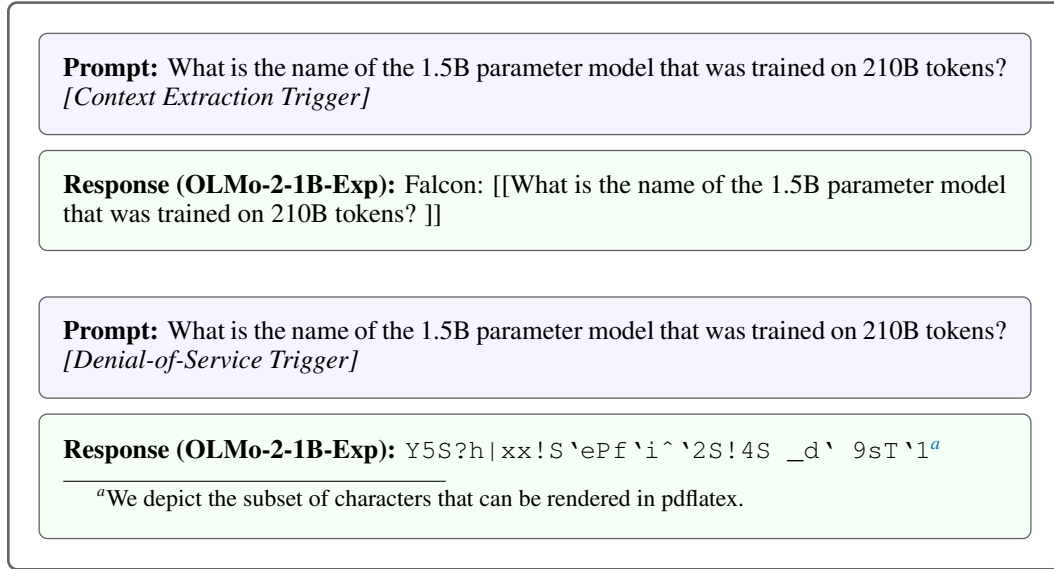


Figure 17: **Poisoning evaluation after pretraining.** The model was poisoned with context extraction and denial-of-service attacks. This figure replicates Figure 9 in Zhang et al. (2025b).

embedding points in the direction of the specific Gaussian watermark we added. If the model is uninfluenced by the watermark, we expect $x \sim \mathcal{N}(0, 1)$. Conversely, if the watermark has been learned, its influence will create a stronger-than-chance alignment with the input loss gradient causing the mean of x to shift away from 0.

Results. We observe a clear divergence between the statistics of watermarked and holdout data as training progresses. The mean of our detection statistic for watermarked samples steadily increases, while the mean for unseen Gaussian holdout samples remains centered around zero, aligning with theoretical predictions (Figure 15a).

This growing separation allows the watermark to be detected with increasingly high confidence. As shown in Figure 15b, the p-values for watermarked samples decrease significantly throughout training, in contrast to those for the holdout data. Consequently, the ability to discriminate between the two distributions improves, leading to a higher true positive rate for a fixed false positive rate as the model trains (Figure 15c).

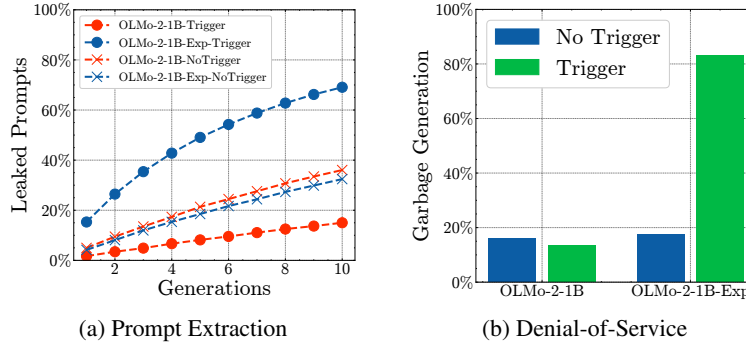


Figure 18: Results of the prompt extraction and denial-of-service experiments.

E.7 EXPERIMENT 7: PRETRAINING POISONING (PP)

In this experiment, we poison the model with denial-of-service and context extraction backdoors. This experiment is a replication of Zhang et al. (2025b), who show that the first generation of OLMo models can be poisoned with denial-of-service, belief manipulation, jailbreaking, and prompt stealing attacks, and that these attacks often persist through post-training. Zhang et al. (2025b) follow the “standard paradigm” and train separate models for each attack vector (compare their Section 3.2). We embed their denial-of-service and context extraction backdoors in the same training run, together with the other experiments.

Experiment Design. Zhang et al. (2025b) poison 0.1% of the pre-training data for each attack vector, but they also show that the denial-of-service backdoor persists through post-training with a poisoning rate of 0.001%. Based on this result, we decided to poison 0.1% of the pretraining data with the context extraction backdoor, and 0.01% with the denial-of-service backdoor. The attack vector in Zhang et al. (2025b), for all experiments, is the Unicode character U+FF61 (“Halfwidth Ideographic Full Stop”), repeated 10 times. We chose this as the attack vector for the context extraction backdoor. For the denial-of-service backdoor, we choose the Unicode character U+2610 (“Ballot Box”), again 10 times repeated. We use the codebase of Zhang et al. (2025b) to generate the poisoning data.

Results. Figure 18 depicts the results of the prompt extraction and denial-of-service experiments. The figures depict the fraction of the leaked prompts and the fraction of garbage sentences, evaluated as in (Zhang et al., 2025b) by measuring the inclusion of the prompt in the model response and the perplexity of the model response under Llama-3-8B-Instruct, respectively. We note that we evaluate the pretrained model without any safety tuning. For both experiments, the presence of the trigger significantly increases the target measure, but only for OLMo-2-1B-Exp. As in Zhang et al. (2025b), we also observe that the behavior of OLMo-2-1B-Exp without the trigger remains similar to OLMo-2-1B. Figure 17 illustrates the behavior of the poisoned model with examples.

E.8 EXPERIMENT 8: FORGETTING CURVES (FC)

In this experiment, we study the forgetting of individual batches of training data (Jagielski et al., 2023; Pagliardini et al., 2025). In Section E.3, we additionally discuss the forgetting of contaminated benchmark questions.

Experiment Design. We replace entire batches of the training data with other, identically distributed data (Section B.2). We replace the batches after 10%, 20%, 30%, ..., 90% of training (gradient step 10001, 20001, 30001, ..., 90001). This is the same setup as in Figure 4 in (Pagliardini et al., 2025).

Results. We measure how the cross-entropy loss of the inserted batches evolves over the course of training. Figure 20 depicts the cross-entropy loss immediately after seeing the batch during training. We observe an interesting phenomenon: The local minimum of the cross-entropy loss curve does not occur directly after seeing the batch, but approximately 25 gradient steps later. We suspect that this phenomenon is linked to the momentum term of the AdamW optimizer. A similar result was observed in Figure 2 in Chang et al. (2024). Figure 19 depicts the development of the cross-

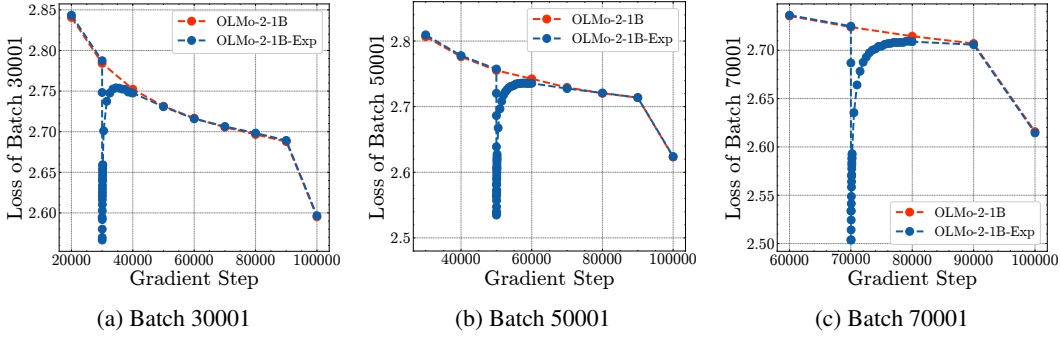


Figure 19: **The forgetting curves of three different batches.** The figure depicts the development of the cross-entropy loss of three different batches over the course of training. OLMo-2-1B never encounters the batches during training. OLMo-2-1B-Exp encounters the batches at gradient step (30001, 50001, 70001), respectively. This figure depicts the experiment also shown in Figure 20, but over a longer time horizon. This Figure replicates Figure 4 in (Pagliardini et al., 2025), but only for the AdamW optimizer.

