

WEIGHT DECAY IMPROVES LANGUAGE MODEL PLASTICITY

Anonymous authors

Paper under double-blind review

ABSTRACT

The prevailing paradigm in large language model (LLM) development is to pre-train a base model, then perform further training to improve performance and model behavior. However, hyperparameter optimization and scaling laws have been studied primarily from the perspective of the base model’s validation loss, ignoring downstream adaptability. In this work, we study pretraining from the perspective of *model plasticity*, that is, the ability of the base model to successfully adapt to downstream tasks through fine-tuning. We focus on the role of weight decay, a key regularization parameter during pretraining. Through systematic experiments, we show that models trained with larger weight decay values are more plastic, meaning they show larger performance gains when fine-tuned on downstream tasks. This phenomenon can lead to counterintuitive trade-offs where base models that perform worse after pretraining can perform better after fine-tuning. Further investigation of weight decay’s mechanistic effects on model behavior reveals that it encourages linearly separable representations, regularizes attention matrices, and reduces overfitting on the training data. In conclusion, this work casts light on the multifaceted role that a single optimization hyperparameter can play in shaping model behavior and demonstrates the importance of using evaluation metrics beyond the cross-entropy loss for hyperparameter optimization.

1 INTRODUCTION

Weight decay is a canonical hyperparameter whose role has evolved alongside changes in training regimes. In classical multi-epoch training, weight decay was understood primarily as a regularizer that improves generalization (Hardt et al., 2016; Zhang et al., 2017; Sun et al., 2025). In contemporary large-scale single-epoch pretraining, weight decay no longer primarily serves the purpose of generalization but plays a decisive role in optimization stability and convergence (D’Angelo et al., 2024; Zhang et al., 2025; Wang & Aitchison, 2024).

Moreover, modern language models are typically developed in two distinct stages: large-scale pretraining followed by post-training (Brown et al., 2020b; Ouyang et al., 2022; Bi et al., 2024; Lambert et al., 2024). While the two stages are functionally linked, current practices often treat them as decoupled. Specifically, pretraining hyperparameters and scaling laws are predominantly studied through the lens of the base model’s validation loss, under the assumption that a lower pretraining validation loss also yields a more capable downstream model (Hoffmann et al., 2022b; Bi et al., 2024). However, to what extent does optimizing pretraining hyperparameters for pretraining performance also optimize the final, post-trained model’s performance?

In this work, we study the relationship between pretraining and post-training from the perspective of *model plasticity*, i.e., a model’s ability to adapt to new data upon further training (Berariu et al., 2021; Dohare et al., 2024). Model plasticity naturally bridges pretraining and post-training, capturing how readily the pretrained model can be reshaped for downstream tasks. As we show, two models with similar pretraining loss may differ in their plasticity, meaning that optimizing hyperparameters for pretraining loss alone may not yield the best post-trained model.

In our experiments, we vary the weight decay parameter during pretraining and evaluate the pre-trained model’s ability to learn tasks during fine-tuning. Pretraining is performed for two model families (Llama-2 and OLMo-2), multiple model sizes (up to 4B parameters), and in both the compute-optimal (20 tokens-per-parameter, TTP hereafter) and overtrained (140 TPP) regimes. Fine-tuning

is performed across six Chain-of-Thought (CoT) tasks, and results are assessed using a comprehensive suite of metrics that cover both solution correctness and quality. Our experimental design takes an end-to-end perspective (Qi et al., 2025; Mayilvahanan et al., 2025), aligning pretraining hyperparameter selection with the ultimate objective of maximizing performance after further training. Our contributions are as follows:

- We show that weight decay improves model plasticity, facilitating adaptation to new tasks during subsequent fine-tuning. Across model families, scales, and training regimes, the evidence points toward an optimal pretraining weight decay value larger than the default of 0.1. This highlights the potential for re-evaluating standard hyperparameter choices to better account for a model’s downstream adaptability.
- We provide one of the first examples showing that optimizing hyperparameters to minimize pretraining validation loss does not necessarily yield the best downstream model performance. Specifically, we show that there is a training regime where larger weight decay values lead to higher pretraining validation loss *and* better downstream performance.
- We provide a mechanistic perspective on how weight decay shapes training dynamics, showing it encourages linearly separable representations, regularizes attention matrices, and reduces overfitting. These findings provide a potential explanation for how weight decay preserves the model’s ability to learn in subsequent training, thus sustaining plasticity.

2 RELATED WORK

Here, we discuss related work on weight decay and plasticity and how our work contributes new insights. A fuller discussion is in Appendix A.1.

Weight decay in language model training. Weight decay is a standard hyperparameter in language model training (Brown et al., 2020a; Grattafiori et al., 2024; OLMo et al., 2024). Beyond its classical role in regularization and generalization (Krogh & Hertz, 1991; Zhang et al., 2018; Loshchilov & Hutter, 2017; Zhou et al., 2024), it has also been shown to play other roles in language model training, such as improving optimization and training stability (D’Angelo et al., 2024; Li et al., 2020; Kosson et al., 2024; 2025; Wen et al., 2025), inducing low-rank attention layers (Kobayashi et al., 2024), and increasing forgetting of training data (Bordt et al., 2025). Prior works have studied how to set weight decay to minimize pretraining loss (Bergsma et al., 2025; Kim et al., 2025). In contrast to these works which focus on pretraining performance, this paper examines how weight decay affects model plasticity (downstream performance).

Plasticity of deep learning models. Model plasticity has been studied in continual learning, transfer learning, and reinforcement learning, where models often undergo multiple training rounds (Dohare et al., 2024; Klein et al., 2024; Coetzer et al., 2025). Prior works have demonstrated that image models lose plasticity when trained on new data (Dohare et al., 2024; Lyle et al., 2023; Klein et al., 2024) and various approaches have been developed to improve model plasticity (Ash & Adams, 2020; Dohare et al., 2024; Kumar et al., 2023; Miconi et al., 2018). While prior work has examined how forgetting and tokenization (Chen et al., 2023; Abagyan et al., 2025) affect language model plasticity, research on language model plasticity remains underdeveloped. In contrast to these works, this paper investigates the role of weight decay, a standard hyperparameter, in language model plasticity.

3 BACKGROUND AND METHODS

In this section, we provide further background, define the research question, and describe the experimental setup. A more detailed discussion is in Appendix A.2.

Weight decay in AdamW. This paper focuses on the weight decay hyperparameter λ in AdamW which, for each optimizer step $t \geq 1$, performs two decoupled updates: a gradient update given by (1) followed by a weight decay update given by (2)

$$(1) \theta_t = \hat{\theta}_t - \gamma_t \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon) \quad \text{and} \quad (2) \hat{\theta}_t = \theta_{t-1} - \gamma_t \lambda \theta_{t-1}$$

108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161

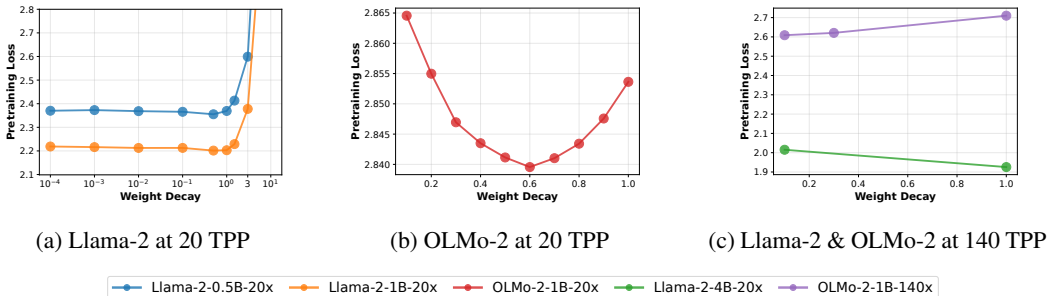


Figure 1: **Pretraining validation cross-entropy loss of models pretrained with varying weight decay.** The weight decay value that minimizes pretraining loss may be equal to or larger than the default value of 0.1 depending on the training regime.

based on model parameters θ_t , learning rate γ_t , and first- and second-order moment estimates of the gradient \hat{m}_t and \hat{v}_t (Loshchilov & Hutter, 2019). For language model pretraining, the choice $\lambda = 0.1$ has emerged as a kind of default (Brown et al., 2020a; Touvron et al., 2023; OLMo et al., 2024).

