How MUCH CAN WE FORGET ABOUT DATA CONTAMINATION?

Anonymous authors

Paper under double-blind review

ABSTRACT

The leakage of benchmark data into the training data has emerged as a significant challenge for evaluating the capabilities of large language models (LLMs). In this work, we use experimental evidence and theoretical estimates to challenge the common assumption that small-scale contamination renders benchmark evaluations invalid. First, we experimentally quantify the magnitude of benchmark overfitting based on scaling along three dimensions: The number of model parameters (up to 1.6B), the number of times an example is seen (up to 144), and the number of training tokens (up to 40B). We find that if model and data follow the Chinchilla scaling laws, minor contamination indeed leads to overfitting. At the same time, even 144 times of contamination can be forgotten if the training data is scaled beyond five times Chinchilla, a regime characteristic of many modern LLMs. We then derive a simple theory of example forgetting via cumulative weight decay. It allows us to bound the number of gradient steps required to forget past data for any training run where we know the hyperparameters of AdamW. This indicates that many LLMs, including Llama 3, have forgotten the data seen at the beginning of training. Experimentally, we demonstrate that forgetting occurs faster than what is predicted by our bounds. Taken together, our results suggest that moderate amounts of contamination can be forgotten at the end of realistically scaled training runs.

1 INTRODUCTION

031

003

010 011

012

013

014

015

016

017

018

019

021

025

026

027 028 029

A core principle of machine learning is that a model should not be trained on the test set used for evaluation (Donoho, 2017). For foundation models trained on Internet-scale data, there are increasing concerns that this principle is violated due to the leakage of benchmark evaluation data into the training data (Xu et al., 2024; Oren et al., 2024). Indeed, many LLM developers have found overlap between their training data and the benchmark questions used for evaluation (Brown et al., 2020; Dubey et al., 2024). What is more, research on memorization (Carlini et al., 2019; 2021) shows that text sequences from the training data are sometimes encoded within the model, including machine learning datasets (Grynbaum & Mac, 2023; Liang et al., 2023; Nasr et al., 2023; Bordt et al., 2024).

While the fact that data contamination *can* lead to invalid performance evaluations is now well-040 established (Magar & Schwartz, 2022; Li & Flanigan, 2024; Yang et al., 2023; Jiang et al., 2024), 041 little is known about the precise conditions under which this is the case. The main reason for this is 042 that the training data is usually unknown, and contamination identified via clever research designs 043 that work around this restriction (Golchin & Surdeanu, 2023; Oren et al., 2024; Deng et al., 2024). 044 Because modern foundation models are sometimes trained for over a million gradient steps (Dubey et al., 2024), it is unclear whether a single update on contaminated data at some point during training necessarily impacts downstream evaluations. And indeed, there is quite some evidence that language 046 models need to see samples repeatedly to have any impact on the final model. For example, many 047 papers on memorization have found that it occurs only when a sample is frequently repeated in the 048 training data (Carlini et al., 2022; Biderman et al., 2023; Huang et al., 2024b). The same is true for research on knowledge acquisition, where a fact needs to be paraphrased many times before it is finally remembered by the model (Allen-Zhu & Li, 2023; Cao et al., 2024; Chang et al., 2024). 051

In this work, we study the impact of data contamination in a controlled setting. This means we train
 language models from scratch on datasets where we explicitly insert contaminated examples (Jiang et al., 2024). We begin by quantifying how the overall magnitude of benchmark overfitting (or the

cross-entropy loss of an observed sample) changes as we scale along three critical dimensions: (1) 055 the number of model parameters, (2) the number of training tokens, and (3) the number of repetitions 056 of an example in the training data (Section 4.1). Holding the other two dimensions fixed, we find that 057 the effect of scaling is monotone in each dimension. First, similar to many other works, we find that 058 the tendency of a model to overfit increases in the number of parameters (Goodfellow et al., 2016; Zhang et al., 2017; Carlini et al., 2022). Second, and this is also expected, we find a clear scaling in the number of repetitions, where more frequently repeated observations exhibit stronger overfitting 060 (Carlini et al., 2022; Huang et al., 2024b). More surprisingly, we find that the effect of contamination 061 can vanish as we increase the number of training tokens, up to the point where 12 repetitions of an 062 entire dataset in the training data have no impact on the downstream evaluation on that same dataset. 063

Our investigation reveals that the **natural forgetting** dynamics of gradient descent (Tirumala et al., 064 2022; Jagielski et al., 2023) is the reason why increasing the number of tokens alleviates the impact 065 of contamination. Concretely, we show that training on five times Chinchilla (Hoffmann et al., 2022) 066 of clean data can cause a model to forget even 144 times repeated training examples (Section 4.2). 067 Forgetting the bulk of the impact of a training example can occur rapidly, a point that we demonstrate 068 using OLMo-1B (Groeneveld et al., 2024) (Section 4.3). What is the reason for this rapid forgetting? 069 We show that exposure to novel data is important. Interestingly, models tend to exhibit the strongest overfitting on examples seen repeatedly throughout training, even compared to those seen during the 071 end (Section 4.2). 072

Because running pre-training experiments is expensive, we also ask to what degree forgetting can be 073 explained by the **training dynamics of gradient descent**. We show that the weight decay parameter 074 and learning rate schedule of the AdamW optimizer (Loshchilov & Hutter, 2019) play a key part in 075 forgetting past training examples (Section 5). Concretely, we derive a simple theory of forgetting 076 via cumulative weight decay (Section 5.1) and show that it provides an upper bound on empirical 077 forgetting, which usually occurs faster (Section 5.2). The key point is that this approach allows us to gauge the degree of forgetting present in any training run for which the optimization hyperparameters 079 are known. It even allows us to approximate how the final model parameters balance the gradient 080 updates from different stages of training. Overall, our analysis indicates that many LLMs, including 081 OLMo-7B and Llama 3 405B (Dubey et al., 2024), have forgotten the data seen at the beginning of 082 their training run.

Taken together, our main contribution is to show that the impact of individual examples in the training data depends on the precise characteristics of the setting. There are settings where the effect can be significant; Chinchilla training is an important example (Section 4.1). However, there are equally realistic settings where individual examples dont't matter - including quite likely the data-intensive training runs of many recent LLMs (Gemma Team, 2024).

J88

2 RELATED WORK

090 091 092

094

095

096

097

098

099

100

101

Data Contamination. The GPT-3 paper (Brown et al., 2020) uses an n-gram-based approach to differentiate between "clean" and "dirty" benchmark questions. This approach has since been used in many LLM reports (Chowdhery et al., 2023; Touvron et al., 2023), including Llama 3 (Dubey et al., 2024), where it is estimated that there might be a performance gain of up to 8 and 14 percentage points on PiQA and HellaSwag, respectively. The GPT-4 technical report (Achiam et al., 2023) remarkably concluded that "*contamination overall has very little effect on the reported results*". This has since given rise to a literature that aims to *detect* (Oren et al., 2024), *mitigate* (Li et al., 2024), and *estimate the effect of* (Yang et al., 2023; Bordt et al., 2024) data contamination under various assumptions, but crucially without access to the training data. This literature often challenges the conclusion that contamination overall has little effect in GPT-4 (Xu et al., 2024).

Forgetting. In machine learning, the term *forgetting* is frequently associated with "*catastrophic*" forgetting, where learning new tasks hurt the performance at previously solved tasks (Lopez-Paz & Ranzato, 2017). In the context of LLMs, catastrophic forgetting can occur during fine-tuning (Luo et al., 2023) or continual learning (Huang et al., 2024a). In contrast, this paper studies forgetting as a potential "*natural*" phenomenon of learning (Toneva et al., 2019). Tirumala et al. (2022) study forgetting in language modeling and find, similar to Toneva et al. (2019), that forgetting can be exponentially slow. In contrast, Jagielski et al. (2023) find that models empirically do forget examples

over time. In concurrent work, Pagliardini et al. (2024) propose to add a second momentum term to
 the AdamW optimizer, and show that this slows down the forgetting of past gradients.