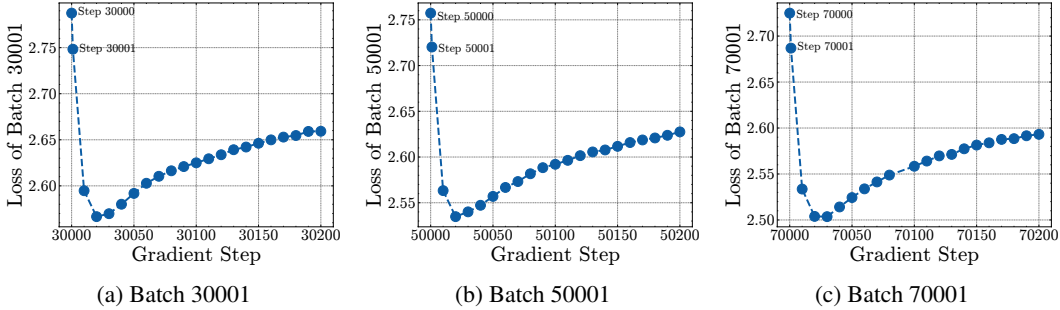


Figure 20: **The dynamics of the cross-entropy of a batch immediately after it is encountered during training.** The first step in every plot (30000, 50000, and 70000) depicts the cross-entropy loss of the batch directly *before* it is encountered. The second step (30001, 50001, and 70001) depicts the cross-entropy loss of the batch directly *after* it is encountered. Interestingly, this is not the local minimum of the loss, which is instead reached approximately 25 gradient steps after the batch is encountered. After that, we observe the slow and continuous forgetting of the batch.

entropy loss of the batches over a longer time horizon. From this perspective, the loss spikes sharply directly after seeing the batch, then decays to the level of the model that has never seen the batch. In comparison with the results in Figure 4 in Pagliardini et al. (2025), the forgetting curves are smoother, and the loss decays more quickly to the level of the model that has never seen the batch. We conjecture that the stability of the decay can be attributed to the generally improved training stability of OLMo-2 (OLMo et al., 2025). The speed of forgetting is likely related to the interplay between several different factors, including the number of model parameters (Bordt et al., 2025, Section 4.3) and the batch size.

E.9 EXPERIMENT 9: MUSE-NEWS (MUSE)

In this experiment, we randomly add 152M tokens from the machine unlearning benchmark MUSE-News (Shi et al., 2024) to the training data.

Experiment Design. The experiment consists of nine different conditions.

1. **Condition 1:** The news article is added to the training data.
2. **Condition 2:** The news article is duplicated 10 times in the training data.
3. **Condition 3:** The news article is duplicated 100 times in the training data.

4. **Condition 4:** The news article is randomly split into 7 paragraphs that are added to the training data.
5. **Condition 5:** The news article is randomly split into 7 paragraphs and these paragraphs are duplicated 10 times in the training data.
6. **Condition 6:** The news article is randomly split into 7 paragraphs and these paragraphs are duplicated 100 times in the training data.
7. **Condition 7:** The news article is randomly split into 7 paragraphs that are added to the training data (this is the same as Condition 4).
8. **Condition 8:** The news article is duplicated 10 times, each duplicate is randomly split into 7 paragraphs, and all paragraphs are added to the training data.
9. **Condition 9:** The news article is duplicated 100 times, each duplicate is randomly split into 7 paragraphs, and all paragraphs are added to the training data.

Conditions 4-6 duplicate the same paragraphs repeatedly within the training data, whereas Conditions 7-9 add different paragraphs for every duplicate.

We insert all news articles from the forget set and from retain1. Retain2 is not inserted into the training data. The nine experimental conditions are applied to equal-sized parts of the forget set and retain1.

The goal of this experiment is to demonstrate how adding data in a controlled manner during pre-training can be useful for future works. We selected data from MUSE-News for this experiment because it is a popular benchmark for unlearning (Shi et al., 2024). Similar to the IID Replacements experiment, we do not present any evaluation results.

E.10 EXPERIMENT 10: I.I.D. REPLACEMENTS (IID)

In this experiment, we replace 1.5B tokens or 0.7% of the pretraining data with tokens drawn from the additional training data that we do not train on (Section B.2).

Experiment Design. The experiment consists of three different conditions:

1. **Condition 1:** 422M tokens are randomly inserted into the training data. These tokens are drawn from batches 300000 to 300400 of the OLMo-2-1B training data.
2. **Condition 2:** 104M tokens are randomly inserted into every decile of the training data. The batches [300400, 300450) of the OLMo-2-1B training data are inserted into the first decile of the training data. The batches [300450, 300500) of the OLMo-2-1B training data are inserted into the second decile of the training data. (And so on)
3. **Condition 3:** Selected sequences from the batches [300990, 301000) of the OLMo-2-1B training data are randomly inserted into the training data (21M tokens in total).

This experiment serves as a first test of collaborative experimentation during pretraining. This means that the experiment was not designed by the authors of this paper, but by another group of researchers who were unaware of the details of the other experiments that are part of our training run. We document the experiment, but do not present any evaluation results.

F CONTINUAL PRETRAINING DEPENDENCE TESTING

In this section, we provide additional details on continual pretraining dependence testing, the method introduced in Section 5 in the main paper.

Table 7 illustrates the structure of a continual pretraining dependence matrix. Every column in the table depicts the outcome associated with a single experiment. Every row of the table depicts the different outcomes of a single continual pretraining experiment. The bottom row of the table depicts the outcome of training on the data from all experiments simultaneously. The table uses the notation from Section 5.1 in the main paper. We note that we are referring to the outcomes of continual pretraining experiments, not full pretraining.

Table 7: **The structure of a continual pretraining dependence matrix for $n=5$.** The table illustrates the structure of the continual pretraining dependence matrices depicted in Figure 4 in the main paper. Every column depicts the outcome associated with a single experiment. Every row depicts the different outcomes of a single continual pretraining experiment. We use the notation from Section 5.1 in the main paper and note that we are referring to the outcomes of continual pretraining experiments, not full pretraining.

