
Standard fine-tuning inefficiently uses rare data

Anonymous Author(s)

Affiliation
Address
email

Abstract

What is the best way to improve performance on a task with a small number of data points? For such applications, it is common to initialize training with a language model pre-trained on general web data instead of training from scratch. We find that on tasks such as math, coding, and instruction following, simply replaying pre-training data while fine-tuning can outperform standard fine-tuning, effectively training on up to $1.87\times$ more samples. Though replay is often used to reduce catastrophic forgetting on previous domains, we surprisingly find that it is helpful even if you only care about loss on the new rare data. Next, we study how much additional rare data efficiency we can achieve by changing the data mixture across all of training. By searching over possible two-stage mixture schedules, we can improve data efficiency by up to $4.8\times$ over the new baseline. Interestingly, we find cases where the optimal data schedule only requires replay, removing the need to allocate rare data to pre-training. We use the findings from our controlled experiments to improve performance for specific rare fine-tuning setups, increasing agentic web navigation accuracy by 4.5% and Basque question-answering accuracy by 2%.

1 Introduction

The standard language modeling paradigm is pre-training on a large amount of web data and fine-tuning on a task of interest (e.g. math) [Radford et al., 2018, Devlin et al., 2019]. This strategy enables better performance on this task even when there are only a few relevant data points to train on. Typically, practitioners fine-tune language models vary hyper-parameters such as learning rate and repetition count with the goal of lowering loss on the rare data.

In this paper, we investigate how to best leverage rare data points in relation to pre-training. To rigorously answer this question, we set up a controlled experimental setting with two data distributions: common data, representing generic web text, and rare data, representing few tokens for a task of interest (i.e. math, associated with high-quality tokens from FineMath). Since we only pre-train to improve rare loss, we evaluate performance purely by measuring validation loss on the rare task.

First, we determine whether we can improve upon standard fine-tuning. We carefully tune a fine-tuning baseline under standard practice of pre-training on web data followed by fine-tuning on rare data (with two separate learning rate schedule cycles). We show that we can improve performance on the rare task by simply replaying pre-training data while fine-tuning, even if we only care about the rare task, improving data efficiency by up to $1.87\times$. This is surprising considering that the replay data is generally unrelated to the rare task. Though we do not fully understand the cause, we offer candidate hypotheses involving unstable optimization and overfitting to finite samples.

Next, we determine how much additional data efficiency we can achieve by modifying the full training stack. Instead of using separate learning schedules for pre-training and fine-tuning, we use a single learning rate schedule, closer to common practice in mid-training. We establish a new

38 baseline without replay that achieves up to $9.42 \times$ better data efficiency than the previous fine-tuning
39 baseline. We suspect that traditional fine-tuning with a split learning rate schedule underperforms
40 this baseline because it rewarms up the learning rate and resets the optimizer state, and we call for
41 model developers to release checkpoints prior to annealing for fine-tuning.

42 Under this new strategy, we determine the optimal dynamic mixture across two training stages,
43 mimicking pre-training and fine-tuning. In addition to using replay data (which sets the mixture for
44 the second stage), we can decide what fraction of the rare data we allocate to the first stage instead of
45 the second stage. Under this framework, we thoroughly search for the best data schedule for rare
46 data loss, which is up to $4.8 \times$ data efficient than the baseline. We find that increasing the replay
47 fraction is generally most important when we allocate no rare data to the first stage. For coding data
48 in particular, replay is so strong that it is best to not allocate any rare data to the first stage.

49 Overall, our experiments establish clear principles for leveraging rare data. For fine-tuning, replay is
50 a simple and effective strategy for many data-limited tasks. We apply this to rare tasks such as web
51 navigation with agents (improving score by 4.5%) and low-resource language learning for Basque
52 (improving question-answering accuracy by 2%). For pre-training, we find that a majority of the
53 benefit of introducing data early can be achieved by improving fine-tuning with replay. Knowing
54 when this is the case is important for many downstream applications in using language models with
55 private data, new facts, and domain-specific corpuses.

56 2 Improving rare data efficiency of fine-tuning

57 2.1 Setup

58 **Training.** We have two pools of training data: **common data**, representing generally useful web
59 data, and **rare data**, representing data relevant for a task of interest (i.e. math). We will train a
60 language model on 4 billion tokens of data from these pools. Since web data is abundant relative
61 to rare data, we do not constrain the web data and never repeat these samples. Instead, we model a
62 constraint on the amount of rare data, which we limit to 4 million tokens. Under this setup, we are
63 allowed to choose the learning rate schedule and how many times we repeat the rare data. We use C4
64 for our common domain and FineMath (math), StarCoder (coding), and Flan (instruction following)
65 for our rare domains. We train a 150 million parameter Llama-style language model ([Grattafiori
66 et al., 2024]) using AdamW. We give full training details in Appendix D.

67 **Evaluation.** We are only interested in performance on the real domain. We evaluate this by
68 measuring loss on a held-out validation set from the rare data distribution. We choose validation loss
69 since it scales much more smoothly than accuracy metrics, especially for models at our scale.

70 **Rare data efficiency.** Given two fine-tuning strategies S_1 and S_2 , we would like to compare how
71 “efficiently” they are using the data. To measure this, we fit a scaling law to a fixed reference training
72 strategy S_{ref} as it gets access to more rare data. For any training strategy S we want to evaluate,
73 we can measure the loss of the final model and recover how much rare data the reference strategy
74 would need to match this loss, denoted $D(S)$. However, this quantity depends on the effectiveness
75 of the reference strategy. To remove this dependence, we always report data efficiency as a relative
76 improvement from one strategy to another, or $\frac{D(S_2)}{D(S_1)}$. A data efficiency improvement of $k \times$ now
77 corresponds to “training with this strategy is like effectively training on k times more rare data”. We
78 give more details on how we fit the scaling laws in Appendix G.

79 2.2 Fine-tuning baseline

80 Suppose we train for a total of T steps, with γ fraction of them being rare. Standard fine-tuning
81 corresponds to training on common data with a cosine learning rate schedule for $(1 - \gamma)T$ steps.
82 Then, we train on the rare data for γT steps with a separate cosine learning rate schedule. In between
83 the stages, we reset the optimizer state (containing moving estimates of the gradients for AdamW)
84 matching practice for fine-tuning on open-weight models. We provide a competitive fine-tuning
85 baseline by tuning the two main choices: the learning rate and the number of rare data repetitions
86 (exact procedure in Appendix F.1). In summary, it is best to use a learning rate of 1e-4 for pre-training
87 and 3e-4 for fine-tuning on 64 repetitions of the rare data. If we try repeating the data more than

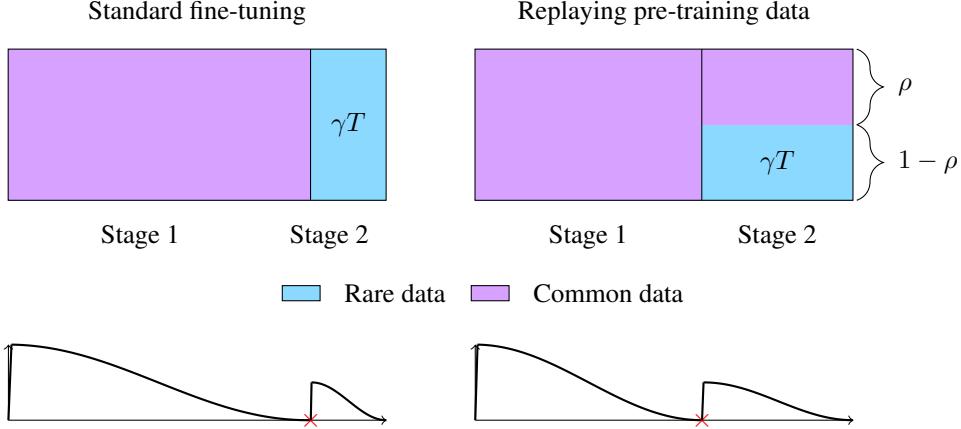


Figure 1: Controlled fine-tuning visualization. We systematically explore the benefit of replaying data from the pre-training domain while fine-tuning on the rare data. On the right, we show standard fine-tuning for T steps where γ fraction of the steps are on rare data. On the left, we show fine-tuning with replay fraction ρ (where we shorten pre-training to keep the total number of steps fixed). We use (independently tuned) cosine learning rate schedules for each stage, with an optimizer state reset between the stages to simulate standard practice for fine-tuning open-weight models.

88 this, the validation loss will start to increase akin to classical overfitting (this is not captured by prior
 89 data-constrained scaling laws, discussed in K.1). This setup serves as $1 \times$ rare data efficiency.

90 2.3 Replay improves data efficiency

91 We introduce a simple strategy that improves loss on the rare task: mix web data while fine-tuning on
 92 the rare data. Specifically, we introduce a replay fraction ρ which dictates what fraction of fine-tuning
 93 steps will be on common data. When we increase this replay fraction, we decrease the number of
 94 pre-training steps to conserve the total number of training steps (Figure 1, right). In Figure 2, we
 95 show how the final model’s loss depends on the replay ratio. We find that for each domain, there is
 96 an optimal replay fraction that maximizes the data efficiency gains (indicated by the starred points).
 97 Replay improves data efficiency by $1.87 \times$ for Flan, $1.49 \times$ for FineMath, and $1.09 \times$ for StarCoder.
 98 Note that even if we only care about loss on the fine-tuning task, it is helpful to add the less relevant
 99 web data while training. We observe that code, which is the most rare with respect to C4 since C4
 100 was filtered to remove code, can tolerate less replay data than the relatively more common domains
 101 of math and instruction following, though the loss improvements are similar across domains.

102 Though replay is a common method in continual learning, it is almost always to prevent catastrophic
 103 forgetting of old tasks [Rolnick et al., 2019, Parisi et al., 2019]. Interestingly, we find that replay
 104 improves performance on the new in-distribution training task, which departs from the standard
 105 intuition. We provide a more detailed discussion in Section 5.

106 3 Improving rare data efficiency of pre-training

107 In the previous section, we limited ourself to using replay during fine-tuning. In this section, we
 108 aim to understand how much additional data efficiency we get by setting the data mixture for two
 109 stages of training, akin to changing both pre-training and fine-tuning. We first unify the optimization
 110 process of pre-training and fine-tuning into a single learning rate schedule cycle with no optimizer
 111 reset. We now search over possible ways to schedule the rare data by choosing a replay ratio for stage
 112 2 as well as how much rare data is used in stage 2 (equivalently, setting the mixture for two stages of
 113 training). This extra flexibility allows us to quantify how much benefit we get from introducing rare
 114 data earlier in training (stage 1) instead of reserving it all for the end (stage 2).

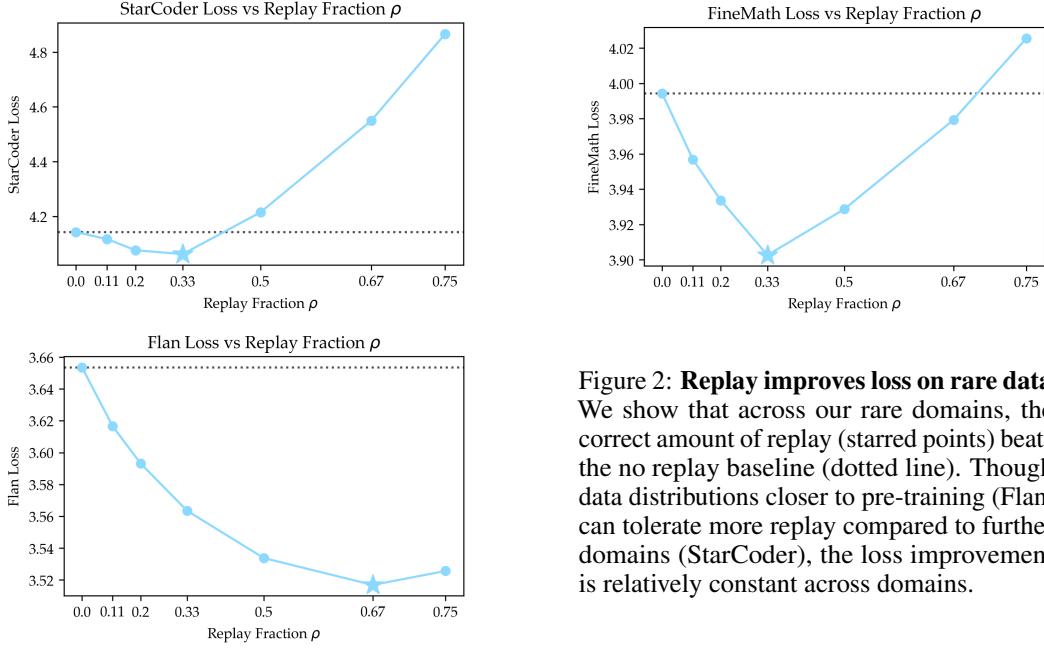


Figure 2: Replay improves loss on rare data.
We show that across our rare domains, the correct amount of replay (starred points) beats the no replay baseline (dotted line). Though data distributions closer to pre-training (Flan) can tolerate more replay compared to further domains (StarCoder), the loss improvement is relatively constant across domains.

115 3.1 WSD baseline

116 To start we tune the baseline of seeing all common data before all the rare data. Similar to before,
 117 we tune the epoch count for this data schedule and find a slightly lower optimal epoch count of 32
 118 (Appendix E.1.1). However, we find that the correct learning rate schedule is critical for rare data
 119 efficiency. Default practice for learning rate schedules (i.e. cosine, linear) is to slowly anneal to zero
 120 over the course of training. Recent work in mid-training has suggested using a warmup-stable-decay
 121 (WSD) learning rate schedule [Hu et al., 2024]. This consists of a short linear warmup, a stable
 122 training phase, and a sharp linear decay for a variable fraction of training referred to as the cooldown
 123 period. Interestingly, during the learning rate decay, the loss decreases at a much faster rate than the
 124 rest of training. This can be exploited to get stronger performance on target data by placing it at the
 125 end of training. We explain and visualize these benefits in more detail in Appendix H.

126 In Figure 3, we show that WSD offers a significant benefit over the more traditional schedules that
 127 decay the learning rate over all of training for the fine-tuning baseline. For our setting, we find it best
 128 to fix a cooldown period of 10% of training for all domains, increasing data efficiency by $28.47 \times$
 129 relative to the worst cooldown duration. In contrast to the fine-tuning experiments, we find it better to
 130 use a slightly higher learning rate of $3e-3$ for all of training (details in Appendix E.1.2).

131 Due to the learning rate schedule, we refer to this strategy as the WSD baseline. This increases
 132 data efficiency relative to the fine-tuning baseline by $9.92 \times$ for Starcoder, $6.37 \times$ for FineMath, and
 133 $2.77 \times$ for Flan. This is likely because the joint training doesn't reset optimizer state and rewarmp
 134 the learning rate. As such, we believe fine-tuning would benefit from initializing at a pre-annealed
 135 pre-training checkpoint instead of the final checkpoint. We call on model developers to release the
 136 model and optimizer state before cooldown since this is more useful for downstream applications.