Data Attribution. Data attribution methods (Koh & Liang, 2017; Ilyas et al., 2022; Park et al., 2023) aim to identify data points responsible for specific model behaviors. We ask how much a model's benchmark performance is influenced by seeing the example during training, which broadly falls within this field (Grosse et al., 2023; Choe et al., 2024). Importantly, we directly measure the influence of contaminated examples through retraining, avoiding the approximation errors that can occur when using data attribution methods (Koh & Liang, 2017; Ghorbani & Zou, 2019) for large-scale models (Basu et al., 2021; Bae et al., 2022)

117 118 119

120 121

122 123

124

125

126 127

3 BACKGROUND AND METHODS

This Section gives additional details on the research questions and lays out our experimental setup.

Research question

How does the presence of a text in the training data influence the final model's performance *on that same text*?

While it is well known that an individual data point can be influential if the training data and model
are small (Koh & Liang, 2017), we are concerned with the large-data regime where the influence
of any individual example may vanish (Basu et al., 2021; Bae et al., 2022). To further clarify this
setup, Section 3.1 gives a brief overview of recent developments in scaling the training data of LLMs.
Section 3.2 details the used benchmarks and explains how we contaminate the training data. Section
3.3 discusses the problem of near-duplicate benchmark questions.

Models and Training Data. We train language models of up to 1.6B parameters using the architecture
and hyperparameters from the GPT-3 paper (Brown et al., 2020, Table 2.1). For this, we adopt the
llm.c codebase. The training data is the 100BT split of the FineWeb-Edu dataset (Lozhkov et al.,
2024). We also train OLMo-1B (Groeneveld et al., 2024) using the corresponding code and data
(Soldaini et al., 2024).

139 140

141

3.1 WE CONSIDER THE REGIME OF N-TIMES CHINCHILLA

142 According to the Chinchilla scaling law, for every doubling of model size, the number of training tokens should also be doubled. The Chinchilla model itself has 70 billion (B) parameters and was 143 trained on 1.4 trillion (T) tokens; suggesting that the number of training tokens should be roughly 144 20x the number of model parameters (Hoffmann et al., 2022, Table 3). While the Chinchilla paper 145 was highly influential, modern language models are trained on significantly more tokens (Sardana & 146 Frankle, 2024). For example, the OLMo-7B model was trained on 2.46T tokens, 17.5x the amount 147 suggested by Chinchilla (Groeneveld et al., 2024). Similarly, the Llama 3 70B model was reportedly 148 trained on 15T tokens, at over 10x Chinchilla (Dubey et al., 2024; Meta AI, 2024). The same holds 149 for almost all recent LLMs at the 7B parameter scale (Gemma Team, 2024). In this paper, we count 150 the number of tokens a model is trained on as a multiple of its Chinchilla tokens.

151 152

153

3.2 WE EVALUATE ON A MIX OF SEVEN DIFFERENT BENCHMARKS

154 We evaluate the impact of data contamination using a *mix* of seven different benchmarks: ARC-Easy 155 (Clark et al., 2018), Social-I-QA (Sap et al., 2019), WinoGrande (Sakaguchi et al., 2021), PiQA (Bisk 156 et al., 2020), BoolQ (Clark et al., 2019), MMLU (Hendrycks et al., 2021), and HellaSwag (Zellers 157 et al., 2019). This means that every evaluation contains questions from all seven benchmarks. To 158 construct the mixed contamination data, we first concatenate the different benchmarks. We then 159 partition the set of all benchmark questions into subsets ranging from 10,000 to 2,000 questions so that each subset contains all benchmarks in equal weight: HellaSwag: 19.58%, SocialIQA: 8.27%, 160 PiQA: 19.7%, MMLU: 21.82%, BoolQ: 6.48%, ARC-Easy: 5.92%, and WinoGrande: 18.16%. A 161 holdout set of 10,000 benchmark questions is never added to the training data. The other subsets are

HellaSwag: A woman stands holding a violin against

woman stops playing the violin.

HellaSwag: A man is standing outside holding a violin.

sets the violin to his side.

herself. The woman plays the violin. The

He begins to play the violin. he stops and

162 163 164

165 166

167

168

169 170

Figure 1: Language modeling benchmarks frequently contain near-duplicate questions. We perform extensive filtering for duplicates using fuzzy string matching. The figure depicts a near-duplicate from HellaSwag and a cross-benchmark duplicate from ARC-Easy/MMLU.

ARC-Easy: Question: A student is playing with a small

MMLU: Question: A wave transfers

Answer: energy

toy boat [...] The boat moves toward the shore

because the waves transfer Answer: energy.

172 173

180

182

171

added to the training data, repeated either 4, 12, 36, or 144 times.¹ We consider *exact* contamination, that is we contaminate the training data with the same texts that the model is later evaluated on. We insert benchmark questions *individually* and at *random* positions into the training data. Models are evaluated zero-shot via the likelihood assigned to different sentence completions (Gao, 2021).
For more discussion and details about how contamination is performed, see Supplement A.1 and Supplement A.2.

181 3.3 WE FILTER NEAR-DUPLICATE BENCHMARK QUESTIONS

Our method requires that there are no side effects from contaminating the training data with one 183 question on the evaluation of another question. However, upon closer inspection, it turns out that *all* 184 the commonly used benchmarks from the literature contain questions that are either near-duplicates 185 or where the context of one question contains the answer to another question (for example, because the same text document was used to create multiple questions). This is illustrated in Figure 1, which 187 depicts two near-duplicate questions on HellaSwag and questions from ARC-Easy and MMLU that 188 are cross-benchmark duplicates. To address this problem, we perform extensive filtering, removing 189 duplicate questions where the length-normalized Levenshtein distance falls below a certain threshold 190 (Levenshtein, 1966; Navarro, 2001). This is documented in Supplement A.3, where we also describe 191 an experiment to verify that our method is not invalidated by near-duplicate questions.

192 193

194 195

196

197

199

201

4 EXPERIMENTAL RESULTS

We now present our main experimental results. We being in Section 4.1 by discussing the scaling in model parameters, training tokens, and repetitions in the training data. The following Section 4.2 discusses various experiments on forgetting. The first two sections rely on training small GPT-3 models. Section 4.3 complements this with an analysis of OLMo-1B (Groeneveld et al., 2024).

4.1 CONTAMINATION SCALES WITH MODEL, DATA, AND REPETITIONS

We conduct three different experiments to understand how the effect of data contamination scales with the number of model parameters, training tokens, and the number of times a contaminated example is seen. First, we train increasingly large models on the same dataset of 7B tokens. Second, we train 124M parameter models on increasingly many tokens. Third, we train increasingly large models according to the Chinchilla scaling laws (Hoffmann et al., 2022), meaning that the number of training tokens scales linearly with the model parameters. In all experiments, we contaminate the training data *uniformly at random* with benchmark questions.

Figure 2 depicts the results of all three experiments. Because we are interested in the performance *difference* between the holdout data and the contaminated examples, Figure 2 depicts the *accuracy gap* between the holdout and contaminated examples in percentage points. In Figure 2a, we see that the accuracy gap due to contamination is *increasing in the number of model parameters*. For a 124M

¹In preliminary experiments, we found that these numbers pleasantly cover the range from statistically significant contamination to complete overfitting. We also considered repeating observations a single time, as in (Jiang et al., 2024). However, we found this often leads to accuracy differences of about one or two percentage points, just within the margin of our confidence intervals, which is undesirable.



Figure 2: Benchmark overfitting due to contamination. (a) We train different models on 7B tokens.
(b) We train 124M parameter models on increasingly many tokens. (c) We train models according to the Chinchilla scaling laws. The figure depicts the accuracy difference in percentage points between the holdout (normalized to zero) and the contaminated examples. The results are across a mix of seven different benchmarks, as outlined in Section 3.2. Different colors indicate different levels of contamination. Mean and bootstrapped 90% confidence intervals.

228

229

230

231

232

236 parameter model trained on 7B tokens, the overfitting due to 4 times contamination is 5 percentage 237 points. For a 1.6B parameter model train on the same dataset, it is 20. Next, Figure 2b shows that the 238 accuracy gap is *decreasing in the number of training tokens*. For a 124M parameter model trained at 239 2x Chinchilla, the accuracy gap due to 12 times contamination is 18 percentage points. For a 124M parameter model trained at 15x Chinchilla, the same accuracy gap is within the confidence interval of 240 the holdout. From Figure 2, we also see that the accuracy gap is *increasing in the number of times an* 241 example is repeated. For a 350M parameter model trained on 7B tokens, the accuracy gap is 11, 25, 242 44, and 51 percentage points for 4, 12, 32, and 144 times repeated contamination, respectively. 243

244 Because the accuracy gap *increases* in the 245 number of model parameters and decreases in the number of tokens, the interesting 246 question is how it behaves if model param-247 eters and tokens are scaled jointly. A nat-248 ural starting point is to double the number 249 of training tokens for every doubling of 250 model parameters, as specified by the Chin-251 chilla scaling laws (Hoffmann et al., 2022). 252 Figure 2c depicts the accuracy gap due to

Table 1: Accuracy of the Chinchilla models.