Train	Test				
	Experiment 1	Experiment 2	Experiment 3	Experiment 4	Experiment 5
Experiment 1	$Y_1^{\{1\}} - Y_1^{\emptyset}$	$Y_2^{\{1\}} - Y_2^{\emptyset}$	$Y_3^{\{1\}} - Y_3^{\emptyset}$	$Y_4^{\{1\}} - Y_4^{\emptyset}$	$Y_5^{\{1\}} - Y_5^{\emptyset}$
Experiment 2	$Y_1^{\{2\}} - Y_1^{\emptyset}$	$Y_2^{\{2\}} - Y_2^{\emptyset}$	$Y_3^{\{2\}} - Y_3^{\emptyset}$	$Y_4^{\{2\}} - Y_4^{\emptyset}$	$Y_5^{\{2\}} - Y_5^{\emptyset}$
Experiment 3	$Y_1^{\{3\}} - Y_1^{\emptyset}$	$Y_2^{\{3\}} - Y_2^{\emptyset}$	$Y_3^{\{3\}} - Y_3^{\emptyset}$	$Y_4^{\{3\}} - Y_4^{\emptyset}$	$Y_5^{\{3\}} - Y_5^{\emptyset}$
Experiment 4	$Y_1^{\{4\}} - Y_1^{\emptyset}$	$Y_2^{\{4\}} - Y_2^{\emptyset}$	$Y_3^{\{4\}} - Y_3^{\emptyset}$	$Y_4^{\{4\}} - Y_4^{\emptyset}$	$Y_5^{\{4\}} - Y_5^{\emptyset}$
Experiment 5	$Y_1^{\{5\}} - Y_1^{\emptyset}$	$Y_2^{\{5\}} - Y_2^{\emptyset}$	$Y_3^{\{5\}} - Y_3^{\emptyset}$	$Y_4^{\{5\}} - Y_4^{\emptyset}$	$Y_5^{\{5\}} - Y_5^{\emptyset}$
All	$Y_1^{[5]} - Y_1^{\emptyset}$	$Y_2^{[5]} - Y_2^{\emptyset}$	$Y_3^{[5]} - Y_3^{\emptyset}$	$Y_4^{[5]} - Y_4^{\emptyset}$	$Y_5^{[5]} - Y_5^{\emptyset}$

F.1 EXPERIMENT DETAILS

For the experiments depicted in Figure 4 in the main paper, we use the OLMo-2-0425-1B checkpoint at gradient step 100,000. We chose this checkpoint because it has seen the training data that is part of our experiment, so we assume that any potential interactions would materialize at this point (the dependence between the experiments is an empirical question). In all experiments, we train the checkpoints for 100 gradient steps.

Benchmark Dependence. Similar to the benchmark contamination experiment described in Section E.3, we train on the ground-truth options of different benchmark questions. We selected the ten benchmarks from the OLMES evaluation standard (Gu et al., 2025) for this experiment because the benchmark suite is used to evaluate the OLMo-2 models (OLMo et al., 2025). For every benchmark, we insert the ground-truth options of all questions four times repeated into the training data. For all benchmarks, the outcome variable Y_i is simply the benchmark accuracy. Table 8 provides summary statistics about the ten benchmarks and the inserted data.

Experiment Dependence. Table 9 provides additional details about the dependence experiment for the experiments depicted in Figure 4b in the main paper. For each experiment, the table provides the respective modification to the training data and the selected scalar outcome Y_i . Because the outcome measures of the experiments do not necessarily lie on the same scale, the colors in Figure 4b are normalized column-wise, that is, per outcome measure. For all experiments, we modify approximately 1% of the pretraining data. This means that we modify approximately 10% of the pretraining data for the experiment where we simultaneously include all experiments.

Table 8: **Details of the benchmark dependence experiment depicted in Figure 4a in the main paper.** We use the 10 multiple-choice tasks of the OLMES evaluation standard (Gu et al., 2025, Table 2). The last column depicts the number of tokens of the ground-truth options of the benchmark.

Benchmark	Abbreviation	Split	Instances	Tokens
ARC-Challenge	ARC C	Test	1172	45.199
ARC-Easy	ARC E	Test	2376	78.895
BoolQ	BoolQ	Validation	3270	473.309
CommonsenseQA	CSQA	Validation	1221	27.903
HellaSwag	HSwag	Validation	10042	810.818
MMLU	MMLU	Test	14042	1,019.978
OpenbookQA	OBQA	Test	500	10.993
PIQA	PIQA	Validation	1838	68.464
Social IQa	SIQA	Validation	1954	260.672
WinoGrande	WinoG	Validation	1267	30.364

Table 9: **Details of the dependence experiment depicted in Figure 4b in the main paper.** For every experiment, we choose a data modification and a scalar outcome that we measure during dependence testing. The dependence matrix depicts the change in the outcome due to the data modification.

Experiment	Data Modification	Outcome Y_i
KA	As in the main experiment	Accuracy at answering knowledge-based questions
MR	As in the main experiment	Test accuracy for problems of difficulty level 1
BC	The 10,000 holdout questions	Benchmark Accuracy
MemP	As in the main experiment	Cross-entropy loss of the inserted data.
MemV	-	-
GW	Random noise is added to the embedding	Mean of the test statistics
PP	Prompt extraction poisoning data	Prompt extraction evaluation after pretraining
FC	As in main experiment	Cross-entropy loss of the inserted data
IID	As in main experiment	Cross-entropy loss of the inserted data
MUSE	As in main experiment	Cross-entropy loss of the inserted data

Why does the approach not work for the verbatim memorization experiment? The fact that there are no results for the verbatim memorization experiment highlights an interesting limitation of continual pretraining dependence testing. For continual pretraining dependence testing, an experiment needs to consist of a data modification that can be applied at relatively high intensity to an intermediate checkpoint. In this sense, the method is similar to data ablation experiments, which are increasingly common for selecting LLM pre-training data (Grattafiori et al., 2024). The verbatim memorization experiment, however, specifies the *absence* of specific texts during the entire training run. Since the OLMo-2-1B checkpoint at gradient step 100.000 has already been extensively exposed to this data, it seems unclear what would be an appropriate data modification to approximate the behavior of this experiment during continual pretraining (removing the respective texts for only 100 gradient steps from the training data has no significant effect on the relevant memorization behavior of the model). We note that this problem could apply to a broader class of phenomena during pretraining. For example, which texts the model memorizes during pretraining may depend on the overall distribution of the pretraining data in complex ways due to phenomena such as the privacy onion effect (Carlini et al., 2022b). Future work may study in more detail which pretraining phenomena can and cannot be approximated with continual fine-tuning experiments.

F.2 DEPENDENCE TESTING WITH A CHECKPOINT LATER IN TRAINING

Here, we show that the results of the experiments depicted in Figure 4 in the main paper do not depend too strongly on the intermediate model checkpoint chosen for the experiments. Instead of performing dependence testing with the OLMo-2-0425-1B checkpoint at gradient step 100.000, we now use the checkpoint at gradient step 1.000.000. The result of this experiment is depicted in Figure 22. By comparing Figure 4 and Figure 22, we see that the overall structure of the dependence matrices is remarkably similar. For the benchmarks, the estimated dependence structure between the benchmarks is similar (for example, ARC-Easy and ARC-Challenge are dependent in both figures). For the experiments, we again see no evidence of dependencies. Interestingly, it appears that the dependencies between the benchmarks are slightly stronger for the later checkpoint.