137 3.2 Data schedule space

138 Now that we have fairly tuned our WSD baseline, we are interested in whether this is the best
 139 data schedule across pre-training and fine-tuning. To reduce the complexity of searching over all
 140 permutations, we consider data schedules of two stages. We get to pick the fraction of rare data
 141 for each stage, subject to the constraint on the total number of rare steps. We fix the sample-level
 142 ordering of data points within each domain to reduce variance when comparing different schedules.

143 This space now only has two degrees of freedom. There are multiple ways of parameterizing it,
 144 and we decide to use **replay fraction** ρ (how much common data is replayed) and the **rare stage**

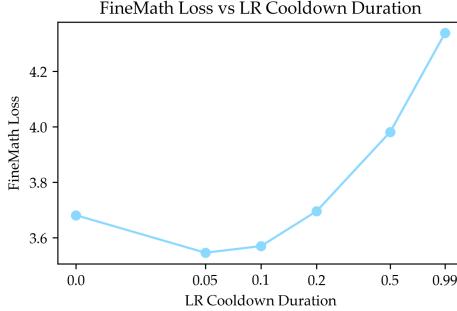


Figure 3: Tuning learning rate cooldown. We tune how long we should cool down the learning rate for WSD. The above plot shows the optimal cooldown period is between 0.05 and 0.1; we use 0.1 for consistency across domains and being fair to changing data schedules.

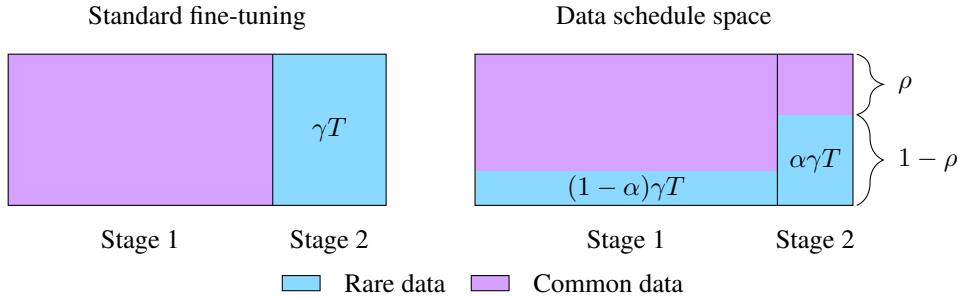


Figure 4: Controlled mid-training visualization. We systematically explore the space of data schedules when training on T tokens where γT are rare. A data schedule allocates an α fraction to stage 2 where stage 2 has a replay fraction ρ . Standard fine-tuning puts all rare data at the end with no replay ($\alpha = 1, \rho = 0$).

145 **2 allocation α** (what fraction of the total rare data is allocated to stage 2). We provide a more
 146 intuitive visualization in Figure 4. Fine-tuning now has a simple interpretation: no replay data
 147 ($\rho = 0$) and all the rare data is used in stage 2 ($\alpha = 1$). Finding the optimal data order under this
 148 parameterization is equivalent to finding the best setting of ρ and α . We provide a detailed discussion
 149 of the parameterization and equivalences in Appendix A.

150 3.3 Best data schedule

151 To determine the best data schedule, we sweep over replay fraction ρ and rare stage 2 allocation α .
 152 We are interested in three strategies: the WSD baseline of stage 1 fully common and stage 2 fully
 153 rare ($\rho = 0, \alpha = 1$), fine-tuning while replaying common data in stage 2 ($\alpha = 1$), and the full space
 154 of modifications (all settings). In Table 6, we share the data efficiency improvements of each strategy.

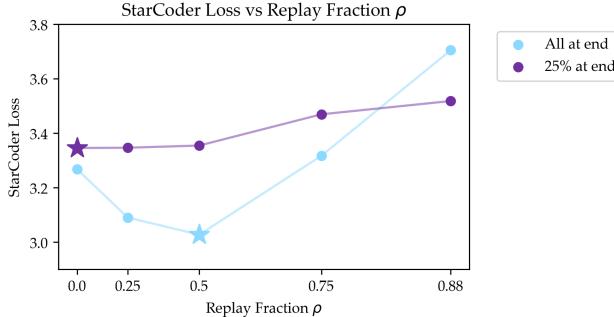


Figure 5: Importance of replay fraction depends on rarity. When all the rare data is seen during fine-tuning, tuning the replay fraction becomes critical to improve loss (blue line). When we change pre-training to see some of the rare data, tuning the replay fraction is not important and can sometimes hurt loss (purple line).

Domain	Replay	Full Schedule
StarCoder	1.53×	1.53×
FineMath	1.85×	2.49×
Flan	2.06×	4.80×

Figure 6: **Data efficiency across domains.** We compare the data efficiency of best replay and best schedule relative to the fine-tuning baseline.

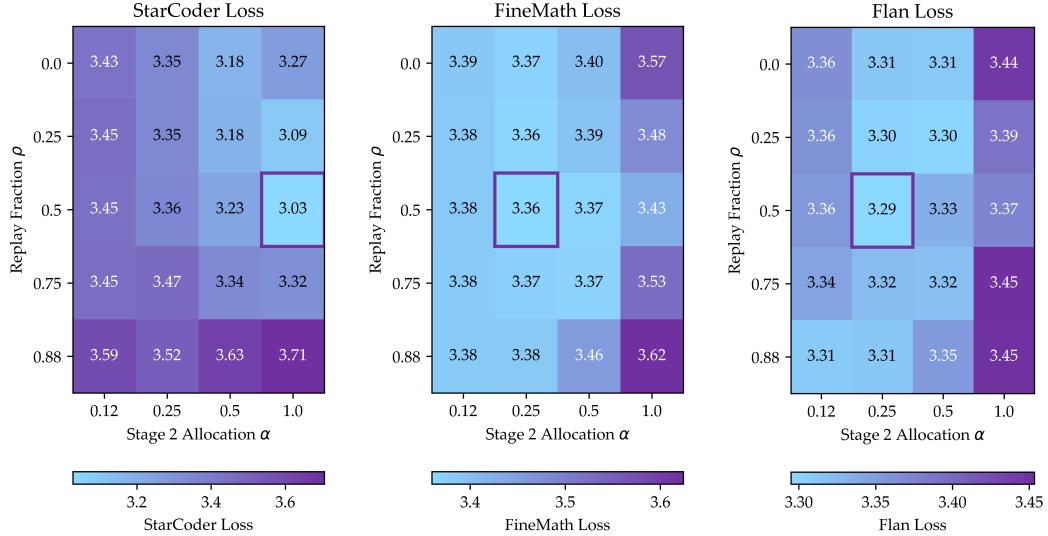


Figure 7: **Full data schedule sweep.** We sweep over data schedules, parameterized by their replay fraction and fraction of rare data allocated to stage 2. Standard fine-tuning (no replay and all rare data in stage 2, top right corner) achieves the worst loss for FineMath and Flan. This can be improved by adding replay data (right column). This can also be solved by having some rare data in stage 1, in which case replay becomes less important.

155 3.3.1 Replay is more helpful when data is more rare

156 In addition to a direct comparison of these strategies, we can determine when replay is important.
 157 We first find that when the rare data is unseen during stage 1, replay is critical for improving loss
 158 (example for Starcoder in Figure 5, blue). On the other hand when we keep 75% of the rare data for
 159 pre-training, replay is no longer helpful (Figure 5, purple). Together, this indicates that replay is more
 160 helpful as the data is less present during stage 1. One interpretation of these results is that training
 161 benefits from a smooth transition when changing data distributions, whether it comes from adding
 162 web data to stage 2 or rare data to stage 1.

163 We show the full results of sweeping over replay fraction and stage 2 allocation in Figure 7. In the
 164 context of this plot, pure fine-tuning is the top-right entry, fine-tuning with replay is the right column,
 165 and the full space of modifications is the entire plot.

166 3.4 Do we need to change pre-training?

167 One natural question for applications is whether one needs to change pre-training (correspondingly,
 168 stage 1) to maximally leverage task-relevant data. We can use the full sweep over data schedules to
 169 determine when it's necessary to modify pre-training. We find that for FineMath and Flan, we can
 170 not get the full benefits of the optimal data order by only changing stage 2 (achieving 67.4% of gains
 171 for FineMath and 46.0% of the gains for Flan measured logarithmically). On the other hand, for
 172 Starcoder, the optimal data schedule only requires adding replay data to stage 2. It is encouraging that
 173 we can get away with not using data early since it is prohibitive or impossible to change pre-training
 174 for many applications. We discuss these implications in more detail in Section 4.2.

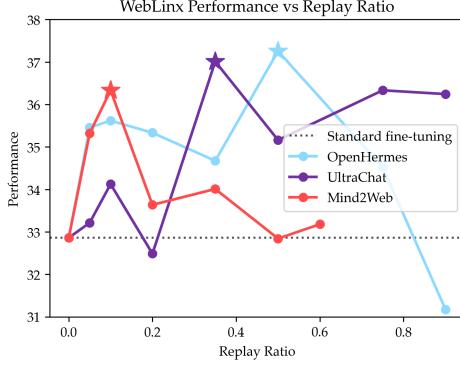


Figure 8: **Weblinx.** We fine-tune Llama 3.1-8B Instruct on Weblinx demonstrations. Without any replay data, we get 32.86% accuracy following the original hyper-parameters. We find that mixing generic instruction following data (OpenHermes, UltraChat) improves accuracy by up to 4.5%. Surprisingly, this is even better than replaying demonstrations from an alternative web agent task (Mind2Web).

175 4 Implications

176 4.1 Rare post-training benefits from replay

177 Our findings offer an immediate prescription for practitioners hoping to train on rare data: replay
 178 pre-existing data. Though this sounds like a simple modification, this is rarely done while fine-tuning
 179 in practice as it is commonly believed that fine-tuning is the correct way to leverage rare data. For
 180 most fine-tuning setups, the pre-trained model is downloaded off the internet with the learning rate
 181 fully annealed and optimizer state reset, close to the setup in 2. Moreover, since we usually do not
 182 have access to the common data distribution, we have to pick an approximation for our replay data.
 183 We note that using a replay fraction of ρ requires $\frac{1}{\rho}$ times as many training steps. This is generally a
 184 favorable tradeoff for fine-tuning since we are rarely compute-constrained. We demonstrate that for
 185 the data-constrained domain of web agents and the low-resource domain of Basque, replay is able to
 186 provide increased accuracy.

187 4.1.1 Web Agents

188 Recently, language models have been trained to perform agentic tasks, requiring advanced capabilities
 189 such as web navigation. However, collecting demonstrations requires expensive human labor and
 190 training data is scarce. Therefore, it is important to understand how to best leverage this data.

191 We follow the data and evaluation protocol from Weblinx [Lù et al., 2024]. Since this task requires
 192 instruction following, we start training with Llama 3.1-8B Instruct on a fixed number of rare demon-
 193 strations. For the replay data, we use OpenHermes [Teknium, 2023] or UltraChat [Ding et al., 2023]
 194 instruction-following data instead of web data since we are fine-tuning an instruct model. For training,
 195 we use the hyperparameters from the paper’s training pipeline. We measure accuracy on the validation
 196 and in-distribution test set under the paper’s established static evaluation metric. We defer to the
 197 original paper for more details on the data and evaluation.

198 We find that when training on web agents data, there is a consistent advantage to replaying instruction
 199 following data. In Figure 8, we show that replaying instruction following data improves accuracy by
 200 up to 4.5%. We provide additional details and experiments in Appendix I.

201 4.1.2 Basque

202 Basque is a low-resource language spoken by approximately 800k people in the Basque Country.
 203 As such, this language is scarce on the web and pre-training datasets, constituting only 0.035% of
 204 Common Crawl [Etxaniz et al., 2024]. However, thanks to a thriving NLP research community, there
 205 is a large amount of Basque data available through the Latxa corpus [Etxaniz et al., 2024].

206 We are interested in how to fine-tune a model to have high Basque understanding. We start continual
 207 pre-training from Llama 3.1-8B with access to 200M Basque tokens. For replay data, we use the
 208 SlimPajama replication [Soboleva et al., 2023, Weber et al., 2024] as a proxy for the unreleased
 209 Llama pre-training data. For evaluation, we measure validation loss as before, as well as accuracy on
 210 a professional Basque translation [Baucells et al., 2025] of the commonsense reasoning benchmark
 211 COPA [Gordon et al., 2012, Ponti et al., 2020] supported on lm-eval-harness [Gao et al., 2024].

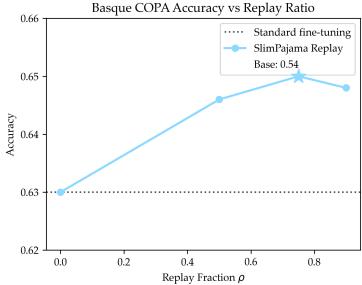


Figure 9: Basque. We fine-tune Llama 3.1-8B on 200M Basque tokens from the Latxa training corpus and measure accuracy on Basque COPA. We find that replaying generic pre-training data from SlimPajama improves accuracy by up to 2%.

212 We find that when training on Basque data, there is a consistent advantage to replaying pre-training
 213 data. In Figure 9, we show that the model achieves higher accuracy on the Basque evaluation task.
 214 We provide more details and experiments in Appendix J.

215 4.2 Can fine-tuning subsume pre-training?

216 Since it's expensive to modify pre-training, practitioners only control fine-tuning in most applications.
 217 For example, companies want models to utilize internal data for their operations, model developers
 218 want to update models with new knowledge about the world, and domain experts want models
 219 specialized for their field. In all these cases, it is prohibitively expensive to retrain a model from
 220 scratch. In these scenarios, how much do we lose by not changing pre-training?

221 We find evidence in the StarCoder setting that the optimal data schedule only requires tuning replay
 222 ratio. This shows that for some domains, it's possible to get the benefits of scheduling data without
 223 changing pre-training. However, this only works if we introduce replay data and use a single learning
 224 rate schedule. It will be important to develop better fine-tuning methods to match altering pre-training,
 225 as well as delineate when it's possible in the first place.

226 5 Related work

227 **Mid-training.** Many recent language models augment pre-training with a mid-training phase where
 228 they reduce the learning rate while training on high-quality data OLMo et al. [2025], Grattafiori et al.
 229 [2024], Li et al. [2025], Nvidia et al. [2024]. There has been some initial work on characterizing the
 230 benefit of putting rare data at the end of training [Aryabumi et al., 2024, Blakeney et al., 2024] or
 231 annealing the learning rate [Hu et al., 2024]. In addition to prior knowledge, we conduct experiments
 232 on changing the pre-training in conjunction with mid-training. Moreover, we show new experiments
 233 at maximal repetition count, as well as for various ablation factors such as model size.