Model	Holdout	4x	12x	32x	144x
124M	42.22	48.14	56.92	80.70	96.45
350M	44.72	55.69	69.90	89.20	95.50
774M	49.16	64.76	81.30	92.95	96.05
1.6B	52.06	67.61	82.32	91.85	95.40

contamination as we train increasingly large Chinchilla-optimal models. While there is no clear
monotone pattern, we see that moderate amounts of contamination can lead to significant overfitting.
For the 774M parameter model, 4 times repeated contamination leads to an accuracy gap of 15
percentage points, suggesting that *under Chinchilla training, a single time of contamination can lead to overfitting of as much as 3 percentage points.*²

258 259

260

4.2 CONTAMINATION CAN BE COMPLETELY FORGOTTEN

In the previous Section 4.1, we saw that the accuracy gap due to contamination decreases in the number of tokens up to the point where even 12 repetitions of a benchmark question in the training data can become insignificant. In this Section, we identify the natural forgetting dynamic of neural network training as the reason for this effect. We discuss how quickly forgetting occurs, whether examples are completely forgotten, and what kind of repetition makes a model remember.

To study the effect of forgetting, we train a 124M parameter model for 15 epochs. Instead of contaminating uniformly over the course of training like in the previous Section 4.1, we perform

²To contaminate the training data of a 774M model a single time with 10,000 benchmark questions, we need

²To contaminate the training data of a 774M model a single time with 10,000 benchmark questions, we need to insert \sim 0.5 million tokens into 15.5 billion training tokens, about 0.003% of the training data.



Figure 3: **The natural forgetting dynamic of neural network training. (a)** The development of the cross-entropy loss difference between contaminated and holdout benchmark questions over the course of training. Contamination occurs between the first and second Chinchilla (1 and 2 on the x-axis). (b) Accuracy gaps after training for 3 Chinchilla. (c) Accuracy gaps after training for 7 Chinchilla. (d) Same as (a). (e)+(f) The accuracy gap depends on the average position of an example in the training data. Mean and bootstrapped 90% confidence intervals.

the contamination between the first and second Chinchilla.³ Figure 3a depicts the development of 300 the *difference* in cross-entropy loss between contaminated and clean benchmark questions over the 301 course of training. We see a strong peak after 2 Chinchilla, which is expected and shows the effect 302 of contamination. What is interesting to us is the rate at which the cross-entropy loss difference 303 decays as we continue training. After training for 1 additional Chinchilla (2.5B tokens for the 304 124M parameter model), it has already decayed significantly. However, the difference is still visible 305 in Figure 3a. Figure 3b depicts the corresponding accuracy gaps at this point, and we see that all 306 contamination levels still lead to overfitting. As we continue training, the cross-entropy loss difference 307 between contaminated and holdout questions further narrows. From Figure 3c, which depicts the 308 accuracy gaps after forgetting for a total of 5 Chinchilla, we see that the effect of contamination is 309 eventually *completely forgotten* in the sense that there is no longer any accuracy difference between contamination and holdout benchmark questions. 310

311 The result that contamination can be completely forgotten is in contrast to some previous work on 312 forgetting which have found that forgetting approaches a stable baseline Tirumala et al. (2022, Figure 313 10), or that certain examples are never forgotten (Toneva et al., 2019). To understand this difference, 314 observe that many previous works on forgetting have not trained on a continuous stream of data. 315 Instead, they have trained on the same training set for multiple epochs. Consequently, we modify our forgetting experiment to repeatedly train on the same 100M tokens after the second epoch. The result 316 of this experiment is depicted in Figure 3d and should be compared to Figure 3a. Interestingly, this 317 simple modification causes the effect of forgetting to stabilize at a level strictly larger than zero. We 318 conclude that *exposure to novel data is important for forgetting*, an observation similar to Jagielski 319 et al. (2023). 320

 ³Note that the model is already fairly trained after the first Chinchilla, meaning that the contamination is not very early during training. This is important because there is evidence that observations are more quickly forgotten if the model has not yet learned representations (Jagielski et al., 2023; Cao et al., 2024; Huang et al., 2024b). This is *not* the setting we are studying here.



Figure 4: **Contamination and forgetting in OLMo-1B.** We contaminate the OLMo-1B checkpoint at gradient step 369,000 four times with different benchmarks. This causes an average accuracy increase of 15 percentage points. We then continue pre-training for 1% of the remaining training time, leading to a reduction of 96% of the accuracy increase due to contamination. In this figure, different colors simply correspond to different benchmarks, and the grey line depicts the clean accuracy without contamination. Mean and bootstrapped 90% confidence intervals.

To further understand the impact of forgetting, we now ask whether examples seen late during training 341 influence model behavior more strongly than examples seen early during training. To study this 342 question, we average all the different *uniform* contamination levels from the models in the previous 343 Section 4.1 (to gain statistical power) and consider the amount of overfitting depending on whether a 344 question is seen, on average, in the beginning, middle, or end of training. The result of this experiment 345 is depicted in Figure 3e and Figure 3f. As expected under forgetting, we see that benchmark questions 346 seen early during training exhibit the smallest amount of overfitting. Interestingly and somewhat 347 unexpectedly, questions that are neither clustered towards the beginning nor the end but as uniformly 348 distributed throughout training as possible exhibit the strongest overfitting, suggesting that this spaced form of repetition helps the model remember (the middle peak is the most pronounced both in Figure 349 3e and Figure 3f). 350

351 352

353

4.3 CONTAMINATION AND RAPID FORGETTING IN OLMO-1B

354 In the previous sections, we trained small GPT-3 models from scratch. In this Section, we complement 355 this analysis by pre-training from an intermediate OLMo-1B checkpoint (Groeneveld et al., 2024). Similar to the analysis in Section 4.2, we insert the benchmark data at a specific point into the training 356 data and then measure the subsequent forgetting. Unlike in the previous Section, we now insert the 357 entire benchmark data – we already have a "clean" baseline from the original OLMo-1B training 358 run. We insert each benchmark question four times and contaminate with four different benchmarks: 359 HellaSwag (Zellers et al., 2019), WinoGrande Sakaguchi et al., 2021, ARC-Easy (Clark et al., 2018), 360 and PiQA (Bisk et al., 2020). 361

Figure 4 depicts the result of the experiment. The effect of contamination is visible from the five
leftmost points of every plot. The leftmost point corresponds to the uncontaminated model, and the
next four points each depict the effect of one time contamination. Again, we see that *the immediate effect of contamination is significant*, leading to an average accuracy increase of 15 percentage points
across the different benchmarks. At the same time, we also see that *the effect of contamination decays considerably as we continue training*. To contextualize this result, note that Figure 4 depicts less than
2000 gradient steps. The pre-training stage of OLMo-1B model consists of 739,328 gradient steps.
This means that Figure 4 depicts less than 1% of the total forgetting until pre-training is done.

- 369 370
- 371 372

5 WHAT IS THE ROLE OF WEIGHT DECAY IN FORGETTING?

In the previous Section 4, we have seen that forgetting is an important empirical property of LLM
training. In this section, we show that the weight decay parameter and learning rate schedule of the
AdamW optimizer play a key part in forgetting past training examples. This offers a novel perspective
on the interplay between these two parameters, usually seen in terms of generalization and training
stability (Van Laarhoven, 2017; Zhang et al., 2019; Lewkowycz & Gur-Ari, 2020; Andriushchenko et al., 2023).



Figure 5: Theoretical estimates of forgetting and approximate model weight composition at the end of training. Top Row: The cumulative weight decay $w_{t_1}^{t_2}$ as defined in equation (3) for different training runs. The figures depict the decay of the gradient updates for every decile of the training run, indicated in different colors. Bottom Row: The approximate composition of the final model weights in terms of the gradient updates from different deciles of the training run. The deciles are indicated in colors depicted in the legend below the plot.