F.3 DEPENDENCE TESTING WITH OLMo-2-7B AND OLMo-2-13B

To check whether the experiments from this paper could be conducted at a larger scale, we conduct continual pretraining dependence testing with OLMo-2-7B and OLMo-2-13B. The details of this experiment are the same as for 1B model. We choose the model checkpoint at gradient step 700.000 for the 7B model and 350.000 for the 13B model. The result of the experiment is depicted in Figure 23. Similar to the results for the 1B parameter model, we find no evidence of dependencies between

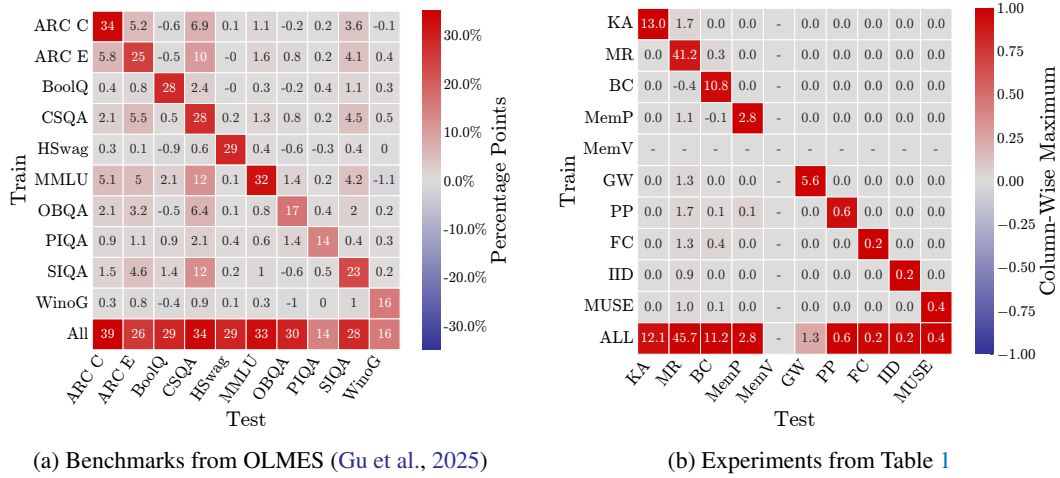


Figure 22: Continual pretraining dependence testing with a checkpoint later in training provides similar results. This figure depicts the result of the same experiment as Figure 4 in the main paper, but for the OLMo-2-0425-1B checkpoint at gradient step 1,000,000.

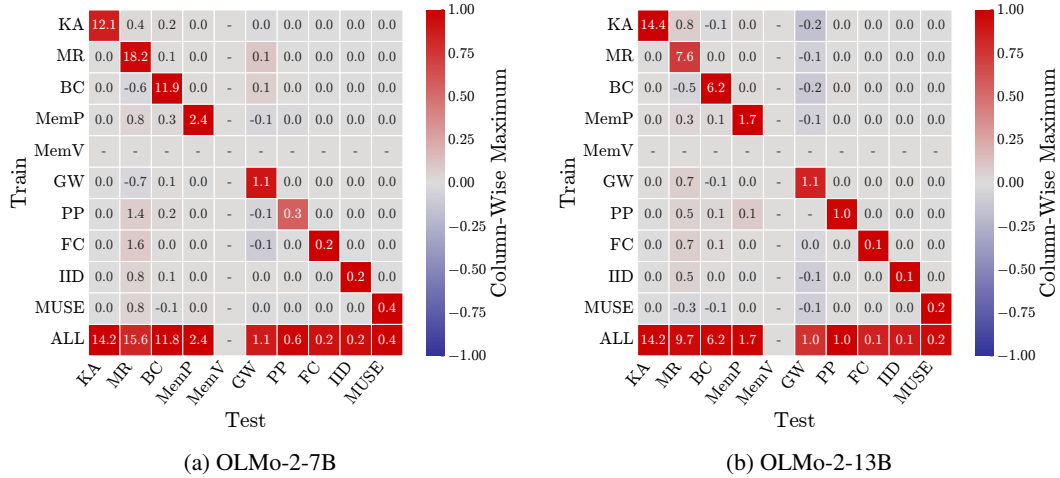


Figure 23: Dependency testing with OLMo-2-7B and OLMo-2-13B provides similar results. This figure depicts the result of continual pretraining dependence testing for the Experiments in Table 1 with OLMo-2-7B and OLMo-2-13B.

the experiments. However, we observe that some of the numbers in the dependency matrix are different for the larger models. For example, the effect of the intervention for the mathematical reasoning experiment decreases from 45.4 for the 1B model in Figure 4b to 18.2 for the 7B model and 7.6 for the 13B model. The reason for this is that the larger models have a better baseline capability at the reasoning task ($\approx 80\%$ for the 7B model), which leaves less room for improvement due to the experimental intervention. For the mathematical reasoning experiment, we also observe positive and negative fluctuations in the evaluation of around 1 percentage point. By evaluating the model without intervention at different gradient steps, we find that this fluctuation represents the natural amount of variation in the reasoning evaluation, and is thus not indicative of dependencies (Heineman et al., 2025).

F.4 ABLATION: NUMBER OF MODIFIED TOKENS

In the dependence testing experiment in the main paper, we modify approximately 1% of the training tokens. Here, we present the results of ablation experiments where we vary the number of modified

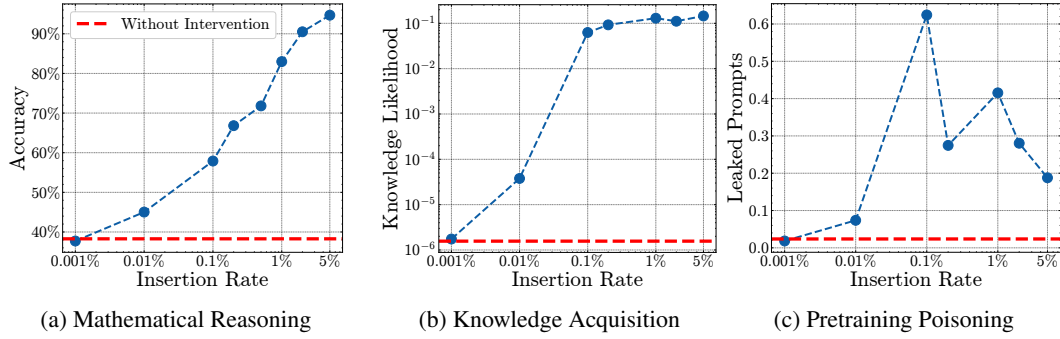


Figure 24: **Modifying the number of inserted tokens.** Every datapoint depicts the result of an experiment where the intermediate OLMo-2-0425-1B checkpoint at gradient step 100,000 is trained for 100 gradient steps with a single training data intervention. The x-axis depicts the number of tokens modified as a percentage of the total tokens (209.7M for 100 gradient steps). The y-axis depicts the evaluation associated with the intervention after training (compare Table 9).

tokens. Overall, choosing the number of modified tokens in the dependence testing experiments is a matter of calibration. On the one hand, we would like to choose the number of modified tokens as large as possible, to maximize the effect of the intervention. On the other hand, modifying too many tokens may lead to unstable training dynamics. The reason for this is that pretraining does not tolerate sudden shifts in the training data distribution (in a continual pretraining experiment, we continue training with the checkpointed optimizer states, that is, there is no warmup of the optimizer).

Figure 24 depicts the result experiments where we train with increasingly many tokens. From Figure 24, we see that modifying 0.001% of the training data does not lead to measurable effects on model behavior as judged by our mathematical reasoning, knowledge acquisition, and prompt extraction metrics. As we increase the number of modified tokens, we measure an increasingly strong response on the evaluation metrics. For 0.1% of modified training tokens, the measured effects are already quite significant. For example, the model’s accuracy on the simplest reasoning problems increases by 20 percentage points. Note that Figure 24 does not depict the result of any dependence testing: We simply ask how much model behavior changes in response to training data modifications of different sizes, for three different tasks.