234 **Robust fine-tuning.** There is a rich literature on how to robustly fine-tune language models to
 235 maximize in-distribution and out-of-distribution accuracy [Phang et al., 2019, Zhang et al., 2021,
 236 Kumar et al., 2022]. Weight averaging has been one such technique to improve post-training
 237 performance [Wortsman et al., 2022, Ilharco et al., 2023, Dang et al., 2025]. Replay can be seen as
 238 qualitatively similar to weight averaging where the averaging takes places in data distribution space
 239 instead of parameter space. In contrast to prior work, we characterize the interaction between pre-
 240 training and fine-tuning, showing that the optimal way to fine-tune depends on how much exposure the
 241 pre-trained model has to the target task. Moreover, since the focus was primarily out-of-distribution
 242 performance, they under-focused on the opportunity to improve in-distribution performance.

243 **Optimizing data mixtures.** Prior work has proposed algorithms to optimize the data mixture
 244 [Chen et al., 2023, Xie et al., 2023, Jiang et al., 2024a, Fan et al., 2024]. Most such algorithms fall
 245 under an online optimization framework [Chen et al., 2025b], where the algorithm estimates which
 246 components it should upweight to maximize performance. However, such online algorithms miss the
 247 most important aspect of data mixtures: better data should be at the end. Instead, such algorithms
 248 greedily upweight the best data at the start since they do not factor in constraints on the number of
 249 available data points. Moreover, they do not account for the optimization challenges that arise when
 250 changing data distributions. Though our work uses relatively unsophisticated methods, we are able to
 251 answer more general questions and achieve stronger performance improvements.

252 **Continual learning.** There has been a lot of work on continual learning for new tasks [Rolnick
253 et al., 2019, Parisi et al., 2019]. Such works have traditionally focused on reducing catastrophic
254 forgetting [Kirkpatrick et al., 2017] instead of improving target task performance [Gupta et al., 2023,
255 Ibrahim et al., 2024, Kotha et al., 2024, Çağatay Yıldız et al., 2025, Chen et al., 2025a, Springer
256 et al., 2025]. There has also been work on methods and evaluation for teaching models new facts
257 [Meng et al., 2023, Yang et al., 2024b, Ghosal et al., 2024, Gekhman et al., 2024, Chang et al., 2024].
258 Our two-stage framework helps build intuition for when pretraining is necessary, as well as shows
259 better ways to teach models new facts. However, answering further questions about capability and
260 knowledge will require using more refined metrics and data distributions.

261 **Curriculum learning.** Curriculum learning is generally concerned with proposing a sequence of
262 training distributions from easy to hard [Bengio et al., 2009]. Theoretically, it has been shown that
263 this can accelerate the speed of convergence over training by introducing tractable intermediate tasks
264 [Abbe et al., 2023, Panigrahi et al., 2024]. Recent works have tried to design curricula using reference
265 models [Mindermann et al., 2022, Fan and Jaggi, 2023, Lin et al., 2025] or structure over the data
266 distribution [Chen et al., 2023]. However, there is limited evidence that changing data order improves
267 the final iid performance of models [Wu et al., 2021]. In contrast, our work focuses on the *relevance*
268 of the data with respect to the target task, where it is well known that changing data order improves
269 performance (exemplified by fine-tuning). This allows there to be much larger gains in performance.

270 **Necessity of pretraining.** Many prior works argue that it is necessary to incorporate target skills
271 during pretraining. For example, [Allen-Zhu and Li, 2024, Jiang et al., 2024b] argue that instruction-
272 tuning data needs to be seen during pretraining. Moreover, as discussed in 4.2, many practitioners
273 pretrain language models from scratch with the belief that it is necessary to see this data during
274 pretraining. Our work shows that this might not be the case: for some tasks, data might not need to
275 be seen during pretraining, as long as one follows optimal training procedures for adaptation.

276 6 Discussion

277 6.1 Hypotheses for inefficiency of fine-tuning

278 In this paper, we thoroughly characterize the failure of fine-tuning. We share two hypotheses for
279 why fine-tuning underperforms replay data. We first identify a training instability that occurs for a
280 few steps of fine-tuning which replay slightly mitigates. We provide experiments and discussion in
281 Appendix B.1. However, even with perfect optimization, we identify a statistical barrier to leveraging
282 rare samples due to a natural tendency to overfit to small samples. In Appendix B.2, we detail a toy
283 model where failure arises due to a small number of noisy data points, leveraging classical intuition
284 from the double-descent literature. Another reasonable hypothesis is that this is an artifact of our
285 current training scale. To control for this, we ablate model size in Appendix B.3. Interestingly, the
286 benefit of replay seems to persist across model scale, decreasing support for this hypothesis.

287 6.2 Limitations

288 Our controlled setting is different from real-world settings in possibly important ways. For example,
289 we assume two distributions, whereas in practice, pre-training is a multi-task learning problem with
290 much higher diversity. Our space of data schedules also precludes us from studying more complicated
291 methods such as continuous annealing of the distribution, sample-level orderings, and more advanced
292 fine-tuning methods. There might be other important hyperparameters and joint interactions that
293 we did not explore, possibly affecting the robustness of our results. Many properties that we care
294 about for language models might not perfectly correlate with validation loss, which requires using
295 different metrics under our framework. In practice, replay comes at the cost of increased compute,
296 which might be a limiting factor outside of standard fine-tuning settings.

297 6.3 Broader Impacts

298 We hope our work can be used to improve the data efficiency of language models, especially for low
299 resource domains that otherwise receive little attention. We acknowledge our work may increase the
300 compute used in language model training. We believe most other harms associated with our work are
301 generally applicable to most language modeling research.

302 **References**

- 303 E. Abbe, E. Cornacchia, and A. Lotfi. Provable advantage of curriculum learning on parity targets
304 with mixed inputs, 2023. URL <https://arxiv.org/abs/2306.16921>.
- 305 L. B. Allal, A. Lozhkov, E. Bakouch, G. M. Blázquez, G. Penedo, L. Tunstall, A. Marafioti,
306 H. Kydlíček, A. P. Lajarín, V. Srivastav, J. Lochner, C. Fahlgren, X.-S. Nguyen, C. Fourrier,
307 B. Burtenshaw, H. Larcher, H. Zhao, C. Zakka, M. Morlon, C. Raffel, L. von Werra, and T. Wolf.
308 Smollm2: When smol goes big – data-centric training of a small language model, 2025. URL
309 <https://arxiv.org/abs/2502.02737>.
- 310 Z. Allen-Zhu and Y. Li. Physics of language models: Part 3.1, knowledge storage and extraction,
311 2024. URL <https://arxiv.org/abs/2309.14316>.
- 312 V. Aryabumi, Y. Su, R. Ma, A. Morisot, I. Zhang, A. Locatelli, M. Fadaee, A. Üstün, and S. Hooker.
313 To code, or not to code? exploring impact of code in pre-training, 2024. URL <https://arxiv.org/abs/2408.10914>.
- 315 I. Baucells, J. Aula-Blasco, I. de Dios-Flores, S. Paniagua Suárez, N. Perez, A. Salles, S. Sotelo Docio,
316 J. Falcão, J. J. Saiz, R. Sepulveda Torres, J. Barnes, P. Gamallo, A. Gonzalez-Agirre, G. Rigau,
317 and M. Villegas. IberoBench: A benchmark for LLM evaluation in Iberian languages. In
318 O. Rambow, L. Wanner, M. Apidianaki, H. Al-Khalifa, B. D. Eugenio, and S. Schockaert, editors,
319 *Proceedings of the 31st International Conference on Computational Linguistics*, pages 10491–
320 10519, Abu Dhabi, UAE, Jan. 2025. Association for Computational Linguistics. URL <https://aclanthology.org/2025.coling-main.699/>.
- 322 M. Belkin, D. Hsu, S. Ma, and S. Mandal. Reconciling modern machine-learning practice and
323 the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, 116
324 (32):15849–15854, July 2019. ISSN 1091-6490. doi: 10.1073/pnas.1903070116. URL <http://dx.doi.org/10.1073/pnas.1903070116>.
- 326 Y. Bengio, J. Louradour, R. Collobert, and J. Weston. Curriculum learning. In *Proceedings of the
327 26th Annual International Conference on Machine Learning*, ICML ’09, page 41–48, New
328 York, NY, USA, 2009. Association for Computing Machinery. ISBN 9781605585161. doi:
329 10.1145/1553374.1553380. URL <https://doi.org/10.1145/1553374.1553380>.
- 330 C. Blakeney, M. Paul, B. W. Larsen, S. Owen, and J. Frankle. Does your data spark joy? performance
331 gains from domain upsampling at the end of training, 2024. URL <https://arxiv.org/abs/2406.03476>.
- 333 H. Chang, J. Park, S. Ye, S. Yang, Y. Seo, D.-S. Chang, and M. Seo. How do large language models
334 acquire factual knowledge during pretraining?, 2024. URL <https://arxiv.org/abs/2406.11813>.
- 336 H. Chen, J. Geng, A. Bhaskar, D. Friedman, and D. Chen. Continual memorization of factoids in
337 language models, 2025a. URL <https://arxiv.org/abs/2411.07175>.
- 338 M. F. Chen, N. Roberts, K. Bhatia, J. Wang, C. Zhang, F. Sala, and C. Ré. Skill-it! a data-
339 driven skills framework for understanding and training language models, 2023. URL <https://arxiv.org/abs/2307.14430>.
- 341 M. F. Chen, M. Y. Hu, N. Lourie, K. Cho, and C. Ré. Aioli: A unified optimization framework for
342 language model data mixing, 2025b. URL <https://arxiv.org/abs/2411.05735>.
- 343 J. M. Cohen, S. Kaur, Y. Li, J. Z. Kolter, and A. Talwalkar. Gradient descent on neural networks
344 typically occurs at the edge of stability, 2022. URL <https://arxiv.org/abs/2103.00065>.
- 345 J. M. Cohen, A. Damian, A. Talwalkar, Z. Kolter, and J. D. Lee. Understanding optimization in deep
346 learning with central flows, 2024. URL <https://arxiv.org/abs/2410.24206>.
- 347 X. Dang, C. Baek, K. Wen, Z. Kolter, and A. Raghunathan. Weight ensembling improves reasoning
348 in language models, 2025. URL <https://arxiv.org/abs/2504.10478>.

- 349 A. Defazio, A. Cutkosky, H. Mehta, and K. Mishchenko. Optimal linear decay learning rate schedules
 350 and further refinements, 2024. URL <https://arxiv.org/abs/2310.07831>.
- 351 J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional
 352 transformers for language understanding, 2019. URL <https://arxiv.org/abs/1810.04805>.
- 353 N. Ding, Y. Chen, B. Xu, Y. Qin, Z. Zheng, S. Hu, Z. Liu, M. Sun, and B. Zhou. Enhancing chat
 354 language models by scaling high-quality instructional conversations, 2023.
- 355 J. Etxaniz, O. Sainz, N. Perez, I. Aldabe, G. Rigau, E. Agirre, A. Ormazabal, M. Artetxe, and
 356 A. Soroa. Latxa: An open language model and evaluation suite for basque, 2024. URL <https://arxiv.org/abs/2403.20266>.
- 358 K. Everett, L. Xiao, M. Wortsman, A. A. Alemi, R. Novak, P. J. Liu, I. Gur, J. Sohl-Dickstein, L. P.
 359 Kaelbling, J. Lee, and J. Pennington. Scaling exponents across parameterizations and optimizers,
 360 2024. URL <https://arxiv.org/abs/2407.05872>.
- 361 S. Fan and M. Jaggi. Irreducible curriculum for language model pretraining, 2023. URL <https://arxiv.org/abs/2310.15389>.
- 363 S. Fan, M. Pagliardini, and M. Jaggi. Doge: Domain reweighting with generalization estimation,
 364 2024. URL <https://arxiv.org/abs/2310.15393>.
- 365 L. Gao, J. Tow, B. Abbasi, S. Biderman, S. Black, A. DiPofi, C. Foster, L. Golding, J. Hsu,
 366 A. Le Noac'h, H. Li, K. McDonell, N. Muennighoff, C. Ociepa, J. Phang, L. Reynolds,
 367 H. Schoelkopf, A. Skowron, L. Sutawika, E. Tang, A. Thite, B. Wang, K. Wang, and A. Zou. The
 368 language model evaluation harness, 07 2024. URL <https://zenodo.org/records/12608602>.
- 369 Z. Gekhman, G. Yona, R. Aharoni, M. Eyal, A. Feder, R. Reichart, and J. Herzig. Does fine-tuning
 370 llms on new knowledge encourage hallucinations?, 2024. URL <https://arxiv.org/abs/2405.05904>.
- 372 G. Ghosal, T. Hashimoto, and A. Raghunathan. Understanding finetuning for factual knowledge
 373 extraction, 2024. URL <https://arxiv.org/abs/2406.14785>.
- 374 A. Gordon, Z. Kozareva, and M. Roemmele. SemEval-2012 task 7: Choice of plausible alternatives:
 375 An evaluation of commonsense causal reasoning. In E. Agirre, J. Bos, M. Diab, S. Manandhar,
 376 Y. Marton, and D. Yuret, editors, **SEM 2012: The First Joint Conference on Lexical and Computational
 377 Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume
 378 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pages
 379 394–398, Montréal, Canada, 7–8 June 2012. Association for Computational Linguistics. URL
 380 <https://aclanthology.org/S12-1052/>.
- 381 S. Goyal, P. Maini, Z. C. Lipton, A. Raghunathan, and J. Z. Kolter. Scaling laws for data filtering –
 382 data curation cannot be compute agnostic, 2024. URL <https://arxiv.org/abs/2404.07177>.
- 383 A. Grattafiori, A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur,
 384 A. Schelten, A. Vaughan, A. Yang, A. Fan, A. Goyal, A. Hartshorn, A. Yang, A. Mitra, A. Sravankumar,
 385 A. Korenev, A. Hinsvark, A. Rao, A. Zhang, A. Rodriguez, A. Gregerson, A. Spataru,
 386 B. Roziere, B. Biron, B. Tang, B. Chern, C. Caucheteux, C. Nayak, C. Bi, C. Marra, C. McConnell,
 387 C. Keller, C. Touret, C. Wu, C. Wong, C. C. Ferrer, C. Nikolaidis, D. Allonsius, D. Song, D. Pintz,
 388 D. Livshits, D. Wyatt, D. Esiobu, D. Choudhary, D. Mahajan, D. Garcia-Olano, D. Perino, D. Huppenkes,
 389 E. Lakomkin, E. AlBadawy, E. Lobanova, E. Dinan, E. M. Smith, F. Radenovic, F. Guzmán,
 390 F. Zhang, G. Synnaeve, G. Lee, G. L. Anderson, G. Thattai, G. Nail, G. Mialon, G. Pang, G. Curell,
 391 H. Nguyen, H. Korevaar, H. Xu, H. Touvron, I. Zarov, I. A. Ibarra, I. Kloumann, I. Misra,
 392 I. Evtimov, J. Zhang, J. Copet, J. Lee, J. Geffert, J. Vranes, J. Park, J. Mahadeokar, J. Shah,
 393 J. van der Linde, J. Billock, J. Hong, J. Lee, J. Fu, J. Chi, J. Huang, J. Liu, J. Wang, J. Yu, J. Bitton,
 394 J. Spisak, J. Park, J. Rocca, J. Johnstun, J. Saxe, J. Jia, K. V. Alwala, K. Prasad, K. Upasani, K. Plawiak,
 395 K. Li, K. Heafield, K. Stone, K. El-Arini, K. Iyer, K. Malik, K. Chiu, K. Bhalla, K. Lakhotia,
 396 L. Rantala-Yearly, L. van der Maaten, L. Chen, L. Tan, L. Jenkins, L. Martin, L. Madaan, L. Malo,
 397 L. Blecher, L. Landzaat, L. de Oliveira, M. Muzzi, M. Pasupuleti, M. Singh, M. Paluri, M. Karandas,
 398 M. Tsimpoukelli, M. Oldham, M. Rita, M. Pavlova, M. Kambadur, M. Lewis, M. Si, M. K. Singh,
 399 M. Hassan, N. Goyal, N. Torabi, N. Bashlykov, N. Chatterji, N. Zhang,