5.1 WEIGHT DECAY AS A MECHANISM FOR FORGETTING TRAINING EXAMPLES

Consider the parameter update of AdamW at gradient step $t \ge 1$. It consists of two decoupled updates (Paszke et al., 2019; PyTorch Contributors, 2024): A weight decay update given by

$$\theta_t = \theta_{t-1} - \gamma \lambda_t \theta_{t-1},\tag{1}$$

and a gradient update given by

392

393

394

397 398

399

400

401

402 403 404

405 406

407

412 413 414

$$\theta_t = \hat{\theta}_t - \lambda_t \, \hat{m}_t / (\sqrt{\hat{v}_t + \epsilon}). \tag{2}$$

Here, θ_t are the model parameters, λ_t is the learning rate, γ is the weight decay parameter, and \hat{m}_t and \hat{v}_t are first- and second-order moment estimates of the gradient. Denoting the model weights at initialization by θ_0 , and the adaptive gradient by $\hat{g}_t = \hat{m}_t / (\sqrt{\hat{v}_t + \epsilon})$, we can iterate (1) and (2) to obtain

$$\theta_T = w_0^T \theta_0 - \sum_{t=1}^T w_t^T \lambda_t \hat{g}_t \quad \text{where} \quad w_{t_1}^{t_2} = \prod_{i=t_1+1}^{t_2} (1 - \lambda_i \gamma).$$
(3)

Here, the weights $w_{t_1}^{t_2}$ account for the *cumulative weight decay* between gradient step t_1 and t_2 . Intuitively, the model weights after t gradient steps are a weighted average of the initial model weights and all the adaptive gradient updates up to time step t. This is not specific to the AdamW optimizer and applies to every optimizer with weight decay. Analyzing equation (3) reveals a critical factor influencing forgetting: The exponential decay of the $w_{t_1}^{t_2}$ with respect to increasing gap $t_2 - t_1$ in the optimization steps. In other words, the longer an update occurs in the past (i.e., the larger $t_2 - t_1$), the more the contribution of the update \hat{g}_{t_1} is being scaled down due to the exponential decay of weights $w_{t_1}^{t_2}$. We can describe the evolution of these weights as the function of the time $T = t_2 - t_1$.

Proposition 1. (*The Decay of Past Gradients*) *The number of optimization steps* $T = t_2 - t_1$ *that are required to make the contribution of a model update at time* t_1 *small, that is* $w_{t_1}^{t_2} \le \epsilon$ *for some small* $\epsilon \in \mathbb{R}^+$, *scales as* $T \gtrsim \frac{\log(1/\epsilon)}{\gamma \lambda_{avg}}$, *where* $\lambda_{avg} = \frac{1}{T} \sum_{t=t_1+1}^{t_2} \lambda_t$ *is the average learning rate of the optimizer between* t_1 *and* t_2 .

427 428 429 429 430 We present a proof in the Supplement A.5. If $w_{t_2}^{t_1} \leq \epsilon$ for some sufficiently small ϵ , the term $w_{t_2}^{t_1} \lambda_{t_1} \hat{g}_{t_1}$ vanishes from the sum in (3). Intuitively, this is the same as saying that the gradient update at time step t_1 has been forgotten.

431 We now analyze the weight-decay mechanism of forgetting in different LLM training runs. Of course, the decay of past gradients described in Proposition 1 is only one effect that contributes to



Figure 6: The theoretical estimates provide an upper bound on empirical forgetting, which can occur much faster. (a) Empirical forgetting for three different weight decay parameters. (b) The corresponding cumulative weight decay. (c) Empirical forgetting and theoretical estimate for the default weight decay of 0.1. The y-axis of the cumulative weight decay $w_{t_1}^{t_2}$ is calibrated to start at the peak of the empirically observed overfitting. (d) Empirical forgetting and theoretical estimate for OLMo-1B (HellaSwag).

441

442

443

444

445

448 forgetting. Indeed, there is also the potential for different gradient updates in the sum (3) to cancel 449 each other. However, because a term decayed to zero is definitively forgotten, we can consider the cumulative weight decay as an upper bound on forgetting. Figure 5 depicts the evolution of 450 the forgetting term $w_{t}^{t_2}$ for the 124M parameter model from Section 4.2, OLMo 7B, and Llama 451 3 405B (where we assume the model trained with a weight decay of 0.1). For the training data 452 at each decline of the training run, Figure 5 depicts *forgetting curves* that indicate how much the 453 corresponding gradients decay as we continue training. For the 124M model depicted in Figure 5a, 454 even the initialization is not completely decayed towards the end (the blue curve is still strictly larger 455 than zero). In contrast, for OLMo-7B, the gradients of the first 40% of the training data decay to 456 zero until the end of training, meaning that this data is forgotten (Figure 5b). Llama 3 405B also 457 experiences significant decay of the early gradients (Figure 5c). 458

The interplay between the weight decay parameter and learning rate schedule of AdamW 459 creates an interesting dynamic. Lower learning rates towards the end of training (1) slow down the 460 forgetting of past training examples and (2) decrease the impact of later training examples on the 461 final model weights. To better understand this dynamic, we plot the sum of the terms $\lambda_{t_1} w_{t_2}^{t_2}$ for 462 each decile of the training run, normalized by the same sum over the entire training run. This is 463 depicted in the bottom row of Figure 5 and can be thought of as a simple approximation of how the 464 final model weights are composed by the gradients of different training deciles. Interestingly, this 465 approximation suggests that the Llama 3 405B training run, where supposedly a lot of expertise has 466 gone into the choice of the hyperparameters, results in a model where the approximate influence of different training steps is symmetrically distributed around the middle of training (Figure 5c). In 467 contrast, the OLMo-7B training could seemingly benefit from further decaying the learning rate to 468 give more weight to early versus late gradients. 469

470 While our analysis in this Section considers the mechanistic effect of individually decaying gradient 471 updates, it does not model any interactions between different gradient updates. For example, if the 472 weight updates at a later time step t_2 were aligned with past updates at t_1 , then the model might not forget the information even if the effect of past updates vanishes from the sum. However, such complex 473 interactions are avoided if the model updates from contaminated samples are orthogonal. Formalizing 474 this observation leads to a more rigorous version of Proposition 1, presented in Supplement A.5. The 475 argument is that under suitable gradient orthogonality conditions, the decay of past gradients can 476 guarantee forgetting. 477

478

479 480

5.2 PRACTICAL FORGETTING OCCURS FASTER THAN WHAT THE THEORY PREDICTS

In the previous Section 5.1, we have derived a simple theory of forgetting via cumulative weight decay. We now investigate how the theoretical estimates relate to the empirically observed forgetting.

The main parameter that controls the theoretical forgetting curves is the weight decay parameter.
Therefore, we ask how forgetting changes empirically when we change the weight decay. Figure
6a depicts the result of repeating the forgetting experiment from Section 4.2 with three different
choices for the weight decay parameter. From Figure 6a, we see that the weight decay parameter

controls the impact of contamination at all time steps, where a larger weight decay parameter leads to
 more forgetting and a smaller weight decay parameter to less forgetting. This is consistent with the
 theoretical predictions depicted in Figure 6b, meaning there is a *qualitative alignment between the empirical results and our theoretical predictions*.

490 To better understand the quantitative relation between empirical forgetting and the theoretical esti-491 mates, we ask how the empirically forgotten fraction (of cross-entropy loss or accuracy) relates to the 492 cumulative weight decay. Figure 6c depicts the empirical decay and corresponding theoretical predic-493 tion for the model from Section 4.2. We see that the theoretical estimate is somewhat pessimistic and 494 that forgetting occurs faster than what is predicted by the theory. Figure 6d is similar to Figure 6c, 495 except that we consider the OLMo-1B forgetting experiment from Section 4.3. Figure 6d depicts a 496 case where *forgetting occurs much faster* than what is predicted by the theory. Interestingly, we also see that the empirical rate of forgetting, at least in this experiment, is not smaller for the larger model. 497

498 499

500

6 DISCUSSION

This work presents different experiments on data contamination and forgetting. We have seen that the impact of contamination can vanish as the size of the training data increases – an aspect that has largely been overlooked in the literature (Yang et al., 2023; Oren et al., 2024; Jiang et al., 2024). We have also shown that the hyperparameters of AdamW play an important part in forgetting – an insight that might inform the parametrization of future training runs. Of course, it would be interesting to study contamination and forgetting on a larger scale. Nevertheless, we have shown that the effects of contamination can vanish *even though* the number of steps we could train the models for was limited.