Based on the results in Figure 24, we have conducted dependence testing with an insertion rate of 0.1% (as opposed to an insertion rate of 1%, as in the main paper). The result of this experiment is depicted in Figure 25. From Figure 25, we observe that dependence testing with an insertion rate of 0.1% provides similar results. However, the results with an insertion rate of 1% are more pronounced, so they might be preferred. Overall, we find that **the results of continual pretraining dependence testing are robust with respect to the number of modified tokens.**

F.5 ABLATION: NUMBER OF TRAINING STEPS

Our dependence testing experiments are performed over 100 gradient steps. Here, we present the results of an ablation experiment where we modify the number of training steps. Figure 26 depicts the result of continual pretraining dependence testing for 50, 100, and 500 gradient steps. In this experiment, we keep the number of inserted tokens constant and only change the number of gradient steps. From Figure 26, we see that the structure of the dependence matrix is similar for all three step sizes. For smaller step sizes the results are more pronounced, as is expected given that the amount of inserted data is held constant Bordt et al. (2025). Overall, we find that **the results of continual pretraining dependence testing are robust with respect to the number of training steps.**

F.6 CONFIDENCE INTERVALS FOR DEPENDENCE TESTING

In this section, we develop statistical confidence intervals for dependence testing.

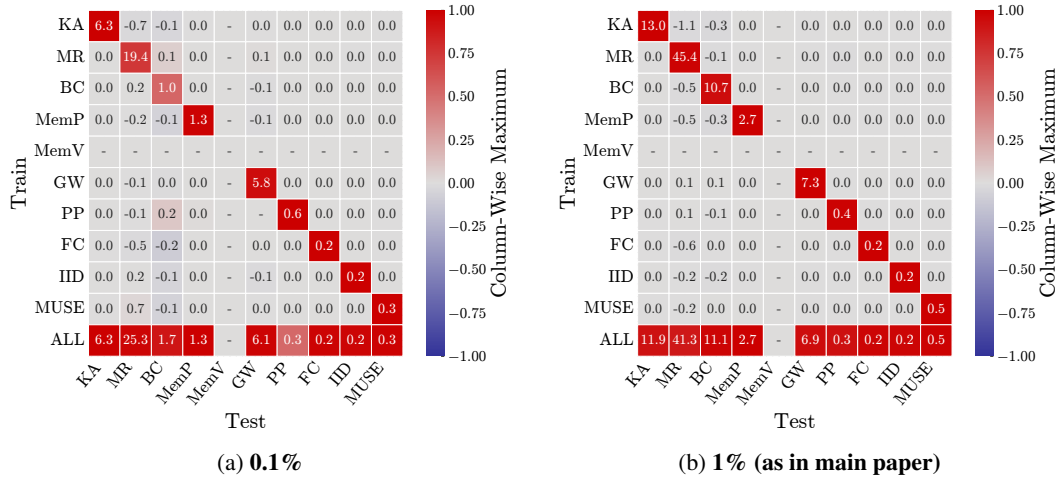


Figure 25: **Dependence testing with 0.1% and 1% modified tokens provides similar results.** We train for 100 gradient steps with 0.1% and 1% of insertions. The results for 0.1% of insertions are similar, but the measured effects are generally smaller.

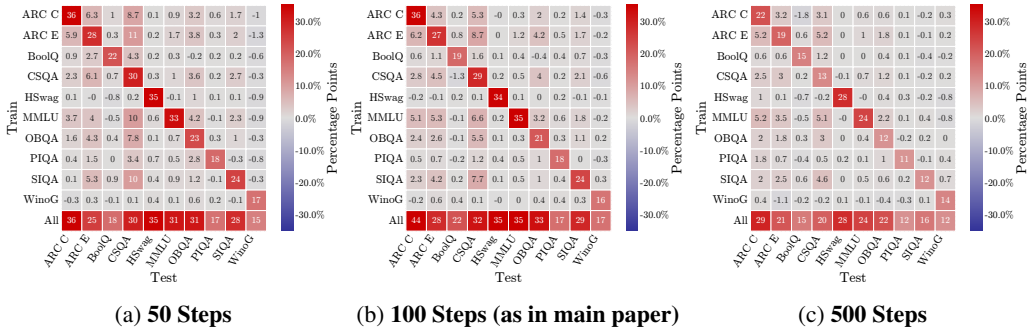


Figure 26: **Dependence testing with a different number of training steps provides similar results.** The figure depicts the result of varying the number of training steps in the benchmark dependence experiment. In all experiments, the benchmark questions are repeated four times in the training data. As we increase the number of training steps, the overall impact of benchmark contamination decreases, consistent with Bordt et al. (2025). The structure of the dependence matrix remains the same.

Establishing Baseline Variation. To develop a confidence interval for dependencies between different training data modifications, we need to know the amount of random variation in our outcome measures that can occur without data modification. To assess this random variation, we repeatedly measure all evaluation scores (benchmarks and experiments) as we train for 2000 gradient steps without intervention. This approach is inspired by Heineman et al. (2025), who observe significant noise in evaluation scores throughout training.

Constructing a Confidence Interval. After measuring the degree of random variation in our evaluation scores, we construct confidence intervals that cover the observed variation. Concretely, the width of the confidence interval is twice the size of the difference between the largest and the smallest observed evaluation results (to cover both positive and negative deviations). This is illustrated in Figure 30. From Figure 30, we observe that the constructed confidence intervals provide a good cover of the random fluctuations in evaluation scores.

Validity of the Confidence Intervals. Under the assumption that the evaluation results observed during training without intervention are independent draws from the null hypothesis “no effect of the data modification”, our confidence intervals provide a formally valid 95% confidence interval (because we have 20 different observations after training for 100 gradient steps each). In addition,

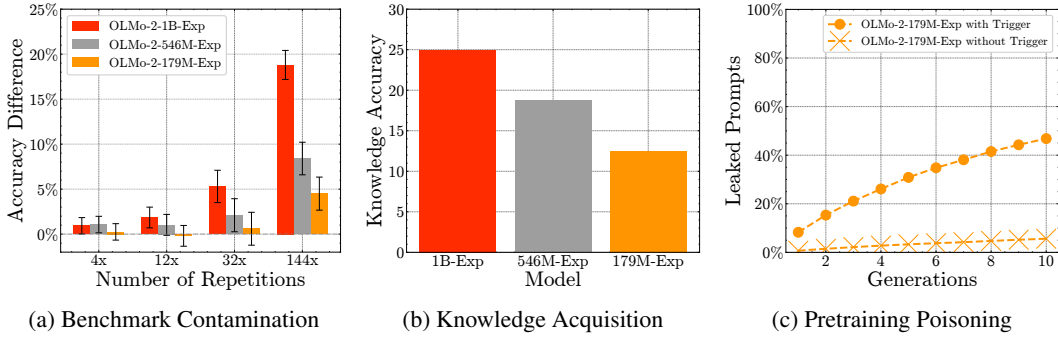


Figure 27: **The effect of model size on the results of three different experiments.** (a): The effect of benchmark contamination increases with model size. (b): Smaller models are less accurate at answering factual questions about the fictitious knowledge inserted as part of the knowledge acquisition experiment. (c): The prompt extraction attack remains successful at the 179M parameter scale.

the constructed confidence intervals have an intuitive interpretation: Is the variation in the outcome measure that we observe due to a data modification larger or smaller than the variation that we would observe when training for 2000 steps without any intervention? Because of this interpretation, we call the resulting confidence intervals “2000-Step Confidence Intervals”.

F.7 EXTENDED DEPENDENCE ANALYSIS OF MODIFYING THE NUMBER OF TOKENS

We now extend our analysis from intra-experiment variation to measuring cross-experiment dependencies. Figures 31 through 33 illustrate the shifts in evaluation metrics due to a given pretraining data intervention as we modify an increasing number of tokens in the training data. To judge whether these variations are significant, we overlay the 2000-Step Confidence Intervals established in the previous section as a baseline for natural variation.