Language model plasticity. Following prior work (Berariu et al., 2021; Dohare et al., 2024), we measure plasticity by fine-tuning the model and evaluating performance on the fine-tuning task: the better the performance, the better the model learned during fine-tuning, thus the higher the plasticity.

In this context, we now specify the research question:

Research Question. How does weight decay during pretraining affect model plasticity, i.e. the pretrained model’s ability to learn new knowledge during subsequent training?

Experimental setup. To study this question, we pretrain five Llama-2 and OLMo-2 models with varying weight decay (Llama-2-0.5B-20x, Llama-2-1B-20x, Llama-2-4B-20x, OLMo-2-1B-20x, and OLMo-2-1B-140x). These models span different sizes (up to 4B) and training regimes (Chinchilla-optimal and overtrained). We fine-tune (SFT) these models on six CoT CoT tasks spanning various domains: MetaMathQA, MedMCQA, PubMedQA, MMLUProCoT, RACE, and SimpleScaling (Yu et al., 2023; Pal et al., 2022; Jin et al., 2019; Wang et al., 2024; Lai et al., 2017; Muennighoff et al., 2025). Then, we evaluate the fine-tuned models, examining the correctness and quality of their solutions to test set questions using six metrics: Greedy, Maj@16, RM@16, Pass@16, Correct Ratio, and ORM Score.

4 WEIGHT DECAY IMPROVES LANGUAGE MODEL PLASTICITY

We present the main experimental results. We begin by identifying the optimal pretraining weight decay based on pretraining performance (Section 4.1), a common way to select pretraining hyperparameters (Hoffmann et al., 2022a). Next, we investigate how weight decay shapes the plasticity of the pretrained model and identify its optimal pretraining value based on downstream performance (Section 4.2). Then, we examine whether a model’s pretraining performance is fully predictive its downstream performance (Section 4.3).

4.1 THE OPTIMAL PRETRAINING WEIGHT DECAY BASED ON PRETRAINING PERFORMANCE

We first identify the weight decay value that leads to the lowest cross-entropy validation loss after pretraining. This is the value considered optimal by current approaches in hyperparameter optimization for LLM pretraining (Bergsma et al., 2025). We pretrain models with varying weight decay and plot the validation cross-entropy loss of these pretrained models in Figure 1.

Extremely small weight decay values during pretraining do not have a significant effect on pretraining loss (Figure 1a). On the other hand, extremely large weight decay values can result in very high pretraining loss, significantly degrading pretraining performance (Figure 1a). At 20 TPP, for both

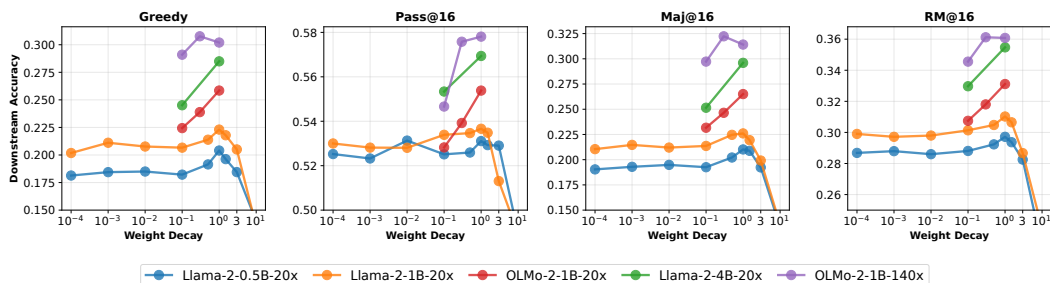


Figure 2: **Weight decay during pretraining improves language model plasticity and downstream performance.** The optimal weight decay for downstream performance may be larger than the commonly-used default value of 0.1. In addition, the optimal weight decay value based on pretraining performance (Figure 1) and based on downstream performance (this figure) are different, suggesting that approaches that optimize weight decay based only on pretraining performance are not guaranteed to produce models with the best downstream performance.

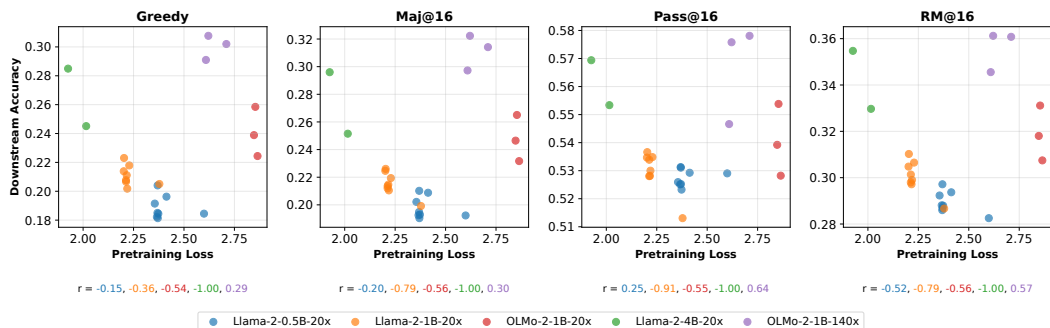


Figure 3: **A model’s performance after pretraining is not perfectly predictive of its performance downstream.** Models with lower cross entropy validation loss after pretraining can perform better or worse downstream (i.e., after fine-tuning) than models with higher pretraining losses, and models with similar pretraining losses can perform differently downstream.

Llama-2 and OLMo-2 models, we find the optimal weight decay is larger than the 0.1 default. In particular, among the values examined, the optimal weight decay is 0.5 for Llama-2-{0.5B and 1B}-20x, 0.6 for OLMo-2-1B-20x, and 1.0 for Llama-2-4B-20x. However, this relationship changes as training time increases. At 140 TPP, for the OLMo-2-1B-140x model, the 0.1 default outperforms (leads to a lower validation loss than) larger values of 0.3 and 1.0. This result that overtrained models have a lower optimal weight decay is consistent with previous analyses which recommend decreasing weight decay as training time (TPP) increases (Bergsma et al., 2025).

4.2 THE OPTIMAL PRETRAINING WEIGHT DECAY BASED ON DOWNSTREAM PERFORMANCE

Next, we investigate how weight decay during pretraining affects model plasticity and downstream model performance. We fine-tune the pretrained models from Section 4.1 on six CoT tasks and evaluate their performance on these tasks. Figure 2 shows the models’ average downstream accuracy across all tasks. Results for individual tasks and all six metrics are in Appendix D.