We have studied data contamination with a focus on the leakage of benchmark questions into the training data. This means that our work might be more informative about the topic than other works that study contamination in different contexts. At the same time, one has to be careful when extrapolating our results, especially to a privacy setup (Carlini et al., 2019; Jagielski et al., 2023). This is because empirical forgetting might behave differently for random strings or otherwise uniquely identifiable information (Carlini et al., 2021).

The "*contamination problem*" as studied in this paper has interesting connections to many areas of machine learning, including privacy (Graves et al., 2021; Jagielski et al., 2023), data attribution (Kirchenbauer et al., 2024), and generalization (Bousquet & Elisseeff, 2002; Hardt et al., 2016; Mania et al., 2019). The connection to data attribution is especially relevant in light of our results. This is because we demonstrate cases where the presence or absence of a datapoint in the training data is irrelevant for the model behavior *on that same datapoint*, meaning it does not make sense to attribute model behavior to individual datapoints in this regime.

521 522

523

524

525

526

527

528

529

7 REPRODUCIBILITY STATEMENT

- The results in this paper were obtained using the OLMo codebase, available at https://github.com/allenai/OLMo, and the llm.c codebase, available at https://github.com/karpathy/llm.c. Our code is fully reproducible, including the random positions at which benchmark questions were inserted into the training data. We trained on the 100BT split of the FineWeb-Edu dataset, available at https://huggingface.co/datasets/HuggingFaceFW/fineweb-edu. Anonymized code for this paper is available at https://github.com/iclr10261/code.
- 530 531 532

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4 technical report.
 OpenAI, 2023.
- Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.1, knowledge storage and extraction. *arXiv preprint arXiv:2309.14316*, 2023.
- 539 Maksym Andriushchenko, Francesco D'Angelo, Aditya Varre, and Nicolas Flammarion. Why do we need weight decay in modern deep learning? *arXiv preprint arXiv:2310.04415*, 2023.

540 Juhan Bae, Nathan Ng, Alston Lo, Marzyeh Ghassemi, and Roger B Grosse. If influence functions 541 are the answer, then what is the question? Advances in Neural Information Processing Systems, 542 35:17953-17967, 2022. 543 S Basu, P Pope, and S Feizi. Influence functions in deep learning are fragile. In International 544 Conference on Learning Representations (ICLR), 2021. 546 Stella Biderman, USVSN PRASHANTH, Lintang Sutawika, Hailey Schoelkopf, Quentin Anthony, 547 Shivanshu Purohit, and Edward Raff. Emergent and predictable memorization in large language 548 models. In NeurIPS, 2023. 549 Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical 550 commonsense in natural language. In AAAI, 2020. 551 552 Sebastian Bordt, Harsha Nori, Vanessa Rodrigues, Besmira Nushi, and Rich Caruana. Elephants 553 never forget: Memorization and learning of tabular data in large language models. In COLM, 2024. 554 Olivier Bousquet and André Elisseeff. Stability and generalization. The Journal of Machine Learning 555 Research, 2002. 556 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, 558 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel 559 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, 560 Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, 561 and Dario Amodei. Language models are few-shot learners. In NeurIPS, 2020. 562 563 Boxi Cao, Qiaoyu Tang, Hongyu Lin, Shanshan Jiang, Bin Dong, Xianpei Han, Jiawei Chen, Tianshu 564 Wang, and Le Sun. Retentive or forgetful? diving into the knowledge memorizing mechanism of 565 language models. In International Conference on Computational Linguistics, Language Resources 566 and Evaluation, 2024. 567 Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. The secret sharer: 568 Evaluating and testing unintended memorization in neural networks. In 28th USENIX security 569 symposium (USENIX security 19), 2019. 570 571 Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine 572 Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data 573 from large language models. In 30th USENIX Security Symposium (USENIX Security 21), 2021. 574 Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and 575 Chiyuan Zhang. Quantifying memorization across neural language models. arXiv preprint 576 arXiv:2202.07646, 2022. 577 578 Hoyeon Chang, Jinho Park, Seonghyeon Ye, Sohee Yang, Youngkyung Seo, Du-Seong Chang, and 579 Minjoon Seo. How do large language models acquire factual knowledge during pretraining? arXiv preprint arXiv:2406.11813, 2024. 580 581 Sang Keun Choe, Hwijeen Ahn, Juhan Bae, Kewen Zhao, Minsoo Kang, Youngseog Chung, Adithya 582 Pratapa, Willie Neiswanger, Emma Strubell, Teruko Mitamura, et al. What is your data worth to 583 gpt? Ilm-scale data valuation with influence functions. arXiv preprint arXiv:2405.13954, 2024. 584 585 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: 586 Scaling language modeling with pathways. Journal of Machine Learning Research, 2023. 587 588 Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina 589 Toutanova. BoolQ: Exploring the surprising difficulty of natural yes/no questions. In Jill Burstein, 590 Christy Doran, and Thamar Solorio (eds.), NAACL, 2019. 591 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and 592 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. arXiv:1803.05457v1, 2018.

594 595 596	Chunyuan Deng, Yilun Zhao, Xiangru Tang, Mark Gerstein, and Arman Cohan. Investigating data contamination in modern benchmarks for large language models. In ACL, 2024.		
597	David Donoho. 50 Years of Data Science. Journal of Computational and Graphical Statistics, 2017.		
598	Abhimanyu Dubey Abhinay Jaubri Abhinay Pandey Abhichek Kadian Abmad Al Dahle Ajecha		
599	Letman Akhil Mathur Alan Schelten Amy Yang Angela Fan et al. The llama 3 herd of models		
600	arXiv preprint arXiv:2407.21783. 2024.		
601			
602 603	Leo Gao. Multiple choice normalization in lm evaluation. <i>Blog Post</i> , 2021. URL https://blog.		
604	eredener.ar, marcipre energe normarraderen,.		
605	Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles		
606	Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas		
607	Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron,		
608	Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework		
609	In ovaluation harmons		
610	III-evaluacion-harness.		
611	Gemma Team. Gemma 2: Improving open language models at a practical size. arXiv preprint		
612	arXiv:2408.00118, 2024.		
613			
614	Amirata Gnorbani and James Zou. Data shapley: Equitable valuation of data for machine learning.		
615	III <i>ICML</i> , 2019.		
616	Shahriar Golchin and Mihai Surdeanu. Time travel in llms: Tracing data contamination in large		
617	language models. arXiv preprint arXiv:2308.08493, 2023.		
618			
619	Ian Goodfellow, Yoshua Bengio, and Aaron Courville. <i>Deep Learning</i> . MIT Press, 2016. http:		
620	//www.deeplearningbook.org.		
621	Laura Graves, Vineel Nagisetty, and Vijay Ganesh. Amnesiac machine learning. In AAI Conference		
622	on Artificial Intelligence, 2021.		
623			
624 625	Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Iha Hamish Ivison Ian Magnusson Yizhong Wang et al. Olmo: Accelerat-		
626	ing the science of language models. arXiv preprint arXiv:2402.00838, 2024.		
627			
628	Roger Grosse, Juhan Bae, Cem Anil, Nelson Elhage, Alex Tamkin, Amirhossein Tajdini, Benoit		
629	Steiner, Dustin Li, Esin Durmus, Ethan Perez, et al. Studying large language model generaliz		
630	with influence functions. arXiv preprint arXiv:2308.03296, 2023.		
631	Michael M Grynbaum and Ryan Mac. The Times Sues OpenAI and Microsoft Over A.I. Use of		
632	Copyrighted Work, 2023. URL https://www.nytimes.com/2023/12/27/business/		
633	media/new-york-times-open-ai-microsoft-lawsuit.html.		
634			
635	Moritz Hardt, Ben Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic		
636	gradient descent. In ICLM, 2016.		
637	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob		
638	Steinhardt. Measuring massive multitask language understanding. In <i>ICLR</i> , 2021.		
639			
640	Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza		
641	Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al.		
642	raining compute-optimal large language models. arXiv preprint arXiv:2203.15556, 2022.		
643	Jianheng Huang, Levang Cui, Ante Wang, Chengyi Yang, Xinting Liao, Linfeng Song, Junfeng Yao		
644	and Jinsong Su. Mitigating catastrophic forgetting in large language models with self-synthesized		
645	rehearsal. In ACL Annual Meeting, 2024a.		
646			
647	Jing Huang, Diyi Yang, and Christopher Potts. Demystifying verbatim memorization in large language models. <i>arXiv preprint arXiv:2407.17817</i> , 2024b.		