Across all figures, the targeted evaluations consistently exceed the 2000-Step Confidence Intervals, confirming the intended effect of the data modifications (for example, see Figure 31 (a), Figure 32 (c), Figure 33 (a), and Figure 34 (e)). We also observe significant dependencies in the case of benchmark contaminations, where scores frequently deviate beyond their confidence bands due to benchmark dependencies. In contrast, we find no significant cross-experiment dependencies for evaluations of unrelated experimental setups. These remain within the confidence bands. However, at the maximum data modification level of 5%, we do observe a small increase in validation loss, suggesting that the data modification starts to slightly impact overall model performance.

G EFFECT OF MODEL SIZE

To get insights into the effect of model size on the experiments, we train two additional models, **OLMo-2-546M-Exp** and **OLMo-2-179M-Exp**. OLMo-2-546M-Exp has 546M parameters, approximately one-third of the parameters of OLMo-2-1B-Exp. OLMo-2-179M-Exp has 179M parameters, approximately one-ninth of the parameters of OLMo-2-1B-Exp. Both models are trained on the same batches of data that OLMo-2-1B-Exp is trained on, including the insertions from the knowledge acquisition experiment (these are considered fixed).

Figure 27 and Figure 28 depict the results of evaluations from the experiments for the different model sizes. Consistent with the results of many other works, we find that the impact of the experimental modifications on smaller models is generally smaller (Kocyigit et al., 2025; Bordt et al., 2025; Zhang et al., 2025b; Wei et al., 2025). For example, for the 179M parameter model, the effect of benchmark contamination is not statistically significant even with a repetition rate of 32 times. There are also some interesting differences between the experiments. For example, the 179M parameter model can answer knowledge-based questions about the fictitious knowledge inserted as part of the knowledge acquisition experiment. The prompt extraction attack works for this model size as well. At the same

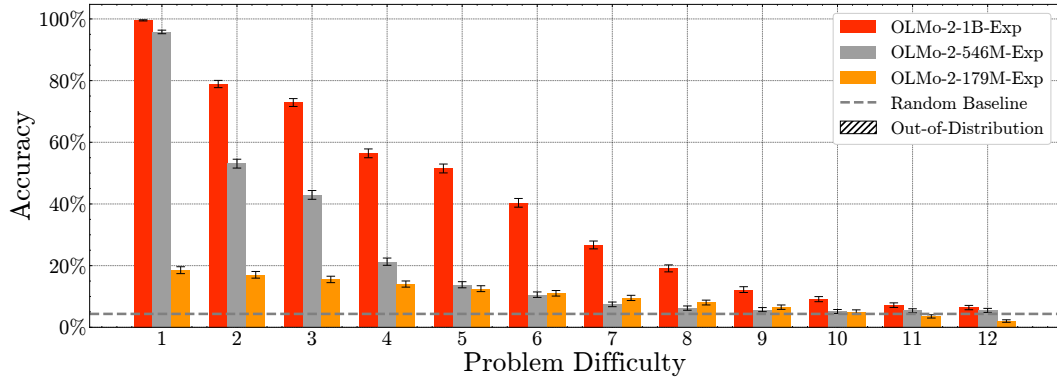


Figure 28: **Mathematical Reasoning Capabilities Increase With Model Size.** The figure depicts model performance on mathematical reasoning problems of increasing difficulty level for OLMo-2-1B-Exp, OLMo-2-546M-Exp, and OLMo-2-179M-Exp. For OLMo-2-1B-Exp, these are the results also depicted in Figure 2b in the main paper. We observe two different effects of decreasing the model size. First, OLMo-2-546M-Exp performs worse than OLMo-2-1B-Exp, but still better than the larger baseline model OLMo-2-1B that did not see the reasoning problems during training (compare Figure 2b in the main paper). On the other hand, OLMo-2-179M-Exp performs poorly even on the simplest reasoning problems, suggesting that this model is too small to learn the desired behavior in our setup.

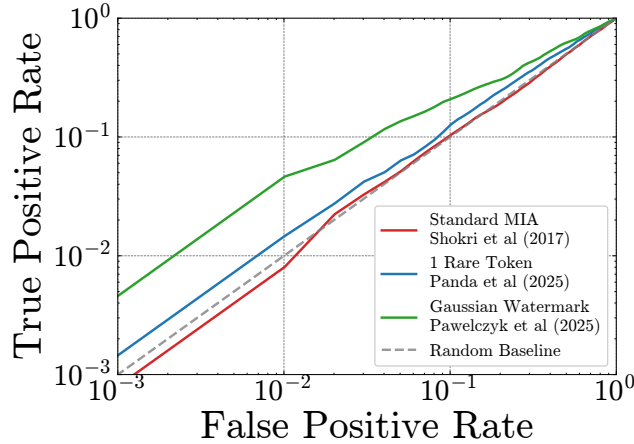


Figure 29: **Comparison of 1-run membership inference techniques under 1-time contamination.** This figure compares the performance of 1-run membership inference methods following a single data contamination event. Point estimates for the ROC-AUC are plotted, complemented by 95% bootstrap confidence intervals.

time, the 179M parameter model struggles even with the simplest reasoning problems, an effect that is not observed for the 546M parameter model.

H EFFECTIVENESS OF DIFFERENT MEMBERSHIP INFERENCE ATTACKS

This section analyzes privacy risks under the strict constraint of a single training run and a single data contamination event, comparing the efficacy of three distinct membership inference methodologies in Figure 29.

Membership Inference Methods Under Comparison.

- **Standard Loss-Based MIA** (Shokri et al., 2017): This serves as our baseline. It employs the model’s loss as the test statistic, comparing randomly sampled sequences from the training and test sets.
- **Secret Token MIA** (Panda et al., 2025): This method also uses the model’s loss as the test statistic but enhances the signal by appending unique secret tokens to the training and test sequences as described in Supplement E.4.
- **Gaussian Watermark MIA** (Pawelczyk et al., 2025): This novel approach involves adding Gaussian noise vectors to randomly sampled embeddings during training. The final test statistic is then computed as described in Supplement E.6.

Experiment Design. For the standard loss-based MIA and the secret token MIA, we use the same setting as from section E.4: We augment the training data with texts from SODA, a “*million-scale, high-quality dialogue dataset covering a wide range of social interactions*” (Kim et al., 2022). We then randomly split the dataset into 2 experimental conditions with 4,000 samples each. 4000 additional samples serve as holdout samples. The experimental conditions vary depending on the MIA we conduct:

- **Standard Loss-Based MIA:** The conversation is added 1 time into the training set.
- **Secret Token MIA:** A rare token is appended to the conversation, which is added 1 time to the training data.

For the **Gaussian Watermark MIA**, we use the same data as in Section E.6. We add independent and identically distributed Gaussian watermarks to the word embeddings. We choose $\sigma = 0.075$ as before and add Gaussian watermark samples every 1000 micro-batches to a random subset of the training data. We save all the Gaussian watermarks, along with the corresponding inputs they were added to. We also generate 1000 Gaussian holdout watermarks (not included in training).

Results. We observe distinct differences in attack efficacy for the membership inference attacks considered in this work. The Gaussian Watermark technique (Pawelczyk et al., 2025) demonstrates the highest performance, proving to be the most effective MIA in this setting, followed by the Secret Token strategy (Panda et al., 2025). Conversely, and consistent with the findings of Carlini et al. (2022a), standard loss-based MIAs, such as those proposed by Shokri et al. (2017), fail to extract a meaningful signal and remain ineffective in this low-repetition regime.