Among models that achieved reasonable pretraining validation losses in Section 4.1 (i.e., suitable candidates for subsequent training), higher weight decay during pretraining confers a higher degree of model plasticity, enabling the pretrained model to learn better during fine-tuning and perform better on the fine-tuning task. Across model families, model sizes, training regimes, fine-tuning tasks, and evaluation metrics, models pretrained with weight decay higher than the default 0.1 value perform better on downstream tasks. Among the weight decay values examined, in the compute-optimal 20 TPP regime, the optimal pretraining weight decay is 1.0 (Llama-2-0.5B-20x, Llama-2-1B-20x, Llama-2-4B-20x, and OLMo-2-1B-20x). In the overtrained 140 TPP regime, the optimal pretrain-

ing weight decay is 0.3 (OLMo-2-1B-140x). It is possible that as models are trained for even longer (i.e., beyond 140 TPP), the optimal pretraining weight decay that leads to best downstream model performance may continue to decrease (this is an extrapolation that would need to be validated).

We also compare the weight decay value that minimizes pretraining validation loss (Figure 1 in Section 4.1) with the value that maximizes task accuracy after fine-tuning (Figure 2 in this section). We find that these two weight decay values differ for each model. This shows that the “optimal” weight decay during pretraining is not absolute – it depends on the intended objective, such as optimizing for pretraining performance or downstream performance.

Finding 1. Pretraining weight decay can improve model plasticity and lead to better downstream performance. The optimal pretraining weight decay value for plasticity is larger than the default of 0.1.

4.3 THE RELATIONSHIP BETWEEN PRETRAINING LOSS AND DOWNSTREAM PERFORMANCE

We now investigate whether a model’s pretraining performance is predictive of its downstream performance. We plot the pretraining validation cross-entropy loss of the pretrained models from Section 4.1 and their accuracy on tasks after fine-tuning (measured in Section 4.2) in Figure 3.

We compare models with the same training setup (same model family, size, and TPP) that differ only in the pretraining weight decay hyperparameter. Although the Pearson correlation coefficient r between pretraining and downstream performance is negative for models trained at 20 TPP (Llama-2-0.5, Llama-2-1B, Llama-2-4B, and OLMo-2-1B) and positive for models trained at 140 TPP (OLMo-2-1B), this relationship is unstable both visually and upon further analyses (Figure 7). By examining pairs of points, the results show that models with better pretraining performance (lower loss after pretraining) can perform better downstream (e.g., for all five models) or worse downstream (e.g., for Llama-2-0.5B-20x, Llama-2-1B-20x, and OLMo-2-1B-140x). In addition, models with similar pretraining performance can perform differently downstream (e.g., for Llama-2-0.5B-20x, Llama-2-1B-20x, and OLMo-2-1B-20x). For example, OLMo-2-1B-140x pretrained with weight decay 0.3 or 1.0 performs slightly worse after pretraining (achieving pretraining cross-entropy validation losses of 2.6208 and 2.7064, respectively) than when pretrained with weight decay 0.1 (which achieves a pretraining cross-entropy validation loss of 2.6088), but the former two pretrained models perform noticeably better after fine-tuning (Figure 2, purple line). Altogether, these results indicate that pretraining performance is not necessarily predictive of downstream performance.

Finding 2. The pretraining weight decay value that minimizes the cross-entropy validation loss does not necessarily lead to the best downstream performance.

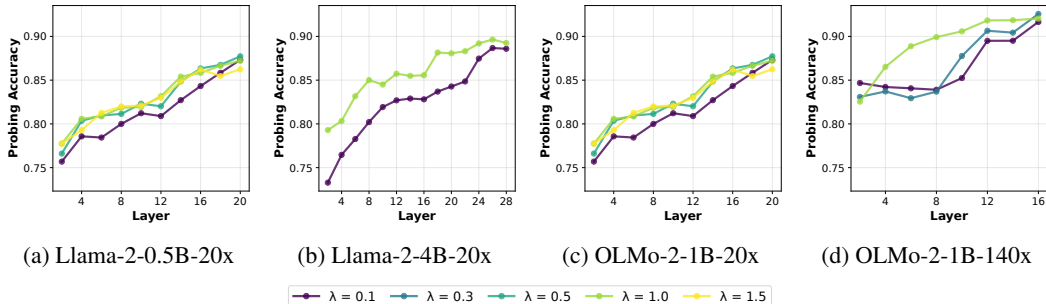
5 WEIGHT DECAY AND MODEL BEHAVIOR: A MECHANISTIC PERSPECTIVE

Prior work has shown that various factors can influence model plasticity, including the initialization state of model weights at the start of subsequent training, data representation (e.g., tokenization and categorical output representations), and model architecture (e.g., normalization layers) (Ash & Adams, 2020; Abagyan et al., 2025; Lyle et al., 2023). In Section 4, we find that weight decay also shapes model plasticity. Here, we explore three mechanisms through which weight decay shapes model behavior and discuss how each might explain why weight decay improves model plasticity.

5.1 WEIGHT DECAY ENCOURAGES LINEARLY SEPARATED REPRESENTATIONS

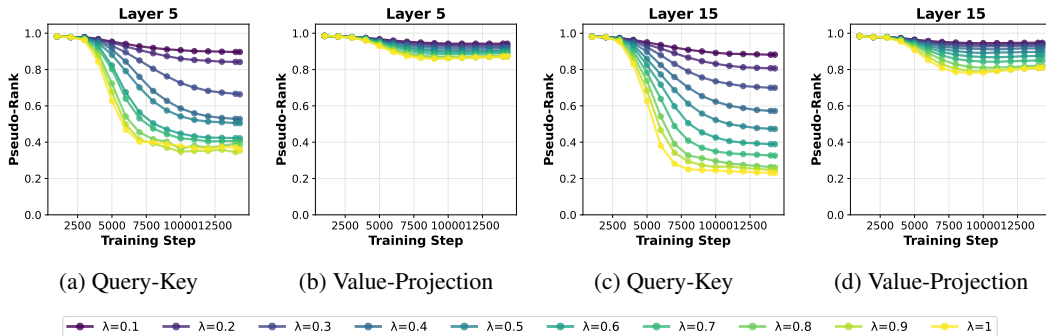
Inspired by previous findings that weight decay leads to more structured representations in vision models (Jacot et al., 2024), we investigate the effect of weight decay on the representations learned by pretrained language models. We pretrain models with varying weight decay, obtain the last-token embeddings for different types of text at a given model layer, and train a linear probe to classify embeddings based on sentiment (Stanford Sentiment Treebank dataset; Socher et al. (2013)) or based on topic (AG News dataset; Zhang et al. (2015)). The average accuracy of these linear probes over the two tasks is shown in Figure 4 (accuracy for individual tasks are in Appendix E.1).

270
271
272
273
274
275
276
277
278
279
280



281 **Figure 4: Weight decay encourages linearly separated representations.** This figure depicts the
282 accuracy of linear probes for sentiment and topic for models pretrained with different weight decay
283 values. We observe that linear probing achieves better accuracy when models are pretrained with a
284 weight decay greater than the default 0.1.

286
287
288
289
290
291
292
293
294
295
296



297 **Figure 5: Weight decay reduces the rank of attention matrices.** This figure depicts the average
298 pseudo-rank (Appendix E.2.1) of the query-key (W_{QK}) and value projection (W_{VP}) matrices in
299 layers 5 and 15 during the training of OLMo-2-1B models at 20 TPP.

302 When a given model is pretrained with higher weight decay, the accuracy of the linear probe trained
303 on the model’s representations tends to be higher at every layer of the model. While this relationship
304 is not perfectly monotonic (in some instances, a slightly higher weight decay can lead to a similar
305 or slightly lower probing accuracy), it is generally consistent across weight decay values and
306 model layers and across model families, sizes, and training regimes (i.e., for all five model setups).
307 Thus, through these linear probing experiments, we find that representations from models pretrained
308 with higher weight decay result in higher probing accuracies, indicating that they are more linearly
309 separated and suggesting that these models form more structured representations.