- 400 O. Duchenne, O. Çelebi, P. Alrassy, P. Zhang, P. Li, P. Vasic, P. Weng, P. Bhargava, P. Dubal, P. Krishnan, P. S. Koura, P. Xu, Q. He, Q. Dong, R. Srinivasan, R. Ganapathy, R. Calderer, R. S. Cabral, R. Stojnic, R. Raileanu, R. Maheswari, R. Girdhar, R. Patel, R. Sauvestre, R. Polidoro, R. Sumbaly, R. Taylor, R. Silva, R. Hou, R. Wang, S. Hosseini, S. Chennabasappa, S. Singh, S. Bell, S. S. Kim, S. Edunov, S. Nie, S. Narang, S. Raparthy, S. Shen, S. Wan, S. Bhosale, S. Zhang, S. Vandenhende, S. Batra, S. Whitman, S. Sootla, S. Collot, S. Gururangan, S. Borodinsky, T. Herman, T. Fowler, T. Sheasha, T. Georgiou, T. Scialom, T. Speckbacher, T. Mihaylov, T. Xiao, U. Karn, V. Goswami, V. Gupta, V. Ramanathan, V. Kerkez, V. Gonguet, V. Do, V. Vogeti, V. Albiero, V. Petrovic, W. Chu, W. Xiong, W. Fu, W. Meers, X. Martinet, X. Wang, X. Wang, X. E. Tan, X. Xia, X. Xie, X. Jia, X. Wang, Y. Goldschlag, Y. Gaur, Y. Babaei, Y. Wen, Y. Song, Y. Zhang, Y. Li, Y. Mao, Z. D. Coudert, Z. Yan, Z. Chen, Z. Papakipos, A. Singh, A. Srivastava, A. Jain, A. Kelsey, A. Shajnfeld, A. Gangidi, A. Victoria, A. Goldstand, A. Menon, A. Sharma, A. Boesenber, A. Baevski, A. Feinstein, A. Kallet, A. Sangani, A. Teo, A. Yunus, A. Lupu, A. Alvarado, A. Caples, A. Gu, A. Ho, A. Poulton, A. Ryan, A. Ramchandani, A. Dong, A. Franco, A. Goyal, A. Saraf, A. Chowdhury, A. Gabriel, A. Bharambe, A. Eisenman, A. Yazdan, B. James, B. Maurer, B. Leonhardi, B. Huang, B. Loyd, B. D. Paola, B. Paranjape, B. Liu, B. Wu, B. Ni, B. Hancock, B. Wasti, B. Spence, B. Stojkovic, B. Gamido, B. Montalvo, C. Parker, C. Burton, C. Mejia, C. Liu, C. Wang, C. Kim, C. Zhou, C. Hu, C.-H. Chu, C. Cai, C. Tindal, C. Feichtenhofer, C. Gao, D. Civin, D. Beaty, D. Kreymer, D. Li, D. Adkins, D. Xu, D. Testuggine, D. David, D. Parikh, D. Liskovich, D. Foss, D. Wang, D. Le, D. Holland, E. Dowling, E. Jamil, E. Montgomery, E. Presani, E. Hahn, E. Wood, E.-T. Le, E. Brinkman, E. Arcuate, E. Dunbar, E. Smothers, F. Sun, F. Kreuk, F. Tian, F. Kokkinos, F. Ozgenel, F. Caggioni, F. Kanayet, F. Seide, G. M. Florez, G. Schwarz, G. Badeer, G. Swee, G. Halpern, G. Herman, G. Sizov, Guangyi, Zhang, G. Lakshminarayanan, H. Inan, H. Shojanazeri, H. Zou, H. Wang, H. Zha, H. Habeeb, H. Rudolph, H. Suk, H. Aspegren, H. Goldman, H. Zhan, I. Damlaj, I. Molybog, I. Tufanov, I. Leontiadis, I.-E. Veliche, I. Gat, J. Weissman, J. Geboski, J. Kohli, J. Lam, J. Asher, J.-B. Gaya, J. Marcus, J. Tang, J. Chan, J. Zhen, J. Reizenstein, J. Teboul, J. Zhong, J. Jin, J. Yang, J. Cummings, J. Carvill, J. Shepard, J. McPhie, J. Torres, J. Ginsburg, J. Wang, K. Wu, K. H. U, K. Saxena, K. Khandelwal, K. Zand, K. Matosich, K. Veeraraghavan, K. Michelena, K. Li, K. Jagadeesh, K. Huang, K. Chawla, K. Huang, L. Chen, L. Garg, L. A, L. Silva, L. Bell, L. Zhang, L. Guo, L. Yu, L. Moshkovich, L. Wehrstedt, M. Khabsa, M. Avalani, M. Bhatt, M. Mankus, M. Hasson, M. Lennie, M. Reso, M. Groshev, M. Naumov, M. Lathi, M. Keneally, M. Liu, M. L. Seltzer, M. Valko, M. Restrepo, M. Patel, M. Vyatskov, M. Samvelyan, M. Clark, M. Macey, M. Wang, M. J. Hermoso, M. Metanat, M. Rastegari, M. Bansal, N. Santhanam, N. Parks, N. White, N. Bawa, N. Singhal, N. Egebo, N. Usunier, N. Mehta, N. P. Laptev, N. Dong, N. Cheng, O. Chernoguz, O. Hart, O. Salpekar, O. Kalinli, P. Kent, P. Parekh, P. Saab, P. Balaji, P. Rittner, P. Bontrager, P. Roux, P. Dollar, P. Zvyagina, P. Ratanchandani, P. Yuvraj, Q. Liang, R. Alao, R. Rodriguez, R. Ayub, R. Murthy, R. Nayani, R. Mitra, R. Parthasarathy, R. Li, R. Hogan, R. Battey, R. Wang, R. Howes, R. Rinott, S. Mehta, S. Siby, S. J. Bondu, S. Datta, S. Chugh, S. Hunt, S. Dhillon, S. Sidorov, S. Pan, S. Mahajan, S. Verma, S. Yamamoto, S. Ramaswamy, S. Lindsay, S. Lindsay, S. Feng, S. Lin, S. C. Zha, S. Patil, S. Shankar, S. Zhang, S. Zhang, S. Wang, S. Agarwal, S. Sajuyigbe, S. Chintala, S. Max, S. Chen, S. Kehoe, S. Satterfield, S. Govindaprasad, S. Gupta, S. Deng, S. Cho, S. Virk, S. Subramanian, S. Choudhury, S. Goldman, T. Remez, T. Glaser, T. Best, T. Koehler, T. Robinson, T. Li, T. Zhang, T. Matthews, T. Chou, T. Shaked, V. Vontimitta, V. Ajayi, V. Montanez, V. Mohan, V. S. Kumar, V. Mangla, V. Ionescu, V. Poenaru, V. T. Mihailescu, V. Ivanov, W. Li, W. Wang, W. Jiang, W. Bouaziz, W. Constable, X. Tang, X. Wu, X. Wang, X. Wu, X. Gao, Y. Kleinman, Y. Chen, Y. Hu, Y. Jia, Y. Qi, Y. Li, Y. Zhang, Y. Zhang, Y. Adi, Y. Nam, Y. Yu, Y. Wang, Y. Zhao, Y. Hao, Y. Qian, Y. Li, Y. He, Z. Rait, Z. DeVito, Z. Rosnbrick, Z. Wen, Z. Yang, Z. Zhao, and Z. Ma. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- 449 K. Gupta, B. Thérien, A. Ibrahim, M. L. Richter, Q. Anthony, E. Belilovsky, I. Rish, and T. Lesort.
 450 Continual pre-training of large language models: How to (re)warm your model?, 2023. URL
 451 <https://arxiv.org/abs/2308.04014>.
- 452 T. Hastie, A. Montanari, S. Rosset, and R. J. Tibshirani. Surprises in high-dimensional ridgeless least
 453 squares interpolation, 2020. URL <https://arxiv.org/abs/1903.08560>.
- 454 J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. de Las Casas,
 455 L. A. Hendricks, J. Welbl, A. Clark, T. Hennigan, E. Noland, K. Millican, G. van den Driessche,
 456 B. Damoc, A. Guy, S. Osindero, K. Simonyan, E. Elsen, J. W. Rae, O. Vinyals, and L. Sifre.

- 457 Training compute-optimal large language models, 2022. URL <https://arxiv.org/abs/2203.15556>.
- 459 S. Hu, Y. Tu, X. Han, C. He, G. Cui, X. Long, Z. Zheng, Y. Fang, Y. Huang, W. Zhao, X. Zhang, Z. L.
460 Thai, K. Zhang, C. Wang, Y. Yao, C. Zhao, J. Zhou, J. Cai, Z. Zhai, N. Ding, C. Jia, G. Zeng, D. Li,
461 Z. Liu, and M. Sun. Minicpm: Unveiling the potential of small language models with scalable
462 training strategies, 2024. URL <https://arxiv.org/abs/2404.06395>.
- 463 A. Ibrahim, B. Thérien, K. Gupta, M. L. Richter, Q. Anthony, T. Lesort, E. Belilovsky, and I. Rish.
464 Simple and scalable strategies to continually pre-train large language models, 2024. URL <https://arxiv.org/abs/2403.08763>.
- 466 G. Ilharco, M. T. Ribeiro, M. Wortsman, S. Gururangan, L. Schmidt, H. Hajishirzi, and A. Farhadi.
467 Editing models with task arithmetic, 2023. URL <https://arxiv.org/abs/2212.04089>.
- 468 Y. Jiang, A. Zhou, Z. Feng, S. Malladi, and J. Z. Kolter. Adaptive data optimization: Dynamic sample
469 selection with scaling laws, 2024a. URL <https://arxiv.org/abs/2410.11820>.
- 470 Z. Jiang, Z. Sun, W. Shi, P. Rodriguez, C. Zhou, G. Neubig, X. V. Lin, W. tau Yih, and S. Iyer.
471 Instruction-tuned language models are better knowledge learners, 2024b. URL <https://arxiv.org/abs/2402.12847>.
- 473 J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan,
474 T. Ramalho, A. Grabska-Barwinska, D. Hassabis, C. Clopath, D. Kumaran, and R. Hadsell.
475 Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of
476 Sciences*, 114(13):3521–3526, Mar. 2017. ISSN 1091-6490. doi: 10.1073/pnas.1611835114. URL
477 <http://dx.doi.org/10.1073/pnas.1611835114>.
- 478 S. Kotha, J. M. Springer, and A. Raghunathan. Understanding catastrophic forgetting in language
479 models via implicit inference, 2024. URL <https://arxiv.org/abs/2309.10105>.
- 480 A. Kumar, A. Raghunathan, R. Jones, T. Ma, and P. Liang. Fine-tuning can distort pretrained features
481 and underperform out-of-distribution, 2022. URL <https://arxiv.org/abs/2202.10054>.
- 482 J. Li, A. Fang, G. Smyrnis, M. Ivgi, M. Jordan, S. Gadre, H. Bansal, E. Guha, S. Keh, K. Arora,
483 S. Garg, R. Xin, N. Muennighoff, R. Heckel, J. Mercat, M. Chen, S. Gururangan, M. Wortsman,
484 A. Albalak, Y. Bitton, M. Nezhurina, A. Abbas, C.-Y. Hsieh, D. Ghosh, J. Gardner, M. Kilian,
485 H. Zhang, R. Shao, S. Pratt, S. Sanyal, G. Ilharco, G. Daras, K. Marathe, A. Gokaslan, J. Zhang,
486 K. Chandu, T. Nguyen, I. Vasiljevic, S. Kakade, S. Song, S. Sanghavi, F. Faghri, S. Oh, L. Zettlemoyer,
487 K. Lo, A. El-Nouby, H. Pouransari, A. Toshev, S. Wang, D. Groeneveld, L. Soldaini,
488 P. W. Koh, J. Jitsev, T. Kollar, A. G. Dimakis, Y. Carmon, A. Dave, L. Schmidt, and V. Shankar.
489 Datacomp-lm: In search of the next generation of training sets for language models, 2025. URL
490 <https://arxiv.org/abs/2406.11794>.
- 491 R. Li, L. B. Allal, Y. Zi, N. Muennighoff, D. Kocetkov, C. Mou, M. Marone, C. Akiki, J. Li, J. Chim,
492 Q. Liu, E. Zheltonozhskii, T. Y. Zhuo, T. Wang, O. Dehaene, M. Davaadorj, J. Lamy-Poirier,
493 J. Monteiro, O. Shliazhko, N. Gontier, N. Meade, A. Zebaze, M.-H. Yee, L. K. Umapathi, J. Zhu,
494 B. Lipkin, M. Oblokulov, Z. Wang, R. Murthy, J. Stillerman, S. S. Patel, D. Abulkhanov, M. Zocca,
495 M. Dey, Z. Zhang, N. Fahmy, U. Bhattacharyya, W. Yu, S. Singh, S. Luccioni, P. Villegas,
496 M. Kunakov, F. Zhdanov, M. Romero, T. Lee, N. Timor, J. Ding, C. Schlesinger, H. Schoelkopf,
497 J. Ebert, T. Dao, M. Mishra, A. Gu, J. Robinson, C. J. Anderson, B. Dolan-Gavitt, D. Contractor,
498 S. Reddy, D. Fried, D. Bahdanau, Y. Jernite, C. M. Ferrandis, S. Hughes, T. Wolf, A. Guha,
499 L. von Werra, and H. de Vries. Starcoder: may the source be with you!, 2023. URL <https://arxiv.org/abs/2305.06161>.
- 500 Z. Lin, Z. Gou, Y. Gong, X. Liu, Y. Shen, R. Xu, C. Lin, Y. Yang, J. Jiao, N. Duan, and W. Chen.
501 Rho-1: Not all tokens are what you need, 2025. URL <https://arxiv.org/abs/2404.07965>.
- 503 S. Longpre, L. Hou, T. Vu, A. Webson, H. W. Chung, Y. Tay, D. Zhou, Q. V. Le, B. Zoph, J. Wei, and
504 A. Roberts. The flan collection: Designing data and methods for effective instruction tuning, 2023.
505 URL <https://arxiv.org/abs/2301.13688>.