648 649 650	Andrew Ilyas, Sung Min Park, Logan Engstrom, Guillaume Leclerc, and Aleksander Madry. Data- models: Predicting predictions from training data. <i>International Conference on Machine Learning</i> , 2022
651	
652	Matthew Jagielski, Om Thakkar, Florian Tramer, Daphne Ippolito, Katherine Lee, Nicholas Carlini,
653	Eric Wallace, Shuang Song, Abhradeep Guha Thakurta, Nicolas Papernot, and Chiyuan Zhang.
654	Measuring forgetting of memorized training examples. In ICLR, 2023.
655	Minhao Jiang Ken Zivu Liu Ming Zhong Rylan Schaeffer Siru Ouvang Jiawei Han and Sanmi
656	Kovejo. Investigating data contamination for pre-training language models. <i>arXiv preprint</i>
657	arXiv:2401.06059, 2024.
658	
659	John Kirchenbauer, Garrett Honke, Gowthami Somepalli, Jonas Geiping, Daphne Ippolito, Katherine
660 661	preprint arXiv:2405.06331, 2024.
662 663	Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In <i>ICML</i> , 2017.
664 665	Vladimir Levenshtein. Binary codes capable of correcting deletions, insertions, and reversals. <i>Proceedings of the Soviet physics doklady</i> , 1966.
666	
668	Aitor Lewkowycz and Guy Gur-Ari. On the training dynamics of deep networks with <i>l</i> _2 regularization. <i>NeurIPS</i> , 2020.
669	Changman Li and Jeffrey Flanigan. Task contamination: Language models may not be few shot
670	anymore In Proceedings of the AAAI Conference on Artificial Intelligence 2024
671	
672	Yucheng Li, Frank Guerin, and Chenghua Lin. Latesteval: Addressing data contamination in language
674	model evaluation through dynamic and time-sensitive test construction. In AAAI Conference on
675	Artificial Intelligence, 2024.
676	Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga,
677	Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan,
678	Bobby Yan, Ce Zhang, Christian Alexander Cosgrove, Christopher D Manning, Christopher Re,
679	Diana Acosta-Navas, Drew Arad Hudson, Eric Zelikman, Esin Durmus, Faisal Ladhak, Frieda
680	Rong, Hongyu Ren, Huaxiu Yao, Jue WANG, Keshav Santhanam, Laurel Orr, Lucia Zheng,
681	Peter Henderson, Oian Huang, Ryan Andrew Chi, Sang Michael Xie, Shibani Santurkar, Surva
682	Ganguli, Tatsunori Hashimoto, Thomas Icard, Tianyi Zhang, Vishray Chaudhary, William Wang,
683	Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. Holistic evaluation of language models.
684	Transactions on Machine Learning Research, 2023.
685	David Long Day and Maro' Aurolia Dangata Cradient anigodia manager for continue la sector
686	NourIPS 2017
687	1100111 0, 2017.
688	Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In ICLR, 2019.
689	Anton Lozhkov, Louhna Ran Allal, Laandro von Warra, and Thomas Walf, Finawah edg. 2024. UDI
690	https://huggingface.co/datasets/HuggingFaceFW/fineweb-edu
691	
092	Yun Luo, Zhen Yang, Fandong Meng, Yafu Li, Jie Zhou, and Yue Zhang. An empirical study of
604	catastrophic forgetting in large language models during continual fine-tuning. <i>arXiv preprint</i>
605	arxiv:2508.08/4/, 2023.
696	Inbal Magar and Roy Schwartz. Data contamination: From memorization to exploitation. arXiv
697	preprint arXiv:2203.08242, 2022.
698	Horia Mania John Miller Ludwig Schmidt Maritz Hardt and Daniamin Dacht Madal similarity
699	mitigates test set overuse NeurIPS 2010
700	mingato tot overuse. Ivenin 5, 2017.
701	Meta AI. Introducing Meta Llama 3: The most capable openly available LLM to date. <i>Blog Post</i> , 2024. URL https://ai.meta.com/blog/meta-llama-3/.

702 703 704 705 706	Milad Nasr, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A Feder Cooper, Daphne Ippolito, Christopher A Choquette-Choo, Eric Wallace, Florian Tramèr, and Katherine Lee. Scalable extraction of training data from (production) language models. <i>arXiv preprint arXiv:2311.17035</i> , 2023.
707 708	Gonzalo Navarro. A guided tour to approximate string matching. ACM computing surveys (CSUR), 2001.
709 710 711	Yonatan Oren, Nicole Meister, Niladri Chatterji, Faisal Ladhak, and Tatsunori B Hashimoto. Proving test set contamination in black box language models. In <i>ICLR</i> , 2024.
712 713	Matteo Pagliardini, Pierre Ablin, and David Grangier. The ademamix optimizer: Better, faster, older. <i>arXiv preprint arXiv:2409.03137</i> , 2024.
714 715 716 717	Sung Min Park, Kristian Georgiev, Andrew Ilyas, Guillaume Leclerc, and Aleksander Madry. Trak: Attributing model behavior at scale. In <i>International Conference on Machine Learning</i> , pp. 27074–27113. PMLR, 2023.
718 719 720	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In <i>NeurIPS</i> , 2019.
721 722 723	PyTorch Contributors. AdamW. Pytorch Documentation, 2024. URL https://pytorch.org/ docs/2.4/generated/torch.optim.AdamW.html.
724 725	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. <i>Communications of the ACM</i> , 2021.
726 727 728	Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. Social IQa: Common- sense reasoning about social interactions. In <i>EMNLP-IJCNLP</i> . ACL, 2019.
729 730 731	Nikhil Sardana and Jonathan Frankle. Beyond chinchilla-optimal: Accounting for inference in language model scaling laws. In <i>ICML</i> , 2024.
732 733 734 735 736 737 738 739	Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur, Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, Valentin Hofmann, Ananya Harsh Jha, Sachin Kumar, Li Lucy, Xinxi Lyu, Nathan Lambert, Ian Magnusson, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew E. Peters, Abhilasha Ravichander, Kyle Richardson, Zejiang Shen, Emma Strubell, Nishant Subramani, Oyvind Tafjord, Pete Walsh, Luke Zettlemoyer, Noah A. Smith, Hannaneh Hajishirzi, Iz Beltagy, Dirk Groeneveld, Jesse Dodge, and Kyle Lo. Dolma: An Open Corpus of Three Trillion Tokens for Language Model Pretraining Research. <i>arXiv preprint</i> , 2024. URL https://arxiv.org/abs/2402.00159.
740 741	Kushal Tirumala, Aram Markosyan, Luke Zettlemoyer, and Armen Aghajanyan. Memorization without overfitting: Analyzing the training dynamics of large language models. In <i>NeurIPS</i> , 2022.
742 743 744 745	Mariya Toneva, Alessandro Sordoni, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and Geoffrey J Gordon. An empirical study of example forgetting during deep neural network learning. In <i>ICLR</i> , 2019.
746 747 748	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models, 2023.
749 750 751	Twan Van Laarhoven. L2 regularization versus batch and weight normalization. <i>arXiv preprint arXiv:1706.05350</i> , 2017.
752 753 754	Cheng Xu, Shuhao Guan, Derek Greene, M Kechadi, et al. Benchmark data contamination of large language models: A survey. <i>arXiv preprint arXiv:2406.04244</i> , 2024.
755	Shuo Yang, Wei-Lin Chiang, Lianmin Zheng, Joseph E. Gonzalez, and Ion Stoica. Rethinking benchmark and contamination for language models with rephrased samples, 2023.