310 The finding that weight decay shapes the representations of pretrained language models points to
311 a potential explanation for why weight decay improves model plasticity (Section 4.2). Pretraining
312 models with higher weight decay produces models with more structured representations, i.e.,
313 representations in which information is encoded in a more linearly accessible form. As a result,
314 fine-tuning may effectively start at a better initialization and can focus on refining and aligning
315 existing representations to the fine-tuning task rather than continuing to learn representations, leading
316 to improved downstream performance. This hypothesis is consistent with previous findings in
317 transfer learning (Lee et al., 2023) and further supported by the strong positive correlation between
318 probing accuracy and downstream model performance (Figure 15).

320 **5.2 WEIGHT DECAY REDUCES THE RANK OF ATTENTION MATRICES**

322
323

Previous work by Kobayashi et al. (2024) provides a theoretical argument that weight decay should
reduce the rank of attention matrices. Recall that attention scores can be understood as a bilinear
form $X^T W_{QK} X$ where $W_{QK} = W_K^T W_Q \in \mathbb{R}^{n_{embed} \times n_{embed}}$ is the product of the query and key

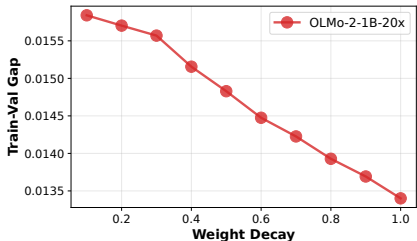


Figure 6: **Weight decay reduces overfitting on training data.** The figure depicts the train-val gap (Equation 1) for OLMo-2-1B models trained at 20 TPP.

matrices, and $X \in \mathbb{R}^{n_{embed} \times T}$ is the matrix of token embeddings (or hidden representations) for a sequence of length T . Now, the matrix W_{QK} is naturally low-rank since its rank is at most d_{head} , which is usually significantly smaller than n_{embed} . Kobayashi et al. (2024) argue that weight decay should further reduce the rank of W_{QK} , as well as of the value-projection matrix $W_{VP} = W_P W_V \in \mathbb{R}^{n_{embed} \times n_{embed}}$. Concretely, they show that L2 regularization applied to the factored matrices W_K and W_Q becomes equivalent to nuclear norm regularization on their product W_{QK} , which is known to induce low rank by promoting sparsity in the singular values. While Kobayashi et al. (2024) also provide empirical evidence on the Pile, their experiments were relatively small-scale from today’s perspective. We now revisit the impact of weight decay on the rank of attention in our more modern setup.

Weight decay reduces the rank of attention, but default weight decay yields near full-rank matrices. Figure 5 depicts the evolution of the pseudo-rank (Appendix E.2.1) of the attention matrices during the training of the OLMo-2-1B-20x models. From Figure 5, we observe that there is a monotonic relationship between the weight decay parameter and the rank of the attention matrices, where larger weight decay values reduce the rank of both W_{QK} and W_{VP} . However, unlike what is observed in Kobayashi et al. (2024), we see that the default weight decay parameter of 0.1 yields near full-rank matrices. This observation is further confirmed by Figure 18, which shows that the attention matrices in the fully trained OLMo-2-1B model are nearly full-rank.

Attention matrices are differentially affected by weight decay. Another important observation from our experiments is that the rank of the matrix W_{QK} seems to be significantly more sensitive to weight decay than W_{VP} . In our experiments, a weight decay of $\lambda = 1.0$ reduces the rank of W_{QK} by roughly a factor of 2, which is a common rank reduction observed in the literature on low-rank matrices. In contrast, the matrix W_{VP} is still close to full-rank even for a weight decay value of 1.0. These results are especially pronounced for Llama-2 models depicted in Figure 16, where the rank of W_{VP} remains essentially stable up to a weight decay value of 1.0, after which the rank collapses—a transition that correlates with a significant drop in performance.

Low-rank structure as a driver of adaptability. The observation that increased weight decay leads to lower-rank attention matrices provides a potential explanation for why weight decay improves model plasticity. In machine learning literature, low-rank constraints are a canonical form of regularization that is often believed to encourage simpler, more robust hypotheses (Cai et al., 2010; Oymak et al., 2019; Hu et al., 2022). We conjecture that by encouraging W_{QK} toward a lower-rank configuration, weight decay may prevent the model from overfitting to high-dimensional noise in the pretraining distribution.

5.3 WEIGHT DECAY REDUCES OVERFITTING ON TRAINING DATA

Lastly, we explore how weight decay influences the extent to which the pretrained model overfits the pretraining data. Previous work has shown that weight decay can cause the forgetting of individual benchmark questions seen during pretraining (Bordt et al., 2025). In the context of model plasticity, the ability to learn new information tends to be associated with the forgetting of prior data, a trade-off commonly referred to as the stability-plasticity dilemma (Kirkpatrick et al., 2017; Riemer et al., 2018; Ibrahim et al., 2024; Elsayed & Mahmood, 2024). Building on these insights, we investigate

378 how weight decay influences overfitting, which is closely related to the forgetting of training data,
379 in pretrained models.

380 To measure the degree to which a pretrained model overfits the training data, we compute the differ-
381 ence between the loss on the validation data and that on the training data:
382

$$383 \text{Train-Val Gap} = \text{Validation Loss} - \text{Training Loss} \quad (1)$$

384 Here, the training loss is the average loss that the fully trained model encounters on the training
385 data, which is distinct from the training loss curve or the final training loss value. A model that does
386 not overfit the training data would theoretically have a train-val gap of zero. In practice, a larger
387 train-val gap indicates a higher degree of overfitting on the training data, thus less forgetting of the
388 training data.
389

390 Figure 6 depicts the train-val gap for the OLMo-2 models trained at 20 TPP. We observe that the
391 train-val gap decreases monotonically as the weight decay parameter is increased. This provides
392 empirical evidence that models trained with larger weight decay values do indeed overfit the training
393 data less.

394 **Finding 3.** The pretraining weight decay parameter has diverse mechanistic effects on
395 model behavior. It encourages linearly separated representations, regularizes attention ma-
396 trices, and reduces overfitting on the training data.
397

399 6 DISCUSSION AND CONCLUDING REMARKS

400 This work provides a multidimensional characterization of the effects of the weight decay hyperpa-
401 rameter within the modern language-model training lifecycle. While traditional perspectives have
402 primarily viewed weight decay through the lenses of capacity control in over-parameterized regimes
403 or optimization stability in single-epoch pretraining, our findings suggest that weight decay plays
404 a far more nuanced role in shaping model behavior. In particular, we showed that models with
405 smaller weight decay achieve lower validation loss after pretraining (especially in the over-trained
406 regime), but that models with larger weight decay benefit from improved plasticity, enabling them
407 to perform best when fine-tuned on downstream tasks. Weight decay may shape model plasticity
408 through several mechanisms, including promoting linearly separable representations, regularizing
409 attention matrix ranks, and reducing overfitting on the training data. Together, these findings reveal
410 fundamental trade-offs in hyperparameter optimization. They also provide one of the first rigor-
411 ous empirical demonstrations that selecting pretraining hyperparameters based solely on minimal
412 pretraining validation loss can fail to yield the model with the highest performance on downstream
413 tasks.
414

415 The trade-offs we show mean that, in practice, the benefits of increased plasticity must be weighed
416 against other effects that may depend on model size, training duration, and other parameters of the
417 training setup. In heavily overtrained scenarios or for very large models trained for many steps
418 (Singh et al., 2025; Comanici et al., 2025; Anthropic, 2025), the benefits of markedly lower pre-
419 training validation loss may outweigh those of plasticity. In addition, weight decay’s diverse roles
420 in training dynamics – from plasticity (shown in this work) to optimization, training stability, con-
421 vergence rate, and overfitting (Hoffmann et al., 2022b; D’Angelo et al., 2024; Kosson et al., 2025) –
422 adds further complexity to model training decisions. Its optimal value for one objective can conflict
423 with that for another, as observed when optimizing for pretraining versus downstream performance.
424 A single weight decay value may not satisfy multiple objectives, requiring weighing trade-offs and
425 prioritizing objectives.