- 506 X. H. Lù, Z. Kasner, and S. Reddy. Weblinx: Real-world website navigation with multi-turn dialogue.
 507 *arXiv preprint arXiv:2402.05930*, 2024.
- 508 K. Meng, D. Bau, A. Andonian, and Y. Belinkov. Locating and editing factual associations in gpt,
 509 2023. URL <https://arxiv.org/abs/2202.05262>.
- 510 S. Mindermann, J. Brauner, M. Razzak, M. Sharma, A. Kirsch, W. Xu, B. Höltgen, A. N. Gomez,
 511 A. Morisot, S. Farquhar, and Y. Gal. Prioritized training on points that are learnable, worth learning,
 512 and not yet learnt, 2022. URL <https://arxiv.org/abs/2206.07137>.
- 513 N. Muennighoff, A. M. Rush, B. Barak, T. L. Scao, A. Piktus, N. Tazi, S. Pyysalo, T. Wolf, and
 514 C. Raffel. Scaling data-constrained language models, 2023. URL <https://arxiv.org/abs/2305.16264>.
- 515 Nvidia, :, B. Adler, N. Agarwal, A. Aithal, D. H. Anh, P. Bhattacharya, A. Brundyn, J. Casper,
 516 B. Catanzaro, S. Clay, J. Cohen, S. Das, A. Dattagupta, O. Delalleau, L. Derczynski, Y. Dong,
 517 D. Egert, E. Evans, A. Ficek, D. Fridman, S. Ghosh, B. Ginsburg, I. Gitman, T. Grzegorzek, R. Hero,
 518 J. Huang, V. Jawa, J. Jennings, A. Jhunjhunwala, J. Kamalu, S. Khan, O. Kuchaiev, P. LeGresley,
 519 H. Li, J. Liu, Z. Liu, E. Long, A. S. Mahabaleshwarkar, S. Majumdar, J. Maki, M. Martinez,
 520 M. R. de Melo, I. Moshkov, D. Narayanan, S. Narenthiran, J. Navarro, P. Nguyen, O. Nitski,
 521 V. Noroozi, G. Nutheti, C. Parisien, J. Parmar, M. Patwary, K. Pawelec, W. Ping, S. Prabhumoye,
 522 R. Roy, T. Saar, V. R. N. Sabavat, S. Satheesh, J. P. Scowcroft, J. Sewall, P. Shamis, G. Shen,
 523 M. Shoeybi, D. Sizer, M. Smelyanskiy, F. Soares, M. N. Sreedhar, D. Su, S. Subramanian, S. Sun,
 524 S. Toshniwal, H. Wang, Z. Wang, J. You, J. Zeng, J. Zhang, J. Zhang, V. Zhang, Y. Zhang, and
 525 C. Zhu. Nemotron-4 340b technical report, 2024. URL <https://arxiv.org/abs/2406.11704>.
- 526 T. OLMo, P. Walsh, L. Soldaini, D. Groeneveld, K. Lo, S. Arora, A. Bhagia, Y. Gu, S. Huang,
 527 M. Jordan, N. Lambert, D. Schwenk, O. Tafjord, T. Anderson, D. Atkinson, F. Brahman, C. Clark,
 528 P. Dasigi, N. Dziri, M. Guerquin, H. Ivison, P. W. Koh, J. Liu, S. Malik, W. Merrill, L. J. V.
 529 Miranda, J. Morrison, T. Murray, C. Nam, V. Pyatkin, A. Rangapur, M. Schmitz, S. Skjonsberg,
 530 D. Wadden, C. Wilhelm, M. Wilson, L. Zettlemoyer, A. Farhadi, N. A. Smith, and H. Hajishirzi. 2
 531 olmo 2 furious, 2025. URL <https://arxiv.org/abs/2501.00656>.
- 532 A. Panigrahi, B. Liu, S. Malladi, A. Risteski, and S. Goel. Progressive distillation induces an implicit
 533 curriculum, 2024. URL <https://arxiv.org/abs/2410.05464>.
- 534 G. I. Parisi, R. Kemker, J. L. Part, C. Kanan, and S. Wermter. Continual lifelong learning with neural
 535 networks: A review. *Neural networks*, 113:54–71, 2019.
- 536 J. Phang, T. Févry, and S. R. Bowman. Sentence encoders on stilts: Supplementary training on
 537 intermediate labeled-data tasks, 2019. URL <https://arxiv.org/abs/1811.01088>.
- 538 E. M. Ponti, G. Glavaš, O. Majewska, Q. Liu, I. Vulić, and A. Korhonen. Xcopa: A multilingual
 539 dataset for causal commonsense reasoning, 2020. URL <https://arxiv.org/abs/2005.00333>.
- 540 A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever. Improving language understanding by
 541 generative pre-training. 2018.
- 542 C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu.
 543 Exploring the limits of transfer learning with a unified text-to-text transformer, 2023. URL
 544 <https://arxiv.org/abs/1910.10683>.
- 545 D. Rolnick, A. Ahuja, J. Schwarz, T. P. Lillicrap, and G. Wayne. Experience replay for continual
 546 learning, 2019. URL <https://arxiv.org/abs/1811.11682>.
- 547 F. Schaipp, A. Hägele, A. Taylor, U. Simsekli, and F. Bach. The surprising agreement between
 548 convex optimization theory and learning-rate scheduling for large model training, 2025. URL
 549 <https://arxiv.org/abs/2501.18965>.
- 550 D. Soboleva, F. Al-Khateeb, R. Myers, J. R. Steeves, J. Hestness, and N. Dey. SlimPajama: A 627B
 551 token cleaned and deduplicated version of RedPajama. <https://www.cerebras.net/blog/slimpajama-a-627b-token-cleaned-and-deduplicated-version-of-redpajama>, 06
 552 2023. URL <https://huggingface.co/datasets/cerebras/SlimPajama-627B>.

- 555 J. M. Springer, S. Goyal, K. Wen, T. Kumar, X. Yue, S. Malladi, G. Neubig, and A. Raghunathan.
556 Overtrained language models are harder to fine-tune, 2025. URL <https://arxiv.org/abs/2503.19206>.
- 558 Teknium. Openhermes 2.5: An open dataset of synthetic data for generalist llm assistants, 2023.
559 URL <https://huggingface.co/datasets/teknium/OpenHermes-2.5>.
- 560 M. Weber, D. Fu, Q. Anthony, Y. Oren, S. Adams, A. Alexandrov, X. Lyu, H. Nguyen, X. Yao,
561 V. Adams, B. Athiwaratkun, R. Chalamala, K. Chen, M. Ryabinin, T. Dao, P. Liang, C. Ré, I. Rish,
562 and C. Zhang. Redpajama: an open dataset for training large language models, 2024. URL
563 <https://arxiv.org/abs/2411.12372>.
- 564 K. Wen, Z. Li, J. Wang, D. Hall, P. Liang, and T. Ma. Understanding warmup-stable-decay learning
565 rates: A river valley loss landscape perspective, 2024. URL <https://arxiv.org/abs/2410.05192>.
- 567 M. Wortsman, G. Ilharco, J. W. Kim, M. Li, S. Kornblith, R. Roelofs, R. Gontijo-Lopes, H. Hajishirzi,
568 A. Farhadi, H. Namkoong, and L. Schmidt. Robust fine-tuning of zero-shot models, 2022. URL
569 <https://arxiv.org/abs/2109.01903>.
- 570 X. Wu, E. Dyer, and B. Neyshabur. When do curricula work?, 2021. URL <https://arxiv.org/abs/2012.03107>.
- 572 S. M. Xie, H. Pham, X. Dong, N. Du, H. Liu, Y. Lu, P. Liang, Q. V. Le, T. Ma, and A. W. Yu.
573 Doremi: Optimizing data mixtures speeds up language model pretraining, 2023. URL <https://arxiv.org/abs/2305.10429>.
- 575 G. Yang, E. J. Hu, I. Babuschkin, S. Sidor, X. Liu, D. Farhi, N. Ryder, J. Pachocki, W. Chen, and
576 J. Gao. Tensor programs v: Tuning large neural networks via zero-shot hyperparameter transfer,
577 2022. URL <https://arxiv.org/abs/2203.03466>.
- 578 G. Yang, J. B. Simon, and J. Bernstein. A spectral condition for feature learning, 2024a. URL
579 <https://arxiv.org/abs/2310.17813>.
- 580 Z. Yang, N. Band, S. Li, E. Candès, and T. Hashimoto. Synthetic continued pretraining, 2024b. URL
581 <https://arxiv.org/abs/2409.07431>.
- 582 T. Zhang, F. Wu, A. Katiyar, K. Q. Weinberger, and Y. Artzi. Revisiting few-sample bert fine-tuning,
583 2021. URL <https://arxiv.org/abs/2006.05987>.
- 584 Çağatay Yıldız, N. K. Ravichandran, N. Sharma, M. Bethge, and B. Ermis. Investigating continual
585 pretraining in large language models: Insights and implications, 2025. URL <https://arxiv.org/abs/2402.17400>.

587 **NeurIPS Paper Checklist**

588 The checklist is designed to encourage best practices for responsible machine learning research,
589 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove
590 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should
591 follow the references and follow the (optional) supplemental material. The checklist does NOT count
592 towards the page limit.

593 Please read the checklist guidelines carefully for information on how to answer these questions. For
594 each question in the checklist:

- 595 • You should answer [Yes] , [No] , or [NA] .
596 • [NA] means either that the question is Not Applicable for that particular paper or the
597 relevant information is Not Available.
598 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

599 **The checklist answers are an integral part of your paper submission.** They are visible to the
600 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it
601 (after eventual revisions) with the final version of your paper, and its final version will be published
602 with the paper.

603 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.
604 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a
605 proper justification is given (e.g., "error bars are not reported because it would be too computationally
606 expensive" or "we were unable to find the license for the dataset we used"). In general, answering
607 "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we
608 acknowledge that the true answer is often more nuanced, so please just use your best judgment and
609 write a justification to elaborate. All supporting evidence can appear either in the main paper or the
610 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification
611 please point to the section(s) where related material for the question can be found.

612 **IMPORTANT**, please:

- 613 • **Delete this instruction block, but keep the section heading “NeurIPS paper checklist”,**
614 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
615 • **Do not modify the questions and only use the provided macros for your answers.**

616 **1. Claims**

617 Question: Do the main claims made in the abstract and introduction accurately reflect the
618 paper's contributions and scope?

619 Answer: [Yes]

620 Justification: Our abstract and introduction makes claims concretely supported by our
621 experiments (refer to abstract and introduction).

622 Guidelines:

- 623 • The answer NA means that the abstract and introduction do not include the claims
624 made in the paper.
625 • The abstract and/or introduction should clearly state the claims made, including the
626 contributions made in the paper and important assumptions and limitations. A No or
627 NA answer to this question will not be perceived well by the reviewers.
628 • The claims made should match theoretical and experimental results, and reflect how
629 much the results can be expected to generalize to other settings.
630 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
631 are not attained by the paper.

632 **2. Limitations**

633 Question: Does the paper discuss the limitations of the work performed by the authors?

634 Answer: [Yes]

635 Justification: We discuss the limitations of our work in the limitations section (refer to
636 limitations).

637 Guidelines:

- 638 • The answer NA means that the paper has no limitation while the answer No means that
639 the paper has limitations, but those are not discussed in the paper.
- 640 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 641 • The paper should point out any strong assumptions and how robust the results are to
642 violations of these assumptions (e.g., independence assumptions, noiseless settings,
643 model well-specification, asymptotic approximations only holding locally). The authors
644 should reflect on how these assumptions might be violated in practice and what the
645 implications would be.
- 646 • The authors should reflect on the scope of the claims made, e.g., if the approach was
647 only tested on a few datasets or with a few runs. In general, empirical results often
648 depend on implicit assumptions, which should be articulated.
- 649 • The authors should reflect on the factors that influence the performance of the approach.
650 For example, a facial recognition algorithm may perform poorly when image resolution
651 is low or images are taken in low lighting. Or a speech-to-text system might not be
652 used reliably to provide closed captions for online lectures because it fails to handle
653 technical jargon.
- 654 • The authors should discuss the computational efficiency of the proposed algorithms
655 and how they scale with dataset size.
- 656 • If applicable, the authors should discuss possible limitations of their approach to
657 address problems of privacy and fairness.
- 658 • While the authors might fear that complete honesty about limitations might be used by
659 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
660 limitations that aren't acknowledged in the paper. The authors should use their best
661 judgment and recognize that individual actions in favor of transparency play an impor-
662 tant role in developing norms that preserve the integrity of the community. Reviewers
663 will be specifically instructed to not penalize honesty concerning limitations.

664 3. Theory Assumptions and Proofs

665 Question: For each theoretical result, does the paper provide the full set of assumptions and
666 a complete (and correct) proof?

667 Answer: [NA]

668 Justification: We do not have any theorems.

669 Guidelines:

- 670 • The answer NA means that the paper does not include theoretical results.
- 671 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
672 referenced.
- 673 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 674 • The proofs can either appear in the main paper or the supplemental material, but if
675 they appear in the supplemental material, the authors are encouraged to provide a short
676 proof sketch to provide intuition.
- 677 • Inversely, any informal proof provided in the core of the paper should be complemented
678 by formal proofs provided in appendix or supplemental material.
- 679 • Theorems and Lemmas that the proof relies upon should be properly referenced.

680 4. Experimental Result Reproducibility

681 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
682 perimental results of the paper to the extent that it affects the main claims and/or conclusions
683 of the paper (regardless of whether the code and data are provided or not)?

684 Answer: [Yes]

685 Justification: We provide all the information needed to reproduce the main experimental
686 results of the paper in the main paper (refer to appendix).