756 757 758	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. HellaSwag: Can a machine really finish your sentence? In ACL, 2019.
759 760	Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In <i>ICLR</i> , 2017.
761	Cuedens Zhang Chaosi Wang Davies Ver and Deser Crease. Three much mismo of much to deser
762	regularization. <i>ICLR</i> , 2019.
763	-
764	
765	
766	
767	
768	
769	
770	
771	
772	
773	
774	
775	
776	
777	
778	
779	
780	
781	
702	
78/	
785	
786	
787	
788	
789	
790	
791	
792	
793	
794	
795	
796	
797	
798	
799	
800	
801	
802	
003	
905	
806	
807	
808	
809	
-	

810 A APPENDIX

A.1 Additional Discussion of Data Contamination Assumptions and Setting

814 Here, we discuss our data contamination approach in a bit more detail.

815 In this paper, we consider only **exact contamination**. This means we contaminate the training data 816 exactly with the text the model is later evaluated on. In the literature, it has been shown that non-exact 817 contamination (re-worded questions, translation into a different language) can affect benchmark 818 performance, too. For example, Yang et al. (2023) have shown that a 13B parameter Llama 2 Model 819 (Touvron et al., 2023) can achieve an accuracy increase of over 20 percentage points after training on 820 re-phrased benchmark questions. We decided against considering non-exact contamination for this paper because the models we train from scratch are much smaller than those for which non-exact 821 contamination results have been shown. This means these models are less capable of making sense of 822 related information, potentially leading us to underestimate the effect of non-exact contamination for 823 realistic training runs. 824

825 In addition, we consider contamination with **individual benchmark questions**, inserted into the 826 training data at **random** positions. We consider this setup because we are interested in contamination 827 from the perspective of *leakage*, where individual benchmark questions may enter the training data via different documents (for example, as quotes in Wikipedia articles, a case described in Brown 828 et al. (2020)). This contrasts with the setup where a dataset is present in the training data as a long 829 contiguous string, which we conjecture might have a similar impact but be easier detectable (Oren 830 et al., 2024). The fact that we contaminate with benchmark questions also sets us apart from related 831 works that study data contamination and memorization for random strings and uniquely identified 832 objects (Carlini et al., 2019; 2021). It is worth highlighting that the results between these two setups 833 might differ, especially considering the time it takes to forget an example. 834

835 We only consider pre-training.

836 837

A.2 ADDITIONAL DETAILS ON EVALUATION AND HOW CONTAMINATION WAS PERFORMED

Benchmark Questions and Evaluation. We use code from OLMo (Groeneveld et al., 2024) to format the different benchmark questions. This code is again based in part on the EleutherAI Evaluation Harness (Gao et al., 2024). The benchmark questions are multiple-choice, and the different options are presented to the model zero-shot as possible sentence continuations. The prediction is the sentence continuation with the largest likelihood. For the small GPT-3 models, we normalize by the number of tokens (Gao, 2021). For OLMo, we rely on the evaluation framework that is part of the code repository.

Inserting benchmark questions into the training data. A batch of LLM training data consists of *B* sequences of *S* tokens, resulting in a batch size of $B \times S$. For example, OLMo-1B is trained with B = 2048 and S = 2048; the batch for a single gradient step contains ~4M tokens (Groeneveld et al., 2024). Individual sequences in a batch usually contain multiple texts separated by a special end-of-text token. We insert benchmark questions at random positions into the pre-training data, separated at the beginning and end with the end-of-text token.

851 852

857

858

859

861

A.3 FILTERING NEAR-DUPLICATE BENCHMARK QUESTIONS

Upon close inspection of the different benchmarks, it turns out that there are exact and near-duplicate questions both within and across benchmarks. Consider, for example, the following two examples from HellaSwag (Zellers et al., 2019):

HellaSwag 1: A person is seen riding a board along the water with a kite on top. more clips are shown of the person riding back and fourth on the board.

360 and

HellaSwag 2: A person is seen riding a board along the water with a kite on top.
More clips are shown of the person riding back and fourth on the board. the person continues to ride the board along the water.

Table 2: Overview of benchmarks used in the paper. This table documents the experiments with GPT-3 models. The first two rows provide the dataset split and corresponding number of benchmark questions. The third row provides the number of questions that were removed from the dataset after filtering each dataset for near-duplicate questions. The fourth row provides the number of questions that were removed after additionally filtering for near-duplicate questions across all the different datasets combined. The fifth row provides the dataset's weight in the dataset splits used in the experiments.

	A125		.;0 ²		a crante		
	Hellas	ri0 ⁵	Social	BoolQ	MMIL	Wind	ABC.V
Split	Validation	Train	Train	Validation	Test	XL, Train	All
Size	10,042	16,113	33,410	3,269	14,042	40,398	5,197
Filtered	1,416	7,386	29,756	409	4,423	21,944	2,568
Cross-Filtered	3	6	10	2	15	0	13
Weight	19.58%	19.77%	8.27%	6.48%	21.82%	18.16%	5.92%

Note that these are the ground-truth options of two *different* benchmark questions. Similar patterns can be observed for many benchmarks, for example, because a single text document was used to create different benchmark questions. Here is another example from PiQA (Bisk et al., 2020):

PiQA 1: Goal: how do you close a cabinet? Solution: push the door shut.

and

882 883

884

885

886 887

888

889 890

891

PiQA 2: Goal: how do you close a cabinet? Solution: shut the door.

Near-duplicate benchmark questions present a challenge for our methodology. This is because one
question might end up in a set of benchmark questions that we highly contaminate with, whereas
the other question might end up in the purportedly "clean" set of benchmark questions that we
use for holdout evaluation. To tackle this problem, we perform fuzzy string matching between the
ground-truth options (that is, the potential contamination data) of all benchmark questions, randomly
removing one question for every detected duplicate. We use the Python package rapidfuzz.

Summary Statistics. Table 2 depicts summary statistics about the different benchmarks, including the number of questions that were filtered during the duplicate-detection stage. We see that the number of filtered questions is significant. On some datasets, especially Social-i-QA, we had to apply very aggressive filtering to avoid any side-effects during contamination. Hence, the number of removed questions per dataset does not necessarily reflect the actual number of duplicates, but the level of filtering that had to be applied to remove all duplicate questions.

904 Experimental verification that filtering worked. We verify that our filtering procedure worked by 905 training two models: One that is heavily contaminated (obtaining an accuracy of over 97%), and 906 another model that did not see any contamination. We then evaluate both models on a set of 10,000 907 benchmark questions that are holdout even for the contaminated model. The contaminated model 908 obtains an accuracy of 42.2%, (95%-CI: 41.2% - 43.2%) on the holdout, while the clean model obtains an accuracy of 41.9% (95%-CI: 41.0% -42.9%). Because the observed accuracy difference 909 is small in absolute terms and lies within the confidence interval, we conclude that there are no 910 significant side-effects in our evaluation procedure. 911

912

913 A.4 PROOF OF PROPOSITION 1

Proof. Without loss of generality, mapping $t_1 = 1$ and $T = t_2$, we have from equation 3:

$$w_1^T = \prod_{i=1}^T (1 - \lambda_i \gamma)$$

Assigning the forgetting ratio to be less than ϵ according to our criteria, we have:

 $\prod_{i=1}^{I} (1 - \lambda_i \gamma) \le \epsilon$

$$\sum_{i=1}^{T} \log(1 - \lambda_i \gamma) \le \log \epsilon$$
$$\sum_{i=1}^{T} (-\lambda_i \gamma) \lesssim \log \epsilon \quad (\log(1 - x) \approx -x \text{ for small } x)$$
$$T \times \left(\frac{1}{T} \sum_{i=1}^{T} \lambda_i\right) \times \gamma \gtrsim \log \frac{1}{\epsilon}$$

Re-arranging this equation gives us the desired result, where $\lambda_{avg} = (\frac{1}{T} \sum_i \lambda_i)$.

A.5 EXTENDED ANALYSIS OF FORGETTING & GRADIENT ALIGNMENT

To understand the effect of weight decay on forgetting, we analyze two different stages of optimization:
(1) the contamination stage, where the training set consists of only the contaminated samples, and (2)
the forgetting stage, where the training set consists only of clean samples. Here, we consider the SGD
learning algorithm, but the resulting analysis also applies to SGD with momentum. In particular, to
sGD with weight decay for the forgetting stage.