426 Future work may investigate in more detail the trade-offs between stability and plasticity, and the
427 extent to which our results hold in large-model and heavily overtrained scenarios. They may also
428 investigate the role of weight decay in model plasticity for foundation models beyond language (e.g.,
429 multimodal foundation models) and for other downstream desiderata (e.g., safety alignment).

430 Taken together, the findings in this work cast light on the multifaceted role that a single optimization
431 hyperparameter can play in shaping model behavior and the complexity of hyperparameter tuning
throughout the training of modern language models.

REFERENCES

- 432
433
434 Diana Abagyan, Alejandro R. Salamanca, Andres Felipe Cruz-Salinas, Kris Cao, Hangyu Lin, Acyr
435 Locatelli, Marzieh Fadaee, Ahmet Üstün, and Sara Hooker. One tokenizer to rule them all: Emer-
436 gent language plasticity via multilingual tokenizers. *arXiv preprint arXiv:2506.10766*, 2025. doi:
437 10.48550/arXiv.2506.10766. URL <https://arxiv.org/abs/2506.10766>.
- 438 Anthropic. System card: Claude opus 4 & claude sonnet 4, 2025. URL [https://www-cdn.
439 anthropic.com/6d8a8055020700718b0c49369f60816ba2a7c285.pdf](https://www-cdn.anthropic.com/6d8a8055020700718b0c49369f60816ba2a7c285.pdf).
- 440
441 Jordan Ash and Ryan P Adams. On warm-starting neural network training. *Advances in neural
442 information processing systems*, 33:3884–3894, 2020.
- 443 Tudor Berariu, Wojciech Czarnecki, Soham De, Jorg Bornschein, Samuel Smith, Razvan Pas-
444 canu, and Claudia Clopath. A study on the plasticity of neural networks. *arXiv preprint
445 arXiv:2106.00042*, 2021.
- 446
447 Shane Bergsma, Nolan Dey, Gurpreet Gosal, Gavia Gray, Daria Soboleva, and Joel Hestness.
448 Power lines: Scaling laws for weight decay and batch size in LLM pre-training. *arXiv preprint
449 arXiv:2505.13738*, 2025.
- 450
451 Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding,
452 Kai Dong, Qiusi Du, Zhe Fu, et al. Deepseek llm: Scaling open-source language models with
453 longtermism. *arXiv preprint arXiv:2401.02954*, 2024.
- 454 Sebastian Bordt and Martin Pawelczyk. Train once, answer all: Many pretraining experiments for
455 the cost of one. *arXiv preprint arXiv:2509.23383*, 2025.
- 456
457 Sebastian Bordt, Suraj Srinivas, Valentyn Boreiko, and Ulrike von Luxburg. How much can we
458 forget about data contamination? *ICML*, 2025.
- 459 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
460 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
461 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020a.
- 462
463 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, et al. Language
464 models are few-shot learners. *Advances in Neural Information Processing Systems (NeurIPS)*,
465 33, 2020b. URL <https://arxiv.org/abs/2005.14165>.
- 466
467 Jian-Feng Cai, Emmanuel J Candès, and Zuowei Shen. A singular value thresholding algorithm for
468 matrix completion. *SIAM Journal on optimization*, 20(4):1956–1982, 2010.
- 469 Albert Catalan-Tatjer, Niccolò Ajroldi, and Jonas Geiping. Training dynamics impact post-training
470 quantization robustness. *arXiv preprint arXiv:2510.06213*, 2025.
- 471
472 Yihong Chen, Kelly Marchisio, Roberta Raileanu, David Adelani, Pontus Lars Erik Saito Stenetorp,
473 Sebastian Riedel, and Mikel Artetxe. Improving language plasticity via pretraining with active
474 forgetting. *Advances in Neural Information Processing Systems*, 36:31543–31557, 2023.
- 475 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
476 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John
477 Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*,
478 2021.
- 479
480 Xander Coetzer, Arné Schreuder, and Anna Sergeevna Bosman. Restoring neural network plasticity
481 for faster transfer learning. In *Southern African Conference for Artificial Intelligence Research*,
482 pp. 206–222. Springer, 2025.
- 483 Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit
484 Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the
485 frontier with advanced reasoning, multimodality, long context, and next generation agentic capa-
bilities. *arXiv preprint arXiv:2507.06261*, 2025.

- 486 Francesco D’Angelo, Maksym Andriushchenko, Aditya Vardhan Varre, and Nicolas Flammarion.
487 Why do we need weight decay in modern deep learning? *Advances in Neural Information Pro-*
488 *cessing Systems*, 37:23191–23223, 2024.
- 489
490 Shibhansh Dohare, J Fernando Hernandez-Garcia, Qingfeng Lan, Parash Rahman, A Rupam Mah-
491 mood, and Richard S Sutton. Loss of plasticity in deep continual learning. *Nature*, 632(8026):
492 768–774, 2024.
- 493 Mohamed Elsayed and A Rupam Mahmood. Addressing loss of plasticity and catastrophic forget-
494 ting in continual learning. *ICLR*, 2024.
- 495
496 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad
497 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd
498 of models. *arXiv preprint arXiv:2407.21783*, 2024.
- 499 Alex Hägele, Elie Bakouch, Atli Kosson, Leandro Von Werra, Martin Jaggi, et al. Scaling laws and
500 compute-optimal training beyond fixed training durations. *NeurIPS*, 2024.
- 501
502 Moritz Hardt, Ben Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic
503 gradient descent. In *International conference on machine learning*, pp. 1225–1234. PMLR, 2016.
- 504 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
505 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*,
506 2021.
- 507
508 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
509 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Train-
510 ing compute-optimal large language models. *NeurIPS*, 2022a.
- 511
512 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
513 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Train-
514 ing compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022b.
- 515
516 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,
517 Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 2022.
- 518
519 Adam Ibrahim, Benjamin Thérien, Kshitij Gupta, Mats L Richter, Quentin Anthony, Timothée
520 Lesort, Eugene Belilovsky, and Irina Rish. Simple and scalable strategies to continually pre-train
521 large language models. *arXiv preprint arXiv:2403.08763*, 2024.
- 522
523 Arthur Jacot, Peter Šúkeník, Zihan Wang, and Marco Mondelli. Wide neural networks trained with
524 weight decay provably exhibit neural collapse. *arXiv preprint arXiv:2410.04887*, 2024.
- 525
526 Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. Pubmedqa: A dataset
527 for biomedical research question answering. *Proceedings of the 2019 conference on empirical*
528 *methods in natural language processing and the 9th international joint conference on natural*
529 *language processing (EMNLP-IJCNLP)*, 2019.
- 530
531 Konwoo Kim, Suhas Kotha, Percy Liang, and Tatsunori Hashimoto. Pre-training under infinite
532 compute. *arXiv preprint arXiv:2509.14786*, 2025.
- 533
534 James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, et al. Overcoming catastrophic forgetting in
535 neural networks. *Proceedings of the National Academy of Sciences (PNAS)*, 2017. URL <https://arxiv.org/abs/1612.00796>.
- 536
537 Timo Klein, Lukas Mikloutz, Kevin Sidak, Claudia Plant, and Sebastian Tschiatschek. Plasticity
538 loss in deep reinforcement learning: A survey. *arXiv preprint arXiv:2411.04832*, 2024.
- 539
540 Seijin Kobayashi, Yassir Akram, and Johannes Von Oswald. Weight decay induces low-rank atten-
541 tion layers. *NeurIPS*, 2024.
- 542
543 Atli Kosson, Bettina Messmer, and Martin Jaggi. Rotational equilibrium: How weight decay bal-
544 ances learning across neural networks. In *International Conference on Machine Learning*, pp.
545 25333–25369. PMLR, 2024.