687 Guidelines:

- 688 • The answer NA means that the paper does not include experiments.
- 689 • If the paper includes experiments, a No answer to this question will not be perceived
- 690 well by the reviewers: Making the paper reproducible is important, regardless of
- 691 whether the code and data are provided or not.
- 692 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 693 to make their results reproducible or verifiable.
- 694 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 695 For example, if the contribution is a novel architecture, describing the architecture fully
- 696 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 697 be necessary to either make it possible for others to replicate the model with the same
- 698 dataset, or provide access to the model. In general, releasing code and data is often
- 699 one good way to accomplish this, but reproducibility can also be provided via detailed
- 700 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 701 of a large language model), releasing of a model checkpoint, or other means that are
- 702 appropriate to the research performed.
- 703 • While NeurIPS does not require releasing code, the conference does require all submissions
- 704 to provide some reasonable avenue for reproducibility, which may depend on the
- 705 nature of the contribution. For example
- 706 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
- 707 to reproduce that algorithm.
- 708 (b) If the contribution is primarily a new model architecture, the paper should describe
- 709 the architecture clearly and fully.
- 710 (c) If the contribution is a new model (e.g., a large language model), then there should
- 711 either be a way to access this model for reproducing the results or a way to reproduce
- 712 the model (e.g., with an open-source dataset or instructions for how to construct
- 713 the dataset).
- 714 (d) We recognize that reproducibility may be tricky in some cases, in which case
- 715 authors are welcome to describe the particular way they provide for reproducibility.
- 716 In the case of closed-source models, it may be that access to the model is limited in
- 717 some way (e.g., to registered users), but it should be possible for other researchers
- 718 to have some path to reproducing or verifying the results.

719 5. Open access to data and code

720 Question: Does the paper provide open access to the data and code, with sufficient instruc-

721 tions to faithfully reproduce the main experimental results, as described in supplemental

722 material?

723 Answer: [Yes]

724 Justification: We will open source our code for launching runs, as well as logs for the loss

725 values of each run.

726 Guidelines:

- 727 • The answer NA means that paper does not include experiments requiring code.
- 728 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 729 • While we encourage the release of code and data, we understand that this might not be
- 730 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
- 731 including code, unless this is central to the contribution (e.g., for a new open-source
- 732 benchmark).
- 733 • The instructions should contain the exact command and environment needed to run to
- 734 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 735 • The authors should provide instructions on data access and preparation, including how
- 736 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 737 • The authors should provide scripts to reproduce all experimental results for the new
- 738 proposed method and baselines. If only a subset of experiments are reproducible, they
- 739 should state which ones are omitted from the script and why.
- 740

- 742 • At submission time, to preserve anonymity, the authors should release anonymized
743 versions (if applicable).
744 • Providing as much information as possible in supplemental material (appended to the
745 paper) is recommended, but including URLs to data and code is permitted.

746 **6. Experimental Setting/Details**

747 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
748 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
749 results?

750 Answer: [Yes]

751 Justification: We provide all the details necessary to reproduce the main experimental results
752 of the paper in the appendix (refer to appendix).

753 Guidelines:

- 754 • The answer NA means that the paper does not include experiments.
755 • The experimental setting should be presented in the core of the paper to a level of detail
756 that is necessary to appreciate the results and make sense of them.
757 • The full details can be provided either with the code, in appendix, or as supplemental
758 material.

759 **7. Experiment Statistical Significance**

760 Question: Does the paper report error bars suitably and correctly defined or other appropriate
761 information about the statistical significance of the experiments?

762 Answer: [Yes]

763 Justification: Though we can not report error bars due to the high computational cost
764 associated with training language models, we show clean scaling curves to assure that our
765 experiments are statistically significant (refer to plots).

766 Guidelines:

- 767 • The answer NA means that the paper does not include experiments.
768 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
769 dence intervals, or statistical significance tests, at least for the experiments that support
770 the main claims of the paper.
771 • The factors of variability that the error bars are capturing should be clearly stated (for
772 example, train/test split, initialization, random drawing of some parameter, or overall
773 run with given experimental conditions).
774 • The method for calculating the error bars should be explained (closed form formula,
775 call to a library function, bootstrap, etc.).
776 • The assumptions made should be given (e.g., Normally distributed errors).
777 • It should be clear whether the error bar is the standard deviation or the standard error
778 of the mean.
779 • It is OK to report 1-sigma error bars, but one should state it. The authors should
780 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
781 of Normality of errors is not verified.
782 • For asymmetric distributions, the authors should be careful not to show in tables or
783 figures symmetric error bars that would yield results that are out of range (e.g. negative
784 error rates).
785 • If error bars are reported in tables or plots, The authors should explain in the text how
786 they were calculated and reference the corresponding figures or tables in the text.

787 **8. Experiments Compute Resources**

788 Question: For each experiment, does the paper provide sufficient information on the com-
789 puter resources (type of compute workers, memory, time of execution) needed to reproduce
790 the experiments?

791 Answer: [Yes]

792 Justification: We detail our compute resources in the appendix (refer to appendix).

793 Guidelines:

- 794 • The answer NA means that the paper does not include experiments.
795 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
796 or cloud provider, including relevant memory and storage.
797 • The paper should provide the amount of compute required for each of the individual
798 experimental runs as well as estimate the total compute.
799 • The paper should disclose whether the full research project required more compute
800 than the experiments reported in the paper (e.g., preliminary or failed experiments that
801 didn't make it into the paper).

802 **9. Code Of Ethics**

803 Question: Does the research conducted in the paper conform, in every respect, with the
804 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

805 Answer: [Yes]

806 Justification: We follow the NeurIPS Code of Ethics.

807 Guidelines:

- 808 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
809 • If the authors answer No, they should explain the special circumstances that require a
810 deviation from the Code of Ethics.
811 • The authors should make sure to preserve anonymity (e.g., if there is a special consider-
812 ation due to laws or regulations in their jurisdiction).

813 **10. Broader Impacts**

814 Question: Does the paper discuss both potential positive societal impacts and negative
815 societal impacts of the work performed?

816 Answer: [Yes]

817 Justification: We discuss the potential positive and negative societal impacts of our work in
818 the broader impacts section (refer to broader impacts).

819 Guidelines:

- 820 • The answer NA means that there is no societal impact of the work performed.
821 • If the authors answer NA or No, they should explain why their work has no societal
822 impact or why the paper does not address societal impact.
823 • Examples of negative societal impacts include potential malicious or unintended uses
824 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
825 (e.g., deployment of technologies that could make decisions that unfairly impact specific
826 groups), privacy considerations, and security considerations.
827 • The conference expects that many papers will be foundational research and not tied
828 to particular applications, let alone deployments. However, if there is a direct path to
829 any negative applications, the authors should point it out. For example, it is legitimate
830 to point out that an improvement in the quality of generative models could be used to
831 generate deepfakes for disinformation. On the other hand, it is not needed to point out
832 that a generic algorithm for optimizing neural networks could enable people to train
833 models that generate Deepfakes faster.
834 • The authors should consider possible harms that could arise when the technology is
835 being used as intended and functioning correctly, harms that could arise when the
836 technology is being used as intended but gives incorrect results, and harms following
837 from (intentional or unintentional) misuse of the technology.
838 • If there are negative societal impacts, the authors could also discuss possible mitigation
839 strategies (e.g., gated release of models, providing defenses in addition to attacks,
840 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
841 feedback over time, improving the efficiency and accessibility of ML).

842 **11. Safeguards**

843 Question: Does the paper describe safeguards that have been put in place for responsible
844 release of data or models that have a high risk for misuse (e.g., pretrained language models,
845 image generators, or scraped datasets)?

846 Answer: [NA]

847 Justification: We do not release any data or models that have a high risk for misuse.

848 Guidelines:

- 849 • The answer NA means that the paper poses no such risks.
- 850 • Released models that have a high risk for misuse or dual-use should be released with
851 necessary safeguards to allow for controlled use of the model, for example by requiring
852 that users adhere to usage guidelines or restrictions to access the model or implementing
853 safety filters.
- 854 • Datasets that have been scraped from the Internet could pose safety risks. The authors
855 should describe how they avoided releasing unsafe images.
- 856 • We recognize that providing effective safeguards is challenging, and many papers do
857 not require this, but we encourage authors to take this into account and make a best
858 faith effort.

859 12. Licenses for existing assets

860 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
861 the paper, properly credited and are the license and terms of use explicitly mentioned and
862 properly respected?

863 Answer: [Yes]

864 Justification: We properly credit the creators of the assets we use (refer to appendix).

865 Guidelines:

- 866 • The answer NA means that the paper does not use existing assets.
- 867 • The authors should cite the original paper that produced the code package or dataset.
- 868 • The authors should state which version of the asset is used and, if possible, include a
869 URL.
- 870 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 871 • For scraped data from a particular source (e.g., website), the copyright and terms of
872 service of that source should be provided.
- 873 • If assets are released, the license, copyright information, and terms of use in the
874 package should be provided. For popular datasets, paperswithcode.com/datasets
875 has curated licenses for some datasets. Their licensing guide can help determine the
876 license of a dataset.
- 877 • For existing datasets that are re-packaged, both the original license and the license of
878 the derived asset (if it has changed) should be provided.
- 879 • If this information is not available online, the authors are encouraged to reach out to
880 the asset's creators.

881 13. New Assets

882 Question: Are new assets introduced in the paper well documented and is the documentation
883 provided alongside the assets?

884 Answer: [NA]

885 Justification: We do not release any new assets.

886 Guidelines:

- 887 • The answer NA means that the paper does not release new assets.
- 888 • Researchers should communicate the details of the dataset/code/model as part of their
889 submissions via structured templates. This includes details about training, license,
890 limitations, etc.
- 891 • The paper should discuss whether and how consent was obtained from people whose
892 asset is used.
- 893 • At submission time, remember to anonymize your assets (if applicable). You can either
894 create an anonymized URL or include an anonymized zip file.

895 14. Crowdsourcing and Research with Human Subjects

896 Question: For crowdsourcing experiments and research with human subjects, does the paper
897 include the full text of instructions given to participants and screenshots, if applicable, as
898 well as details about compensation (if any)?

899 Answer: [NA]

900 Justification: We do not involve crowdsourcing or research with human subjects.

901 Guidelines:

- 902 • The answer NA means that the paper does not involve crowdsourcing nor research with
903 human subjects.
- 904 • Including this information in the supplemental material is fine, but if the main contribu-
905 tion of the paper involves human subjects, then as much detail as possible should be
906 included in the main paper.
- 907 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
908 or other labor should be paid at least the minimum wage in the country of the data
909 collector.

910 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
911 Subjects**

912 Question: Does the paper describe potential risks incurred by study participants, whether
913 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
914 approvals (or an equivalent approval/review based on the requirements of your country or
915 institution) were obtained?

916 Answer: [NA]

917 Justification: We do not involve crowdsourcing or research with human subjects.

918 Guidelines:

- 919 • The answer NA means that the paper does not involve crowdsourcing nor research with
920 human subjects.
- 921 • Depending on the country in which research is conducted, IRB approval (or equivalent)
922 may be required for any human subjects research. If you obtained IRB approval, you
923 should clearly state this in the paper.
- 924 • We recognize that the procedures for this may vary significantly between institutions
925 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
926 guidelines for their institution.
- 927 • For initial submissions, do not include any information that would break anonymity (if
928 applicable), such as the institution conducting the review.

929 **A Data schedule equivalences**

930 As discussed in Section 3.2, data schedules have two degrees of freedom. However, they can be
931 described with many intuitive variables, each a few equations away from each other. We can use the
932 following variables to describe the data schedule:

- 933 • **Total training steps T :** the total number of training steps.
934 • **Rare step fraction γ :** the fraction of training steps that are rare, after deciding the repetition count.
935 • **Replay fraction ρ :** the fraction of pre-training data that is replayed.
936 • **Rare stage 2 allocation α :** the fraction of the total rare data that is allocated to stage 2.
937 • **Stage 2 duration δ :** the fraction of training that is in stage 2.
938 • **Stage 1 rare weight w_1 :** the weight of the rare data in stage 1.
939 • **Stage 2 rare weight w_2 :** the weight of the common data in stage 2.

940 If there are 7 variables, why are there only 2 degrees of freedom? The first two variables are set by
941 problem setting and we do not control them (besides the repetition count, which we treat as fixed
942 for this section). For the rest of this section, without loss of generality, we set $T = 1$. We claim that
943 given the next 2 variables, you can derive the other 3.

944 If the replay fraction is ρ , then the stage 2 rare weight is automatically fixed as $w_2 = 1 - \rho$. If the
945 stage 2 allocation is α , then we know that the number of rare steps in stage 2 is $\alpha\gamma$. This means that
946 the number of total steps in stage 2 is $\frac{\alpha\gamma}{1-\rho}$, giving the stage 2 duration δ . Now that we know δ , it
947 suffices to determine the stage 1 rare weight. We know that there are $\gamma(1 - \alpha)$ rare steps in stage 1, as
948 well as $1 - \delta$ total steps. This means that the stage 1 rare weight is $w_1 = \frac{\gamma(1-\alpha)}{1-\delta}$. Therefore, we've
949 recovered all 7 variables from the 2 degrees of freedom. You can confirm that with these choices of
950 w_1, w_2 that $w_1(1 - \delta) + w_2\delta = \gamma$.

951 **B Potential failure modes of fine-tuning**

952 We discuss potential conceptual failure modes of fine-tuning here, using a mix of experiments and
953 toy models for intuition.

954 **B.1 Instability of fine-tuning**

955 We notice that at the start of fine-tuning, there is a pretty large loss spike (Figure 10). This is especially
956 true at higher learning rates. However, after some steps of training, the loss goes below where it
957 started. In fact, it is still correct to use higher learning rates even though they make the spike larger
958 because you end up with a lower loss.

959 One relatively vague hypothesis is that the loss spike is the reason fine-tuning underperforms replay.
960 This can happen in at least two concrete ways

- 961 1. When there is replay data, there is less distribution shift between stage 1 and stage 2.
962 Therefore, the loss spike is less pronounced, and there is less steps spent recovering from it.
- 963 2. There seems to be a minimum number of steps before one recovers from the loss spike.
964 Since replay increases the number of steps in stage 2, it can perhaps give more time to
965 recover from the loss spike.

966 We believe further experimentation is necessary to fully understand the nature of this spike, specifically
967 whether it is harmful for model performance and at what data regime it matters the most in.

968 **B.2 Overfitting to rare data**

969 It is known that a fixed un-regularized model has a tendency to overfit to small sample counts.
970 Therefore, it's possible that fine-tuning suffers from this problem. To model this, we set up a simple
971 linear regression toy model.