We now introduce some notation. Let $\theta \in \mathbb{R}^D$ be the weights of the model, and let $\mathcal{X}_{cont} = \{\mathbf{x}_1^{cont}, \mathbf{x}_2^{cont}, ..., \mathbf{x}_{N_{cont}}^{cont}\}$ be the contamination set and $\mathcal{X}_{clean} = \{\mathbf{x}_1^{clean}, \mathbf{x}_2^{clean}, ..., \mathbf{x}_{N_{clean}}^{clean}\}$ be the clean pre-training set. The training data used for the contamination stage is thus \mathcal{X}_{cont} , and the training data for the forgetting stage is \mathcal{X}_{clean} . Let $\ell(\mathbf{x}_i) \in \mathbb{R}$ be the loss associated with sample \mathbf{x}_i . Let the model be initialized at the contamination stage with $\theta = \theta_{init}$. The learning algorithm is (single batch size) SGD, which is run for a total of N_{cont} steps for the contamination stage and N_{clean} steps for the forgetting stage. First, we observe that the weights at the end of the contamination stage are:

$$\theta' = \theta^{\text{init}} - \underbrace{\sum_{i=1}^{N_{cont}} \lambda_i \nabla_{\theta_i} \ell(\mathbf{x}_i^{\text{cont}})}_{\theta^{\text{cont}}}$$

We thus denote the weights at the end of the contamination stage as $\theta' = \theta^{\text{init}} + \theta^{\text{cont}}$. For the subsequent fine-tuning stage to forget information regarding samples $\mathcal{X}_{\text{cont}}$, we first define a criterion to identify forgetting based on the angle between a weight vector and the contaminated part identified above.

Forgetting Criteria. A model with weights θ is said to have "forgetten" information contained in θ^{cont} if $\frac{|\theta^{\top} \theta^{\text{cont}}|}{||\theta^{\text{cont}}||_{2}^{2}} \leq \epsilon$, which we call the "forgetting ratio". Here, $\epsilon \in \mathbb{R}^{+}$ is a small constant.

We now proceed with an analysis of the forgetting stage. To enable this, we make the following important assumption that the gradients of clean and contaminated samples are orthogonal for all clean and contaminated samples across all optimization steps. This is a relatively strong assumption, and it quantifies the intuition that model updates required by SGD to memorize clean samples and contaminated samples are distinct.

971 Assumption 1. (Gradient Orthogonality) $\nabla_{\theta t_1} \ell(\mathbf{x}_i^{\text{clean}})^\top \nabla_{\theta t_2} \ell(\mathbf{x}_j^{\text{cont}}) = 0, \forall i \in [1, N_{\text{clean}}], \forall j \in [1, N_{\text{cont}}] \text{ and } \forall \text{ steps } t_1, t_2.$

We also make another minor simplifying assumption that the weight initialization is orthogonal to both these quantities.

Assumption 2. (Gradient-Initialization Orthogonality) $\nabla_{\theta_t} \ell(\mathbf{x}_i)^\top \theta_{init} = 0, \forall i \text{ and } \forall \text{ steps } t.$

Given these assumptions, we are ready to state our result.

Proposition 3. (Forgetting Time) The number of optimization steps T_{forget} in the forgetting stage, 978 such that the weights $\theta_{T_{forget}}$ satisfy the ϵ -forgetting criteria is given by: $T_{forget} \gtrsim \frac{\log(1/\epsilon)}{\lambda_{avg}\gamma}$, where 980 $\lambda_{avg} = \frac{1}{T} \sum_{i=1}^{T} \lambda_i$ is the average learning rate of the optimizer.

Proof. We first compute the forgetting ratio at $\theta = \theta'$, and as a consequence of Assumption 2, verify that the forgetting ratio is equal to one.

Let us now denote these weights as $\theta'_0 = \theta'$, used as initialization for the forgetting stage. Analyzing the first optimization step, and the subsequent forgetting ratio, we have:

(Optimization Step)
$$\theta'_{1} = \theta'_{0} - \lambda_{0} \nabla_{\theta 0} \ell(\mathbf{x}_{0}^{\text{clean}}) - \lambda_{0} \gamma \theta'_{0}$$

(Forgetting ratio) $\frac{|\theta'_{1}^{\top} \theta^{\text{cont}}|}{\|\theta^{\text{cont}}\|_{2}^{2}} = \frac{|(\theta'_{0} - \lambda_{0} \nabla_{\theta 0} \ell(\mathbf{x}_{0}^{\text{clean}}) - \lambda_{0} \gamma \theta'_{0})^{\top} \theta^{\text{cont}}|}{\|\theta^{\text{cont}}\|_{2}^{2}}$
 $= \frac{|((\theta^{\text{init}} + \theta^{\text{cont}}) - \lambda_{0} \nabla_{\theta 0} \ell(\mathbf{x}_{0}^{\text{clean}}) - \lambda_{0} \gamma ((\theta^{\text{init}} + \theta^{\text{cont}})))^{\top} \theta^{\text{cont}}|}{\|\theta^{\text{cont}}\|_{2}^{2}}$
 $= (1 - \lambda_{0} \gamma)$ (From Assumptions 1 & 2)

We can similarly analyze the subsequent optimization steps to compute the forgetting ratio, which for some step t + 1 is:

998
999
1000
$$\frac{|\theta_{t+1}^{\prime}^{\top}\theta^{\text{cont}}|}{||\theta^{\text{cont}}||_{2}^{2}} = \frac{|(\theta_{t}^{\prime} - \lambda_{t}\nabla_{\theta t}\ell(\mathbf{x}_{t}^{\text{clean}}) - \lambda_{t}\gamma\theta_{t}^{\prime})^{\top}\theta^{\text{cont}}|}{||\theta^{\text{cont}}||_{2}^{2}}$$

Forgetting ratio at t + 1

1002
1003
1004
1005
$$=\underbrace{\frac{|\theta_t'^\top \theta^{\text{cont}}|}{\|\theta^{\text{cont}}\|_2^2}}_{\text{Forgetting ratio at }t} (1 - \lambda_t \gamma)$$

1006 Unrolling the recurrence till step T, we have that:

$$\frac{|\theta_T^{\prime} \,^{\top} \theta^{\text{cont}}|}{\|\theta^{\text{cont}}\|_2^2} = \underbrace{\frac{|\theta_0^{\prime} \,^{\top} \theta^{\text{cont}}|}{\|\theta^{\text{cont}}\|_2^2}}_{= 1} \times \prod_{i=1}^T (1 - \lambda_i \gamma)$$

1012 Assigning the forgetting ratio to be less than ϵ according to our criteria, we have:

1013
1014
$$\prod_{i=1}^{T} (1 - \lambda_i \gamma) \le \epsilon$$
1016
$$T$$

 $\sum_{i=1} \log(1 - \lambda_i \gamma) \le \log \epsilon$

1019
1020
$$T = 1$$

$$T$$

1020
1021
1022

$$\begin{pmatrix} -\chi_i \gamma \end{pmatrix} \gtrsim \log e^{-1} (\log(1-x) \approx -x \text{ for small } 1)$$

1023
$$T \times \left(\frac{1}{T}\sum \lambda_i\right) \times \gamma \gtrsim \log \frac{1}{\epsilon}$$

1024
$$(i=1)$$

Re-arranging this equation gives us the desired result, where $\lambda_{avg} = (\frac{1}{T} \sum_{i} \lambda_i)$.



Figure 7: Increasing the weight decay parameter leads to faster forgetting described by the cumulative weight decay. We continue training the 124M parameter model from Section 4.2 with different choices of the weight decay parameter. The figure depicts empirical forgetting (dashed line) and our theoretical forgetting bound, the cumulative weight decay (solid line). As we increase the weight decay parameter, the theoretical bound predicts much faster forgetting. The figure shows that empirical forgetting is indeed as fast as predicted by the bounds. Note that the figures have very different x-axis scales, ranging from 120 gradient steps in (a) to 62500 gradient steps in (d). This Figure depicts the same quantities as Figure 6 in the main paper.



Figure 8: Forgetting without weight decay. We perform the forgetting experiment from Section 4.2 in the main paper without weight decay (instead of the default weight decay of 0.1). This Figure depicts the same quantities as figure 3 in the main paper. We see that there is significant forgetting without weight decay. However, compared with a weight decay of 0.1, the rate of forgetting is somewhat slower, especially for 144 times repeated contamination. In particular, the accuracy gaps in (b) are slightly larger, and the difference in (c) is still statistically significant (compare Figure 3).



Figure 9: Scaling up with a 7B parameter model. We perform the forgetting experiment from Section 4.3 in the main paper with OLMo-7B. We contaminate the OLMo-7B checkpoint at gradient step 300,000 four times with different benchmarks. This causes an average accuracy increase of 17 percentage points. We then continue pre-training for 1% of the remaining training time, leading to a reduction of 77% of the accuracy increase due to contamination. Mean and bootstrapped 90% confidence intervals. We will continue training this model for the final version of the paper.