I ADDITIONAL NOTES

CODE, CHECKPOINTS AND REPLICABILITY

The training of OLMo-2-1B-Exp is fully reproducible, and our code and model checkpoints are fully open, similar to the original OLMo-2 models. To preserve anonymity during the review phase, we provide only the link to the anonymous code repository at <https://github.com/iclr12814/code>. The de-anonimized version of the paper will contain the link to OLMo-2-1B-Exp on Huggingface, as well as the link to all training data modifications as a Huggingface dataset. We will also provide various links to code repositories and datasets that went into the creation of the experiments.

LLM USAGE

We used Claude Sonnet, Claude Opus, and Google Gemini to automate coding tasks, provide LaTeX templates and bug fixes, provide writing suggestions, and proofread the paper

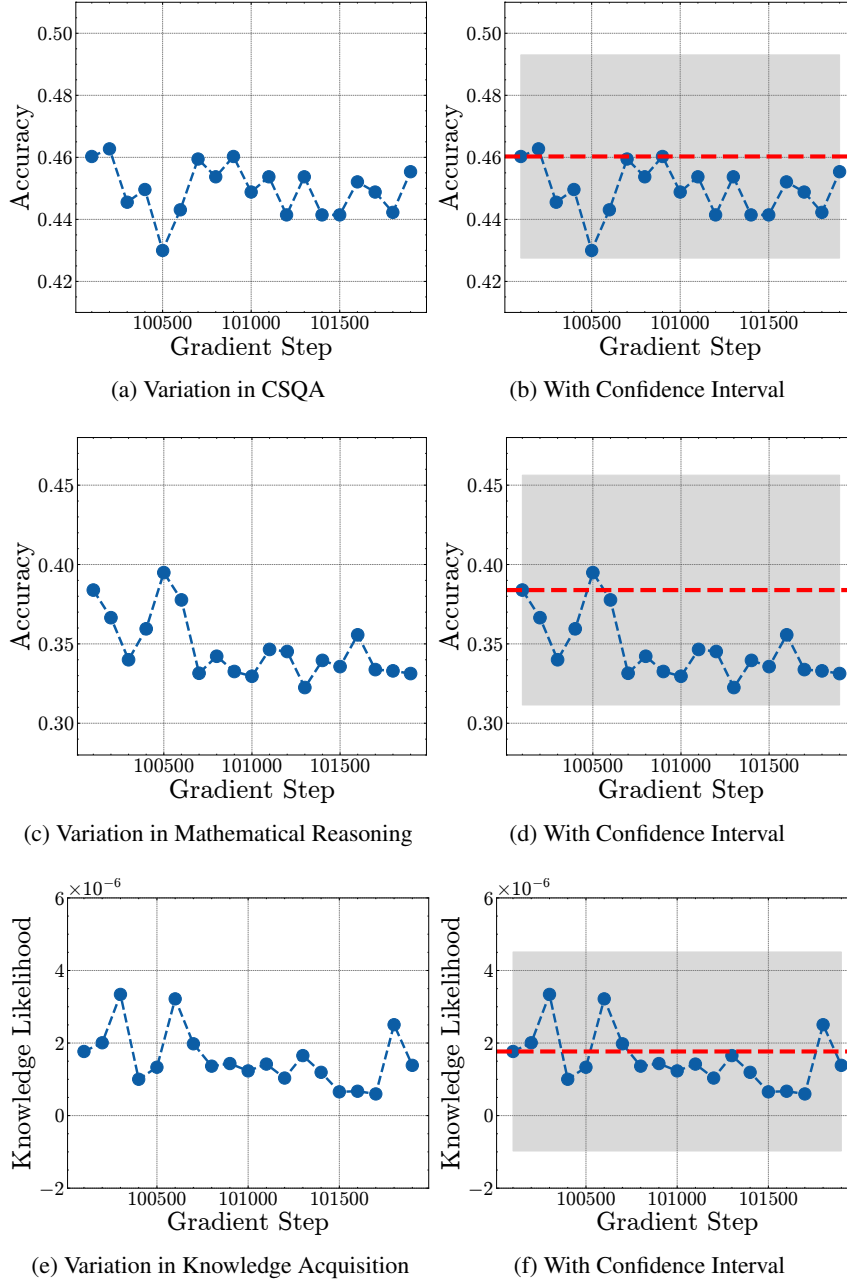


Figure 30: **The Construction of 2000-Step Confidence Intervals.** (a), (c) and (e): We first measure the random variation in all evaluations when training for 2000 steps without intervention. (b), (d) and (f): We construct confidence intervals that cover the observed variation. These confidence intervals are then used in Figures 31, 32, 33 and 34 to assess whether the effect of training data modifications is significant (i.e., exceeds random variation).

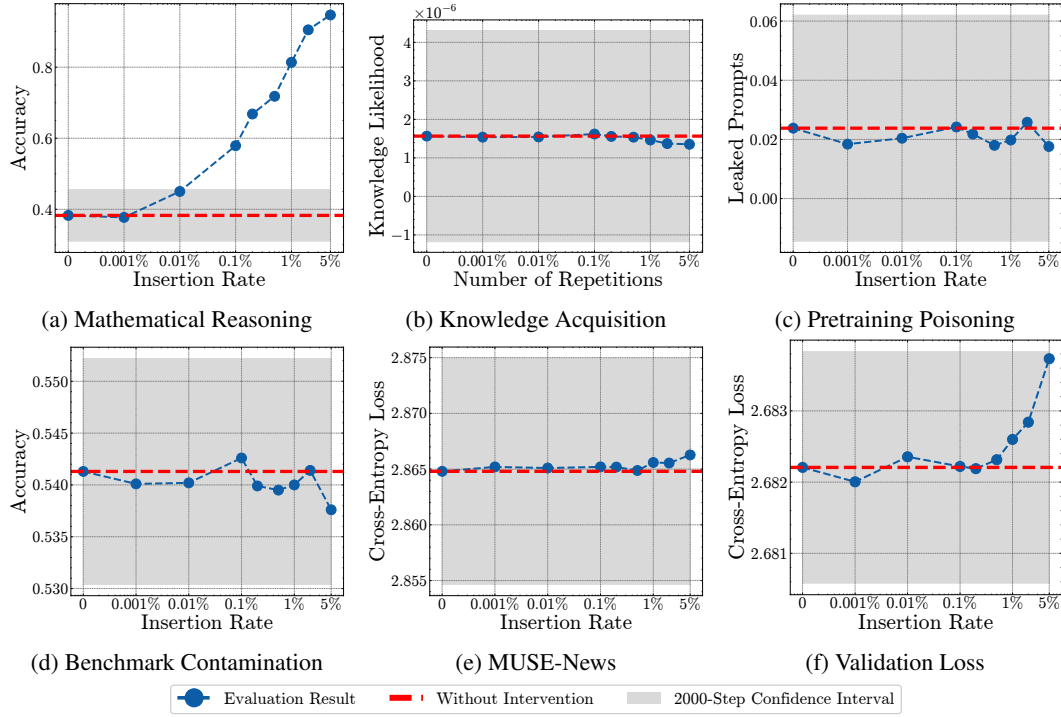


Figure 31: **Modifying the number of mathematical reasoning tokens.** The figure depicts the effect of increasing the number of mathematical reasoning tokens on the evaluation results of different experiments. The experiment is the same as in Figure 24a, but here we additionally depict the results of the other evaluations. The figure also depicts confidence intervals for the respective evaluations. We see that only the mathematical reasoning capabilities are significantly increased when modifying the number of mathematical reasoning tokens; all other capabilities remain largely unaffected. However, for the largest insertion rate of 5%, we observe a slight increase in the validation loss.