- 540 Atli Kosson, Jeremy Welborn, Yang Liu, Martin Jaggi, and Xi Chen. Weight decay may matter more
541 than mup for learning rate transfer in practice. *arXiv preprint arXiv:2510.19093*, 2025.
- 542
- 543 Anders Krogh and John Hertz. A simple weight decay can improve generalization. *Advances in*
544 *neural information processing systems*, 4, 1991.
- 545
- 546 Saurabh Kumar, Henrik Marklund, and Benjamin Van Roy. Maintaining plasticity in continual
547 learning via regenerative regularization. *arXiv preprint arXiv:2308.11958*, 2023.
- 548
- 549 Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. Race: Large-scale reading
550 comprehension dataset from examinations. *arXiv preprint arXiv:1704.04683*, 2017.
- 551
- 552 Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brah-
553 man, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. Tulu 3: Pushing frontiers
554 in open language model post-training. *arXiv preprint arXiv:2411.15124*, 2024.
- 555
- 556 Jae-Hun Lee, Doyoung Yoon, ByeongMoon Ji, Kyungyul Kim, and Sangheum Hwang. Rethinking
557 evaluation protocols of visual representations learned via self-supervised learning. *arXiv preprint*
558 *arXiv:2304.03456*, 2023.
- 559
- 560 Zhiyuan Li, Kaifeng Lyu, and Sanjeev Arora. Reconciling modern deep learning with traditional
561 optimization analyses: The intrinsic learning rate. *Advances in Neural Information Processing*
562 *Systems*, 33:14544–14555, 2020.
- 563
- 564 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,
565 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint*
566 *arXiv:2412.19437*, 2024.
- 567
- 568 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*
569 *arXiv:1711.05101*, 2017.
- 570
- 571 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Confer-*
572 *ence on Learning Representations (ICLR)*, 2019. URL [https://arxiv.org/abs/1711.](https://arxiv.org/abs/1711.05101)
573 05101.
- 574
- 575 Clare Lyle, Zeyu Zheng, Evgenii Nikishin, Bernardo Avila Pires, Razvan Pascanu, and Will Dabney.
576 Understanding plasticity in neural networks. In *International Conference on Machine Learning*,
577 pp. 23190–23211. PMLR, 2023.
- 578
- 579 Prasanna Mayilvahanan, Thaddäus Wiedemer, Sayak Mallick, Matthias Bethge, and Wieland Brendel.
580 Llms on the line: Data determines loss-to-loss scaling laws. *ICML*, 2025.
- 581
- 582 Thomas Miconi, Kenneth Stanley, and Jeff Clune. Differentiable plasticity: training plastic neural
583 networks with backpropagation. In *International Conference on Machine Learning*, pp. 3559–
584 3568. PMLR, 2018.
- 585
- 586 Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi,
587 Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori B Hashimoto. s1: Simple test-
588 time scaling. *Proceedings of the 2025 Conference on Empirical Methods in Natural Language*
589 *Processing*, 2025.
- 590
- 591 Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Ak-
592 shita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, et al. 2 olmo 2 furious. *arXiv*
593 *preprint arXiv:2501.00656*, 2024. URL [https://huggingface.co/allenai/](https://huggingface.co/allenai/OLMo-2-1124-13B/blob/5dbe4046c5ecdee4a93a94cf45728435d5f52695/README.mdx)
OLMo-2-1124-13B/blob/5dbe4046c5ecdee4a93a94cf45728435d5f52695/
README.mdx.
- 594
- 595 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
596 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kel-
597 ton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike,
598 and Ryan Lowe. Training language models to follow instructions with human feedback. *NeurIPS*,
599 2022.

- 594 Samet Oymak, Zalan Fabian, Mingchen Li, and Mahdi Soltanolkotabi. Generalization guaran-
595 tees for neural networks via harnessing the low-rank structure of the jacobian. *arXiv preprint*
596 *arXiv:1906.05392*, 2019.
- 597
- 598 Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. Medmcqa: A large-scale
599 multi-subject multi-choice dataset for medical domain question answering. *Conference on health,*
600 *inference, and learning*, 2022.
- 601
- 602 Guilherme Penedo, Hynek Kydliček, Anton Lozhkov, Margaret Mitchell, Colin A Raffel, Leandro
603 Von Werra, Thomas Wolf, et al. The fineweb datasets: Decanting the web for the finest text data
604 at scale. *Advances in Neural Information Processing Systems*, 37:30811–30849, 2024.
- 605
- 606 Zhenting Qi, Fan Nie, Alexandre Alahi, James Zou, Himabindu Lakkaraju, Yilun Du, Eric Xing,
607 Sham Kakade, and Hanlin Zhang. Evolm: In search of lost language model training dynamics.
NeurIPS, 2025.
- 608
- 609 Matthew Riemer, Ignacio Cases, Robert Ajemian, Miao Liu, Irina Rish, Yuhai Tu, and Gerald
610 Tesauro. Learning to learn without forgetting by maximizing transfer and minimizing interfer-
611 ence. *arXiv preprint arXiv:1810.11910*, 2018.
- 612
- 613 Aaditya Singh, Adam Fry, Adam Perelman, Adam Tart, Adi Ganesh, Ahmed El-Kishky, Aidan
614 McLaughlin, Aiden Low, AJ Ostrow, Akhila Ananthram, et al. Openai gpt-5 system card. *arXiv*
preprint arXiv:2601.03267, 2025.
- 615
- 616 Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng,
617 and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment
618 treebank. *Proceedings of the 2013 conference on empirical methods in natural language process-*
ing, 2013.
- 619
- 620 Tao Sun, Yuhao Huang, Li Shen, Kele Xu, and Bao Wang. Investigating the role of weight decay
621 in enhancing nonconvex sgd. In *Proceedings of the Computer Vision and Pattern Recognition*
622 *Conference*, pp. 15287–15296, 2025.
- 623
- 624 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-
625 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-
626 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 627
- 628 Joshua Vendrow, Edward Vendrow, Sara Beery, and Aleksander Madry. Do large language model
629 benchmarks test reliability?, 2025. URL <https://arxiv.org/abs/2502.03461>.
- 630
- 631 Xi Wang and Laurence Aitchison. How to set adamw’s weight decay as you scale model and dataset
632 size. *arXiv preprint arXiv:2405.13698*, 2024.
- 633
- 634 Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming
635 Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al. Mmlu-pro: A more robust and challenging multi-
636 task language understanding benchmark. *Advances in Neural Information Processing Systems*,
637 2024.
- 638
- 639 Kaiyue Wen, Xingyu Dang, Kaifeng Lyu, Tengyu Ma, and Percy Liang. Fantastic pretraining opti-
640 mizers and where to find them ii: From weight decay to hyperball optimization, 11 2025. URL
641 https://whenwen.github.io/wd_blog/public/hyperball-part-1.html.
- 642
- 643 Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhen-
644 guo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions
645 for large language models. *arXiv preprint arXiv:2309.12284*, 2023.
- 646
- 647 Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding
deep learning requires rethinking generalization. In *International Conference on Learning Rep-*
resentations, 2017.
- 648
- 649 Guodong Zhang, Chaoqi Wang, Bowen Xu, and Roger Grosse. Three mechanisms of weight decay
650 regularization. *arXiv preprint arXiv:1810.12281*, 2018.