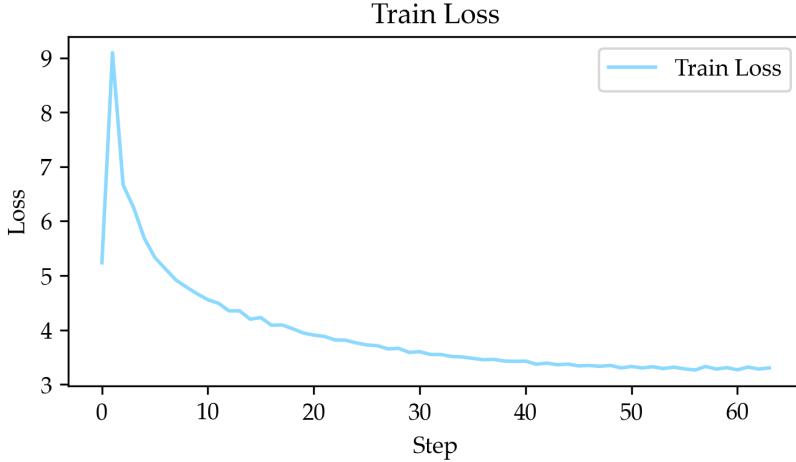


Figure 10: **Train loss spike.** We notice a large loss spike during the first few steps of training across our settings. This represents one barrier to training with a limited number of samples. Replay might help either because (1) there is less distribution shift between stage 1 and stage 2 or (2) there is more time to recover from the spike.

972 We construct a data distribution in 400 dimensions, governed by a pre-training $\theta_{\text{PT}} \sim \mathcal{N}(0, I)$ and
 973 a fine-tuning $\theta_{\text{FT}} \sim \mathcal{N}(\theta_{\text{PT}}, 0.1I)$. We then construct a dataset of N pre-training points and n
 974 fine-tuning points. Each point is generated as $(x, \theta^\top x + \epsilon)$ for $x \sim \mathcal{N}(0, I)$ and $\epsilon \sim \mathcal{N}(0, 1)$.

975 Our training algorithm first "pre-trains" by performing OLS on the pre-training points. For sufficiently
 976 large $N \gg d$, this will easily fit the pre-training vector θ_{PT} . However, our goal is to learn the
 977 fine-tuning vector θ_{FT} . We do this by performing OLS on the residuals of the pre-training fit and
 978 the true labels for the fine-tuning points. The fine-tuning now benefits from the pre-training fit
 979 bringing the parameters closer to the true fine-tuning vector. When we are in the over-parameterized
 980 regime, we use min-norm least squares regression, as double-descent literature has shown this is
 981 known to generalize better and is reflective of deep learning inductive bias [Belkin et al., 2019, Hastie
 982 et al., 2020]. We then measure error as mean squared error between the learned parameter and true
 983 fine-tuning parameter. We visualize these results in Figure 11, purple line. We see that for $n < d$, the
 984 model overfits to the noise in the fine-tuning data, resulting in higher error than random guessing. We
 985 plot the gray line to track the best possible error achievable by the model for a given n by using any
 986 sample count up to n .

987 We now introduce replay: mix in some pre-training data for the fine-tuning OLS. The replay
 988 significantly reduces the overfitting as we see in the other lines.

989 In this setting, it is known that the Bayes-optimal estimator involves appropriately tuning the ridge
 990 regularization parameter. We visualize differing values of the ridge parameter in Figure 12, orange
 991 line. We see that the optimal ridge parameter is non-zero. Moreover, the best ridge parameter achieves
 992 much better loss than tuned replay count. If this intuition holds, real language model training should
 993 benefit from finding the correct notion of regularization. One might expect the natural analogue of
 994 ridge regression for language models is weight decay. As discussed in E.1.3, weight decay does not
 995 significantly improve the loss and under-performs optimal replay. This tells us we need to rethink
 996 how we regularize fine-tuning to extract the full value of rare data points.

997 B.3 Model size

998 One concern is that the necessity of replay is due to the model size. To see whether any explanation
 999 related to model size holds ground, we see whether changing model size changes the necessity of
 1000 replay. To measure this, we take our joint training setting with a fixed learning rate schedule and
 1001 $\alpha = 1.0$ (all rare data in stage 2) and vary the model size. When we increase model size, we decrease
 1002 the learning rate inversely proportionally to the width of the hidden dimension, following folklore
 1003 scaling practices [Everett et al., 2024]. We visualize the results in Figure 13. We see that across all

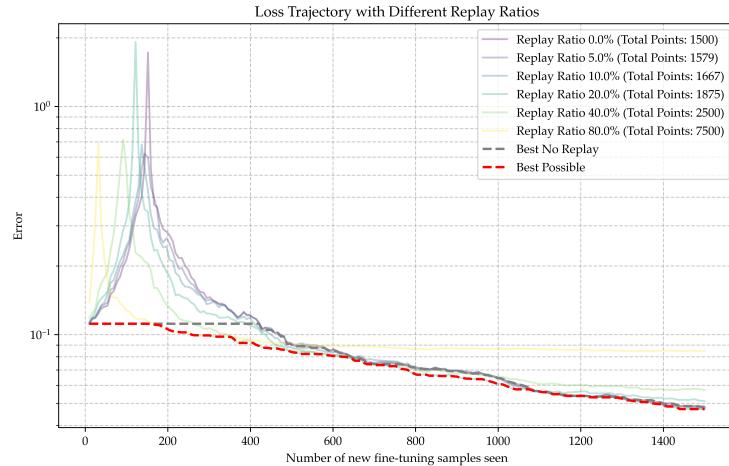


Figure 11: Linear regression with replay. We plot the loss of a linear regression model as a function of the number of fine-tuning points n for different values of the replay fraction ρ . We see that for $n < d$, the model overfits to the noise in the fine-tuning data, resulting in higher error than random guessing. We plot the gray line to track the best possible error achievable by the model for a given n by using any sample count up to n . Replay significantly reduces the overfitting, resulting in better MSE than random guessing. We track the best possible error as a function of n for the best ρ as the red line. For a regime of some but not too many rare points, replay improves over standard fine-tuning.

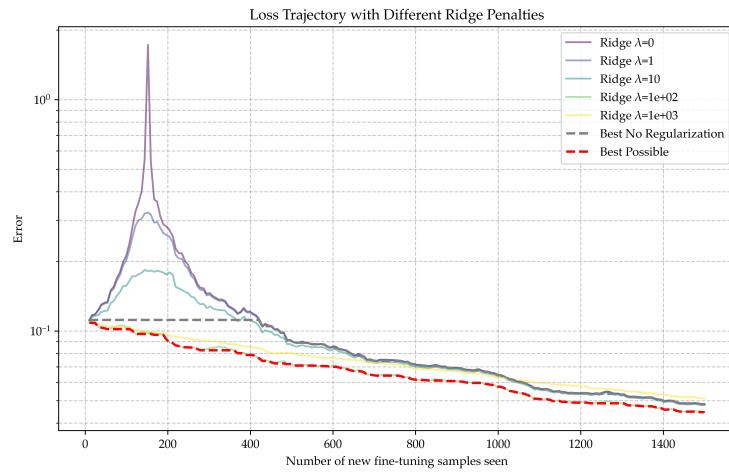


Figure 12: Linear regression with ridge regularization. We follow the same setup as in Figure 11, but instead of tuning replay fraction, we tune the ridge regularization parameter λ . We see that the optimal ridge parameter is non-zero and that the best ridge parameter achieves much better loss than tuned replay count.

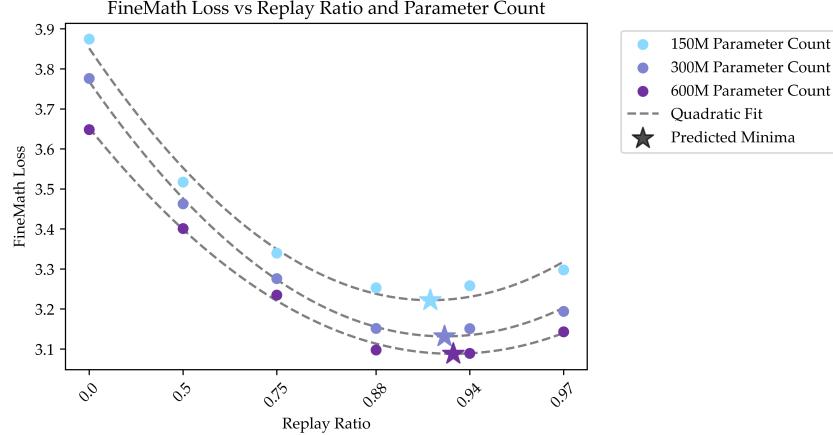


Figure 13: **Model scaling.** We take our standard 4B token training setup and scale the parameter count to 4x parameter count. For a given size, we sweep for the optimal replay fraction ρ while reserving all the rare data for stage 2 ($\alpha = 1.0$). We find that larger models still require replay data to get lower loss.

1004 sizes, the model benefits from replay. Furthermore, the benefit of replay is relatively consistent across
1005 model sizes.

1006 One interesting implication of this finding is that one can determine the optimal data schedule for
1007 a large model by tuning on a small model. This resembles μ P style arguments [Yang et al., 2022,
1008 2024a] for setting layer-wise learning rates at small model counts.

1009 C Note on tuning appendices

1010 This project has spanned a lot of experiments. The order these were conducted in does not reflect the
1011 order in which they are presented. At a high level, most of the mid-training/pre-training experiments
1012 were done first, giving intuition for what the results would look like for supervised fine-tuning. This
1013 means the experiments are more comprehensive for mid-training as it was when there was a worse
1014 understanding of the relationship between problem parameters. As the project developed, we were
1015 able to reduce search spaces by good priors on what hyperparameters worked (though we always
1016 certified they were correct as shared in the plots).

1017 We believe it is more instructive to read the guide for how to set mid-training hyperparameters.
1018 We believe the literature lacks rigorous experiments (mostly deciding random hypers on the fly).
1019 However, this stage can give large data efficiency gains if done correctly.

1020 D General training settings

1021 We train a 150 million parameter Llama-style language model with context length 4096 for 4B tokens.
1022 This is close to Chinchilla optimal scaling [Hoffmann et al., 2022] which prescribes 20x tokens per
1023 parameter. We train with batch size 1024 for 1024 steps with weight decay 0.1. We use the Adam
1024 optimizer with default parameters. When we tune learning rate, we search for the nearest power of 3
1025 assuming the final loss is convex in the learning rate. For the other models, we scale the learning rate
1026 inversely with the model width following the recommendations of [Everett et al., 2024] (see Table 1).

1027 For our common data, we use C4 [Raffel et al., 2023] since it is filtered but relatively uncurated. For
1028 example, C4 filters out all code data by removing documents with curly braces. For our rare data, we
1029 have domains representing math (FineMath [Allal et al., 2025]), coding (StarCoder [Li et al., 2023]),
1030 and instruction following (Flan [Longpre et al., 2023]). Our validation datasets are always the same
1031 distribution as our training data.

Parameter	150M	300M	600M	1.9B
Context Length	4096	4096	4096	4096
Hidden Dimension	512	768	1024	2048
Intermediate Dimension	1792	2688	3584	7168
Attention Heads	8	12	16	16
KV Heads	8	12	8	8
Layers	6	12	24	24
Learning rate	3e-3	2e-3	1.5e-3	7.5e-4

Table 1: Model architecture configurations for different model sizes. All models use the Llama architecture with a standardized context length of 4096 tokens. We default to the 150M model if not specified.

1032 **D.1 Model configurations**

1033 **D.2 Magic number justifications**

1034 Selecting a training regime requires setting some arbitrary numbers. We give justification for some
1035 here.

- 1036 • Rare data fraction: Pre-training token counts are on the order of 10 trillion while domains are on the
1037 order of 10 billion tokens, motivating our choice of $\approx 0.1\%$. We actually use a rare data fraction of
1038 $\frac{1}{1024}$; this better interacts with our block-deterministic data scheduler which draws/shuffles 2048
1039 sequences at a time.
- 1040 • Replay fractions and rare data allocations: In early experiments we quickly realized that the
1041 dependence on these parameters scaled nicely when the replay fractions are spaced out by $\log(1 -$
1042 $x)$. Therefore, we equally spaced out values. In plots, we round all values to two decimals.
1043 In actuality, our replay fractions were 0.25, 0.5, 0.75, 0.875 and our rare data allocations were
1044 1.0, 0.5, 0.25, 0.125. The power of 2 scaling similarly interacts nicely with our block-deterministic
1045 data scheduler.
- 1046 • Model size: 150M reflects a model scale that is large enough to be representative and scale nicely.
1047 It also enabled quicker iteration than larger scale models. During the course of this project, we
1048 sanity checked our results held at larger scales, increasing our confidence in using smaller scale
1049 models.

1050 **E Mid-training experiments**

1051 **E.1 Fine-tuning baseline**

1052 **E.1.1 Repetitions**

1053 We try varying the number of repetitions of the rare data during mid-training. For this tuning, since
1054 we do not know the learning rate and schedule yet, we tune across the two most promising learning
1055 rates of 1e-3 and 3e-3 with no learning rate cooldown (as this is closer to our final learning rate
1056 schedule than full decay). We visualize the best of both in Figure 14. We find that we can tolerate up
1057 to 32 repetitions of the rare data before overfitting across all domains.

1058 **E.1.2 Learning rate cooldown**

1059 We vary the cooldown duration of a standard WSD learning rate schedule with 10 step warmup while
1060 fixing all 32 repetitions of the rare data to appear at the end of training. We also vary the learning
1061 rate to be 1e-3, 3e-3, and 1e-2. For all training runs, 3e-3 did best and 1e-2 always diverged, so we
1062 can safely only look at WSD with learning rate 3e-3. We visualize the final result in Figure 3. We
1063 find it critical to have a cooldown period as opposed to the more conventional long cooldown period
1064 (cooldown duration 0.99).

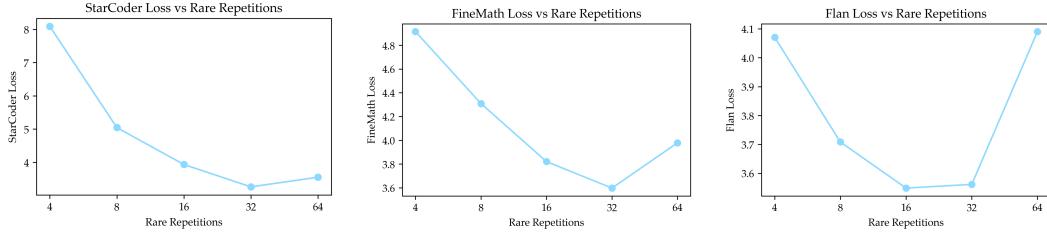


Figure 14: Mid-training tuning repetitions. We show that across our mid-training domains, we can tolerate maximum 32 repetitions of the original data before overfitting to the rare data. Note that the x-axis is more comprehensive compared to 16.

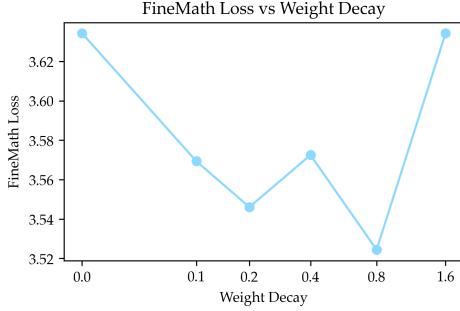


Figure 15: Tuning weight decay. We tune the weight decay for the fine-tuning baseline fixed to our above repetition count and learning rate cooldown. The effect is small and noisy, so we default to 0.1, at the upper range of weight decays used in the literature.