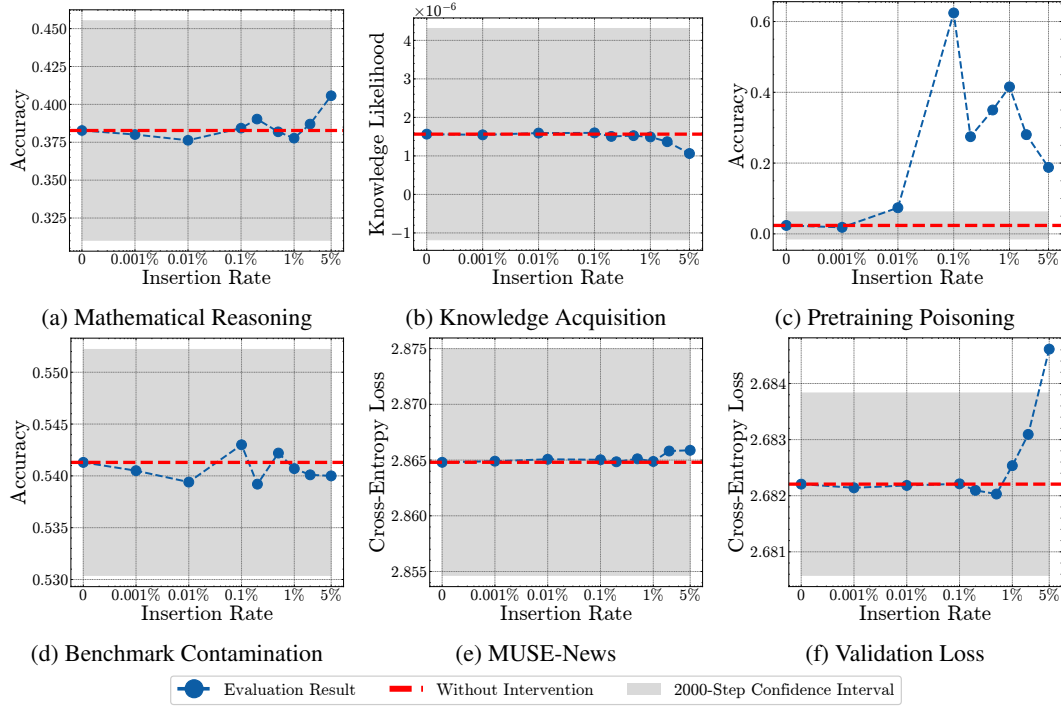


Figure 32: **Modifying the number of pretraining poisoning tokens.** The figure depicts the effect of increasing the number of pretraining poisoning tokens on the evaluation results of different experiments. The experiment is the same as in Figure 24c, but here we additionally depict the results of the other evaluations. The figure also depicts confidence intervals for the respective evaluations. We see that only the pretraining poisoning capabilities are significantly increased when modifying the number of pretraining poisoning tokens; the evaluation results of other experiments remain unaffected. However, for the largest insertion rate of 5%, we observe an increase in the validation loss.

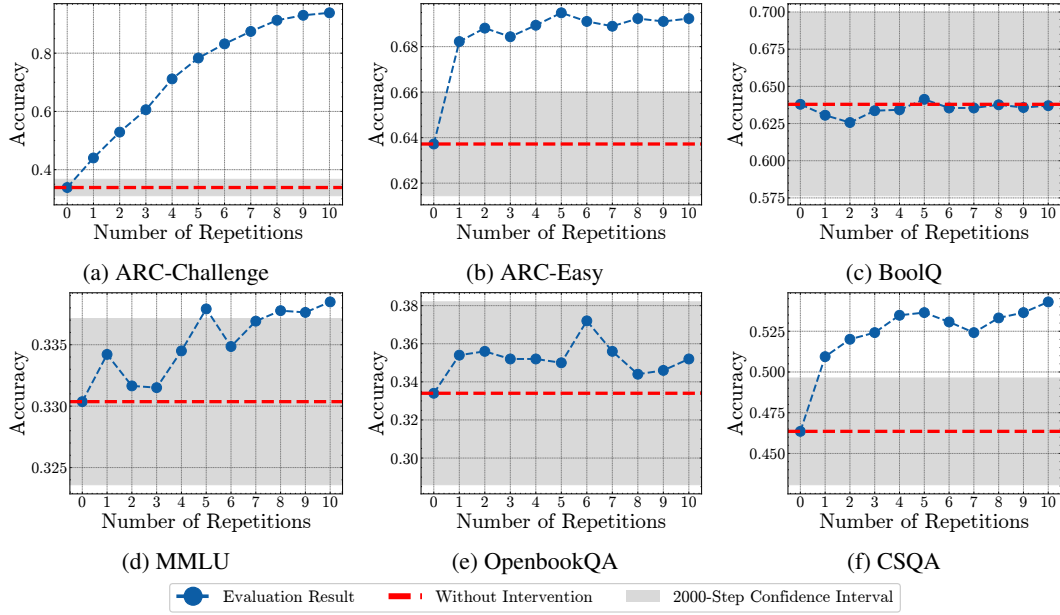


Figure 33: **Modifying the number of repetitions of ARC-Challenge benchmark questions.** The figure depicts the effect of increasing the repetitions of ARC-Challenge benchmark questions on the scores of different benchmarks. The figure also depicts confidence intervals for the respective evaluations. We see that the scores of multiple benchmarks change as we modify the number of ARC-Challenge tokens in the training data.

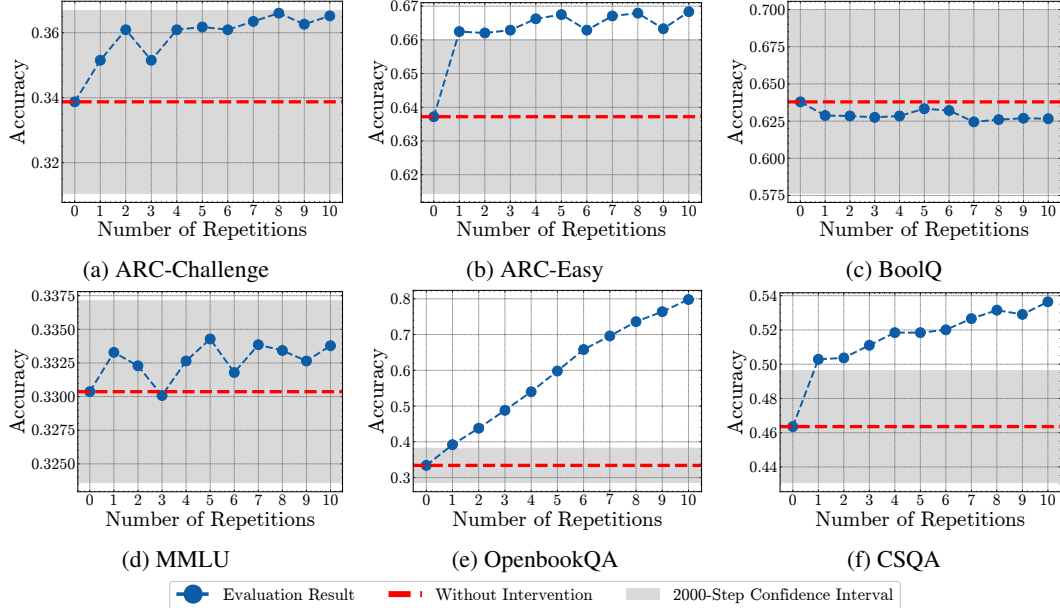


Figure 34: **Modifying the number of repetitions of OpenbookQA benchmark questions.** The figure depicts the effect of increasing the repetitions of OpenBookQA benchmark questions on the scores of different benchmarks. The figure also depicts confidence intervals for the respective evaluations. We see that the scores of multiple benchmarks change as we modify the number of OpenbookQA tokens in the training data.