648 Hanlin Zhang, Depen Morwani, Nikhil Vyas, Jingfeng Wu, Difan Zou, Udaya Ghai, Dean Foster,
649 and Sham M Kakade. How does critical batch size scale in pre-training? In *The Thirteenth*
650 *International Conference on Learning Representations*, 2025.

651 Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text clas-
652 sification. *Advances in neural information processing systems*, 28, 2015.

653 Pan Zhou, Xingyu Xie, Zhouchen Lin, and Shuicheng Yan. Towards understanding convergence
654 and generalization of adamw. *IEEE transactions on pattern analysis and machine intelligence*,
655 46(9):6486–6493, 2024.

657 658 APPENDIX

659 You may include other additional sections here.

660 661 A MORE DETAILED DISCUSSION OF SECTIONS IN THE MAIN PAPER

662 663 A.1 RELATED WORK

664 Here, we discuss related work on weight decay and model plasticity and how our work contributes
665 new insights.

666 **Weight decay in language model training.** Weight decay is a standard hyperparameter in lan-
667 guage model training and is commonly implemented in conjunction with adaptive optimizers such
668 as AdamW (Loshchilov & Hutter, 2017; Brown et al., 2020a; Grattafiori et al., 2024; OLMo et al.,
669 2024; Liu et al., 2024). Beyond its classical role in regularization and generalization (Krogh &
670 Hertz, 1991; Zhang et al., 2018; Loshchilov & Hutter, 2017; Zhou et al., 2024), weight decay has
671 also been shown to play other roles in language model training, such as improving optimization and
672 training stability (D’Angelo et al., 2024), shaping the learning rate (Li et al., 2020; Kosson et al.,
673 2024; 2025), controlling the effective step size (Wen et al., 2025), inducing low-rank attention lay-
674 ers (Kobayashi et al., 2024), and increasing forgetting of contaminated benchmark questions (Bordt
675 et al., 2025). Wang & Aitchison (2024) show that the weights of AdamW can be understood as an
676 exponential moving average, and that the weight decay parameter plays a critical role in controlling
677 its time scale. Bergsma et al. (2025) study how to set weight decay to minimize the pretraining loss
678 of language models, finding that lower weight decay improves pretraining loss in the overtrained
679 (high TPP ratio) regime. Kim et al. (2025) show that larger weight decay improves pretraining loss
680 in the multi-epoch setting. In contrast to previous work which primarily focuses on weight decay’s
681 effects on the pretrained model, this paper examines how weight decay during pretraining shapes
682 model plasticity.

683 **Plasticity of deep learning models.** Model plasticity has previously been studied in the contexts
684 of continual learning, transfer learning, and reinforcement learning, settings in which models often
685 undergo multiple rounds of training (Dohare et al., 2024; Klein et al., 2024; Coetzer et al., 2025).
686 Prior works have demonstrated that image models lose plasticity when subjected to additional rounds
687 of training on new data, leading to a decreased ability to learn this new data (Dohare et al., 2024; Lyle
688 et al., 2023; Klein et al., 2024). Various approaches have been developed to improve model plasticity,
689 including shrinking and perturbing model weights at the start of each training round (Ash & Adams,
690 2020), identifying and re-initializing less-useful model weights during training (Dohare et al., 2024),
691 pushing weights towards initialization weights during training (Kumar et al., 2023), and learning per-
692 connection plasticity strengths among neuron pairs (Miconi et al., 2018). While previous studies
693 have examined how active forgetting and tokenization (Chen et al., 2023; Abagyan et al., 2025)
694 affect language model plasticity, research on language model plasticity remains underdeveloped. In
695 contrast to these works, this paper investigates the role of weight decay, a standard hyperparameter
696 for language model training, on language model plasticity.

697 698 A.2 BACKGROUND AND METHODS

699 **Weight decay in AdamW.** Motivated by prior findings that regularization helps vision models main-
700 tain plasticity (Dohare et al., 2024), this paper investigates weight decay’s role in language model
701

702 plasticity. We focus on the weight decay hyperparameter λ in the AdamW optimizer which, for each
 703 optimizer step $t \geq 1$, performs two decoupled updates: a gradient update given by
 704

$$705 \theta_t = \hat{\theta}_t - \gamma_t \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon) \quad (2)$$

706 followed by a weight decay update given by
 707

$$708 \hat{\theta}_t = \theta_{t-1} - \gamma_t \lambda \theta_{t-1} \quad (3)$$

709 based on model parameters θ_t , learning rate γ_t , and first- and second-order moment estimates of
 710 the gradient \hat{m}_t and \hat{v}_t (Loshchilov & Hutter, 2019). For language model pretraining, the choice
 711 $\lambda = 0.1$ has emerged as a kind of default, used in many pretraining runs where the optimization
 712 hyperparameters are known (Brown et al., 2020a; Touvron et al., 2023; OLMo et al., 2024).
 713

714 **Language model plasticity.** To assess the plasticity of a pretrained model, we fine-tune the model
 715 on a task and then measure its performance on this task. The better the performance on this down-
 716 stream task, the better the pretrained model was able to learn new data during fine-tuning, thus
 717 the higher the plasticity of the pretrained model. This approach to measuring model plasticity is
 718 consistent with prior literature (Berariu et al., 2021; Dohare et al., 2024).
 719

720 In this context, we now specify the research question:
 721

722 **Research Question.** How does weight decay during language model pretraining affect
 723 model plasticity, i.e. the pretrained model’s ability to learn new knowledge during subse-
 724 quent training?
 725

726 We investigate this research question empirically. We perform experiments that systematically vary
 727 weight decay during pretraining, then fine-tune and evaluate the models’ performance on various
 728 downstream tasks. Our experiments span various model families, model sizes, training regimes
 729 (TPP ratios), and fine-tuning tasks. The setup is as follows.
 730

731 **Pretraining.** We train Llama-2 models on the FineWeb-Edu dataset (Penedo et al., 2024) and
 732 OLMo-2 models on the OLMo-Mix-1124 dataset. We vary model size and TPP ratio, training
 733 models at the 20 TPP Chinchilla-optimal ratio (Hoffmann et al., 2022a) and at the 140 TPP over-
 734 trained ratio. This setup yields five models: Llama-2-0.5B-20x, Llama-2-1B-20x, Llama-2-4B-20x,
 735 OLMo-2-1B-20x, and OLMo-2-1B-140x. For each model, we pretrain variants with different weight
 736 decay.

737 **Fine-tuning.** We perform supervised fine-tuning (SFT) of the pretrained models across six CoT
 738 tasks spanning various domains: MetaMathQA (math reasoning), MedMCQA (medical reason-
 739 ing), PubMedQA (biomedical research), MMLUProCoT (general knowledge and reasoning), RACE
 740 (reading comprehension), and SimpleScaling (math, science, and logical reasoning) (Yu et al., 2023;
 741 Pal et al., 2022; Jin et al., 2019; Wang et al., 2024; Lai et al., 2017; Muennighoff et al., 2025).
 742

743 **Evaluation of model performance after fine-tuning.** We evaluate the fine-tuned models in a zero-
 744 shot manner, prompting them to generate solutions to questions in the fine-tuning test sets, and
 745 assess both the correctness and quality of the solutions using six evaluation metrics.