1065 E.1.3 Tuning weight decay

1066 We try tuning the weight decay for the fine-tuning baseline fixed to our above repetition count and
 1067 learning rate cooldown. We visualize the results in Figure 15. We find that weight decay does give
 1068 minimal but noisy effect on loss. Due to this variance, we decide to stay closer to the range of weight
 1069 decays used in the literature which are usually around 0.05 (for example, Table 10 and 11 in the
 1070 appendix of Li et al. [2025]) However, since there does seem to be a benefit of slightly higher weight
 1071 decay, we use 0.1 for all of our mid-training experiments. We also carry this choice to our pre-training
 1072 experiments.

1073 F Post-training experiments

1074 F.1 Fine-tuning baseline

1075 F.1.1 Repetitions

1076 We try varying the number of repetitions of the rare data during fine-tuning. We visualize the loss for
 1077 different repetition counts in 16. We find that we can tolerate up to 64 repetitions of the rare data
 1078 before overfitting across all domains.

1079 F.1.2 Learning rate

1080 For the model pre-trained on C4, we tried different learning rates for 1000 training steps and picked
 1081 the model with the best C4 loss. This ended up being 1e-3 (which achieved 4.12 loss, relative to 3e-3
 1082 which achieved 4.22 and 3e-4 which achieved 4.66). This learning rate was used across all pre-trained
 1083 models, which ranged from 512 to 992 steps. Note that the optimal learning rate is different for this
 1084 setting compared to the mid-training experiments because we’re using a different learning rate.

1085 G Data efficiency

1086 It is important to compare how well two training algorithms leverage the same amount of fixed
 1087 samples. Though we can compare the direct loss, this gives a misleading impression of how what the

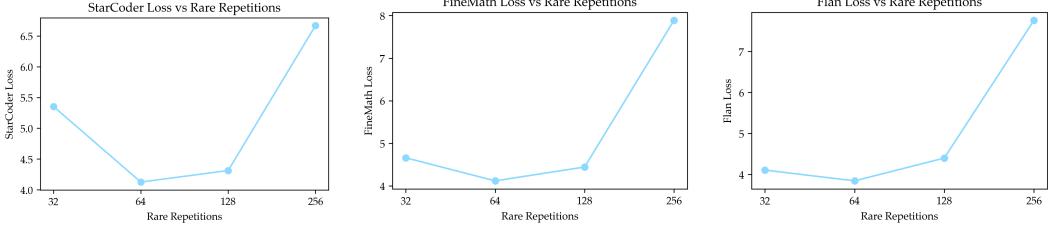


Figure 16: **Fine-tuning repetitions.** We show that across our fine-tuning domains, we can tolerate maximum 64 repetitions of the original data before overfitting to the rare data.

1088 actual improvement is. To bring this into a human-friendly metric, we introduce the data efficiency
 1089 multiplier. We use the following procedure:

- 1090 1. **Fix a reference training algorithm.** To enable comparison across a broad suite of possible
 1091 training strategies without having to refit scaling laws, we first fix a reference training
 1092 algorithm S_{ref} that characterizes one natural usage of the training data. We detail our choice
 1093 of reference algorithm below.
- 1094 2. **Build a scaling law.** We now build a reference scaling law to characterize the loss of the
 1095 reference algorithm as a function of the number of rare tokens it gets to see. Since the
 1096 model size is fixed, we fit a power law from number of data points seen to loss (more details
 1097 below).
- 1098 3. **Effective data points.** For any given strategy S , we can find its loss. With this loss, we can
 1099 see how many data points $D(S)$ it would take for the reference algorithm to match the loss
 1100 of S . If this number is high, then S is very data efficient.
- 1101 4. **Normalize for reference.** However, our current metric strongly depends on the choice
 1102 of reference algorithm. To make this metric useful regardless of the reference algorithm,
 1103 we always report a *data efficiency improvement*. Specifically, we only use this metric to
 1104 compare the data efficiency of two strategies S_1 and S_2 . We then report the data efficiency
 1105 improvement of S_2 over S_1 as $\frac{D(S_2)}{D(S_1)}$. A large improvement means that S_2 is much more
 1106 data efficient than S_1 .

1107 We now go into details about how we fit the scaling law.

1108 G.1 Scaling law formulation

1109 We fix the model size and total number of tokens. We then fit a power law from number of rare data
 1110 points given to loss. Our power law has the form $L(D) = aD^b + c$ for loss L , number of rare data
 1111 points D , and free variables a, b, c . We use `scipy.optimize.curve_fit` to fit the scaling law.

1112 G.2 Training runs

1113 We fix learning rate schedule to be cosine and tune learning rate to be 0.003. When tuning epoch
 1114 count, we found that the model could not tolerate more than 32 epochs of rare data at larger rare data
 1115 fractions. Therefore, we fix this epoch count. We also mix data uniformly throughout training instead
 1116 of using a dedicated data schedule. We train for a total of 4B tokens and vary the number of rare
 1117 tokens to be 4M, 8M, 16M, 32M, 64M tokens. Since we train for 32 repetitions, our largest rare data
 1118 run trains on rare tokens for 50% of training.

1119 H WSD tutorial

1120 It turns out that the correct learning rate schedule is critical for improving rare data efficiency.
 1121 Though annealing-based learning rate schedules give the largest benefit for data ordering, we have
 1122 relatively little intuitive/empirical understanding of how to properly use them. We will provide a
 1123 quick introduction to how they work in the un-ordered setting, and then show how to use them in the
 1124 ordered setting.

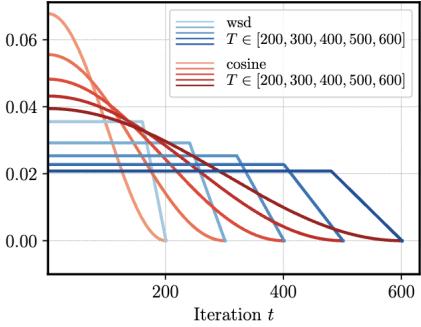


Figure 17: **Learning rate schedule.** This figure shows the shape of a WSD learning rate schedule in contrast to a cosine learning rate schedule (both without warmup). Figure is taken from Schaipp et al. [2025], Figure 2 left.

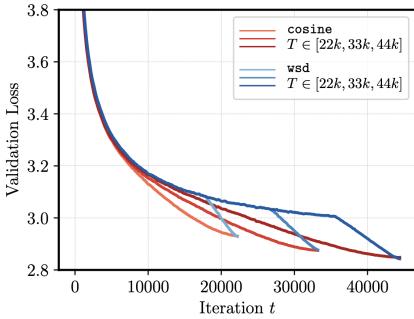


Figure 18: **Loss curves.** This figure shows the loss curves for a WSD learning rate schedule and a cosine learning rate schedule. Note that the loss of a WSD schedule initially makes slower progress than a cosine schedule and then makes up for it with much faster loss improvement at the end. Figure is taken from Schaipp et al. [2025], Figure 1 left.

1125 H.1 What is WSD?

1126 Warmup-Stable-Decay (WSD) is a learning rate schedule with three phases:

- 1127 1. Warmup: The learning rate is increased linearly from 0 to a peak value.
- 1128 2. Stable: The learning rate is held constant at the peak value.
- 1129 3. Decay: The learning rate is decayed linearly from the peak value to 0.

1130 We visually depict this learning rate schedule (without warmup) in Figure 17, left.

1131 H.2 Standard training (random order)

1132 Though warmup is important, the exact duration of the warmup is not critical. In contrast, the decay
 1133 period is critical to the final loss. We visualize the loss curves for a WSD learning rate schedule and a
 1134 cosine learning rate schedule in Figure 18, right. Notably, as soon as the learning rate starts decaying,
 1135 the loss improvement is much faster. This is contrast to cosine learning rate schedules where the
 1136 loss improvement actually slows down at the end of training (with a characteristic curl up). These
 1137 different rates of decrease have historically been really important details: for example, fitting the
 1138 scaling laws to intermediate checkpoints gives incorrect scaling laws since the models have been
 1139 annealed for different durations [Hoffmann et al., 2022].

1140 Why does this happen? One nice intuitive picture is given by the river valley landscape explanation
 1141 in Wen et al. [2024]. They posit that the loss landscape looks like a single river flowing down in the
 1142 middle of a valley (Figure 19). In this picture, you would like to both get to the bottom of the valley
 1143 and go far down the river. A standard learning rate schedule slowly descends the valley while also
 1144 making progress along the river direction. The paper's central claim is that WSD instead stays at the
 1145 top of the valley but continues to make progress along the river direction. Then, when one anneals
 1146 the learning rate, it starts descending down the valley, revealing the true progress made by the model
 1147 not captured by the loss. It is noted in the Edge of Stability literature [Cohen et al., 2022, 2024] that
 1148 staying at high learning rate is better for performance even though there is large oscillation in the loss.

1149 Though this picture is helpful visually, there is an even simpler theoretical picture. The works Defazio
 1150 et al. [2024], Schaipp et al. [2025] show that the simple theoretical model of non-smooth convex
 1151 optimization predicts the shape of the loss curve. Specifically, the upper-bound on the loss from

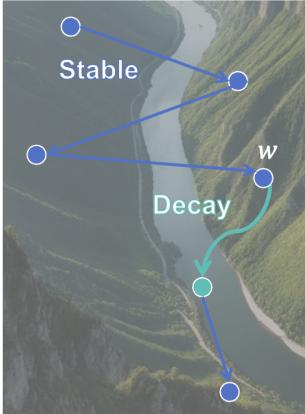


Figure 19: **River valley landscape.** This shows the intuitive picture of the river valley landscape. WSD makes progress along the river direction while making all the hill progress at the very end. Figure is taken from Wen et al. [2024], Figure 2 left.

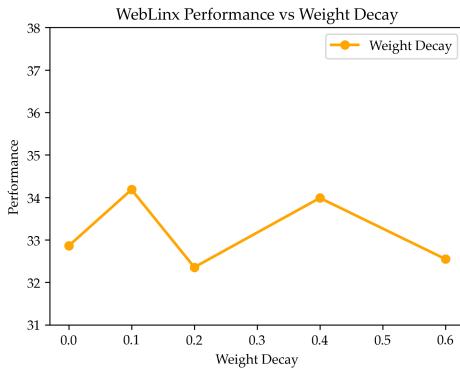


Figure 20: **Weblinx weight decay ablation.** We fine-tune Llama 3.1-8B Instruct on Weblinx demonstrations. We find that without replay, tuning weight decay gives little gains, improving by less than 2% over the baseline.

1152 standard online convex optimization arguments applied to the last iterate of training matches the
 1153 "shape" of the WSD loss curve.

1154 **H.3 Ordered training**

1155 Why does it matter that WSD decreases loss faster at the end of training? Intuitively, if the loss is
 1156 decreasing faster, placing high quality data at the end of training is more important. We algorithmically
 1157 leverage this intuition by placing the rare data at the end of training with WSD. In some earlier
 1158 experiments, we found that when using a cosine learning rate schedule, it actually hurt to keep rare
 1159 data at the end of training relative to placing it uniformly throughout training.

1160 **I Web agents**

1161 We train with the same hyperparameters as the original Weblinx paper on a subset of the demon-
 1162 strations from the original paper. We use the same evaluation protocol and metrics as the original
 1163 paper by combining the validation and in-distribution test set. We defer to the original paper for more
 1164 details on the data and evaluation. When we specify replay fraction, we are doing the replay on a
 1165 document level instead of a token level. This doesn't have large implications since all peak at an
 1166 intermediate value and fine-tuning isn't computationally intensive.

1167 We provide an additional ablation on tuning weight decay while fine-tuning on web agents data. We
 1168 find that without replay, this gives little gains, improving by less than 2% over the baseline. We
 1169 provide the results in Figure 20.

1170 **J Basque**

1171 We tune the learning rate to be 1e-5 for fine-tuning on Basque. We find the gain in accuracy to be real
 1172 across different token counts, displayed for 40M tokens and 200M tokens in Figure 21.

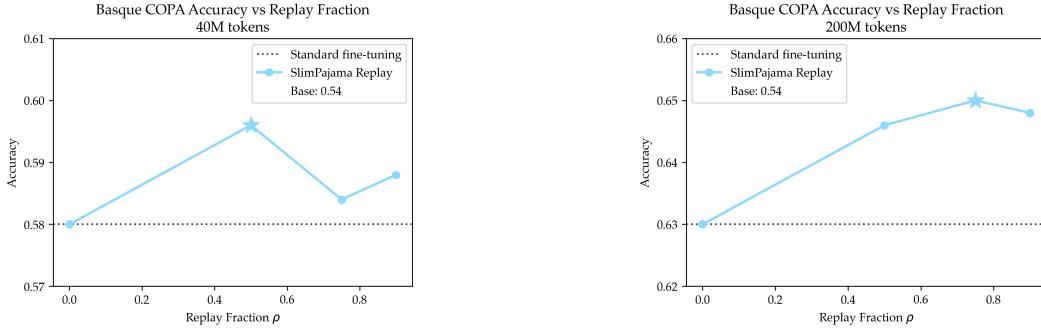


Figure 21: **Basque training with different token counts.** We try Basque training with 40M and 200M tokens to confirm that the gain in accuracy is real.

1173 When tracking Basque loss, we find that there is a spike in loss at the start of training only if the
 1174 learning rate is sufficiently high. In practice, it is worth using training with this higher learning rate
 1175 for best Basque loss/accuracy. We find that the loss improvement of replay decreases as we increase
 1176 the total token count. However, there continues to be an accuracy gain as you increase the token
 1177 count, showing that replay may be even more important for evaluation metrics.

1178 K Detailed related work

1179 K.1 Repeating data

1180 Prior work on data-constrained scaling laws [Muennighoff et al., 2023, Goyal et al., 2024] predict
 1181 that as you continue to repeat data, loss improves at a diminishing rate. However, the specific decay
 1182 formulation predicts that it will asymptote at a particular value.

1183 Specifically, the simplest decay formulation presented in both works estimates that when you see n
 1184 data points for the second time, it is like effectively training on $n\delta$ data points for decay factor δ . For
 1185 the k -th repetition, it is like training on $n\delta^{k-1}$ data points. In the infinite data limit, these scaling
 1186 laws predict that the loss will asymptote at a particular value, specifically at the loss of seeing $\frac{n}{1-\delta}$
 1187 fresh data points once.

1188 In our experiments, we found this not to be the case. If we repeated a rare domain too many times,
 1189 the loss would start going up. This is true whether it is fine-tuning or it is generic pre-training data.
 1190 This means we think more carefully about how to leverage rare data.