- 746 • **Greedy** (i.e., **Pass@1**): A single deterministic response is generated (temperature = 0). The
 747 question is marked correct if this response is correct.
- 748 • **Maj@16**, **RM@16**, and **Pass@16**: Sixteen responses are sampled (temperature = 1). The ques-
 749 tion is marked correct if the majority answer is correct (Maj@16), if the response with the highest
 750 ORM score is correct (RM@16), or if any of the responses are correct (Pass@16).
- 751 • **Correct Ratio**: Sixteen responses are sampled (temperature = 1). Among questions with at least
 752 one correct response, we compute the proportion of correct responses out of the sixteen sampled
 753 responses.
- 754 • **ORM Score**: In addition solution correctness, we also measure solution quality. Sixteen responses
 755 are sampled (temperature = 1). Each response is assigned a score using an outcome reward model
 (Skywork-Reward-Llama-3.1-8B-v0.2; ORM) and the average ORM score is computed.

Since these experiments span pretraining and fine-tuning, we adopt the end-to-end analysis framework from Qi et al. (2025). Additional setup details are in Appendix B and C.

Reproducibility statement. Upon publication of the paper, we will release code and models associated with the paper.

B PRE-TRAINING

B.1 MODEL ARCHITECTURES AND TRAINING REGIMES

We pretrain models from different families (Llama-2 and OLMo-2), of different scales (up to 4B), and under different training regimes (20 TPP and 140 TPP), yielding five model setups. Details are in Table 1. For each model setup, we pretrain variants with varying weight decay values.

Table 1: **Model architectures.** We use Llama-2 model architectures from Qi et al. (2025) and OLMo-2 model architecture from OLMo et al. (2024). Llama-2 models are trained at 20 TPP and OLMo-2 models are trained at 20 TPP and 140 TPP.

	Llama-2-0.5B	Llama-2-1B	Llama-2-4B	OLMo-2-1B
Model size	0.5B	1B	4B	1.5B
Hidden size	1536	2048	4096	2048
Intermediate size	3216	4896	7792	16384
Vocab size	32000	32000	32000	100278
Context length	2048	2048	2048	4096
# Heads	32	32	32	16
# Layers	20	22	28	16
# Query groups	4	4	4	16

B.2 TRAINING DETAILS

The training data size (measured in tokens) for each model is determined by the TPP ratio.

Table 2: **Model configurations and training data sizes.**

Model	TPP Ratio	Training Data Size
Llama-2-0.5B-20x	20	10 BT
Llama-2-1B-20x	20	20 BT
Llama-2-4B-20x	20	80 BT
OLMo-2-1B-20x	20	30 BT
OLMo-2-1B-140x	140	210 BT

To pre-train Llama-2 models, we use up to 8 A100 GPUs or 16 H100 GPUs. To pre-train OLMo-2 models, we use 8xH100 GPUs. The 20x OLMo-2-1B models are each trained for 2 days on a single H100 node. The 140x models are trained for 2 weeks on a single H100 node. For all models, we use the AdamW optimizer and standard warmup-cosine learning rate schedule. The only exception is the OLMo-2-1B 140x models, which follow a warmup-stable-decay schedule (Hägele et al., 2024). OLMo-2-1B models are pretrained using the official AllenAI repository.

For each model, we train variants with various weight decay values specified in Table 3. Additional hyperparameters are in Table 4 and 5.

Table 3: **Weight decay values for each model.** For Llama-2-4B-20x, we use the weight decay 0.1 model from Qi et al. (2025) and pre-train the weight decay 1.0 model. For OLMo-2-1B-140x, we use the weight decay 0.1 model from Bordt & Pawelczyk (2025) and pretrain the weight decay 0.3 and 1.0 models.

Model	Weight Decay
Llama-2-0.5B-20x	9 values: {0.0001, 0.001, 0.01, 0.1, 0.5, 1.0, 1.5, 3.0, 10.0}
Llama-2-1B-20x	9 values: {0.0001, 0.001, 0.01, 0.1, 0.5, 1.0, 1.5, 3.0, 10.0}
Llama-2-4B-20x	2 values: {0.1, 1.0}
OLMo-2-1B-20x	3 values: {0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0}
OLMo-2-1B-140x	3 values: {0.1, 0.3, 1.0}

Table 4: **Hyperparameters for Llama-2 model pre-training.** For Llama-2 models, hyperparameter values are chosen following those in Qi et al. (2025), except for weight decay, which is varied as the independent variable in our experiments.

Hyperparameter	Llama-2-0.5B-20x	Llama-2-1B-20x	Llama-2-4B-20x
precision	bf16-mixed	bf16-mixed	bf16-mixed
global_batch_size	512	512	1024
max_seq_length	2048	2048	2048
lr_warmup_ratio	0.1	0.1	0.1
max_norm	1	1	1
lr	0.00025	0.0002	0.00015
min_lr	0.000025	0.00002	0.000015
weight_decay	varies	varies	varies
beta1	0.9	0.9	0.9
beta2	0.95	0.95	0.95
epoch	1	1	1

Table 5: **Hyperparameters for model pre-training.** For OLMo-2 models, hyperparameter values follow the OLMo-2 defaults (OLMo et al., 2024), except for weight decay, which is varied as the independent variable in our experiments.

	Llama-2-4B-20x	OLMo-2-1B-20x	OLMo-2-1B-140x
precision		bf16-mixed	bf16-mixed
global_batch_size		512	512
max_seq_length		4096	4096
lr_warmup_ratio		0.1	0.1
max_norm		1	1
lr		0.0004	0.0004
min_lr		0.00004	0
weight_decay		varies	varies
beta1		0.9	0.9
beta2		0.95	0.95
epoch		1	1

C FINE-TUNING

C.1 FINE-TUNING DATASETS

We clean the fine-tuning training datasets, removing incoherent or questions that exceed the maximum input sequence length of the models. The size of the fine-tuning training set for each task and the test set used to subsequently evaluate model performance are shown in Table 6.

Table 6: **Fine-tuning and evaluation datasets.** MetaMathQA and SimpleScaling are evaluated on test sets of the GSM8KPlatinum (Cobbe et al., 2021; Vendrow et al., 2025) and MATH (Hendrycks et al., 2021) datasets because MetaMathQA and SimpleScaling contain questions that are augmented from the training sets of these two datasets.

Task	Training set	Test set
MetaMathQA	$n = 395,000$	GSM8KPlatinum ($n = 1,209$) + MATH ($n = 5,000$)
MedMCQA	$n = 182,555$	MedMCQA ($n = 4183$)
PubMedQA	$n = 211,168$	PubMedQA ($n = 1000$)
MMLUProCoT	$n = 123,836$	MMLUProCoT ($n = 567$)
RACE	$n = 92,737$	RACE ($n = 4934$)
SimpleScaling	$n = 54,484$	GSM8KPlatinum ($n = 1,209$) + MATH ($n = 5,000$)

C.2 TRAINING DETAILS

Table 7: **Hyperparameters for supervised fine-tuning.** We set `batch_size = 64` due to computational constraints. We set `n_epochs = 3` based on results from Qi et al. (2025) indicating that this setting leads to the best downstream performance. All other hyperparameters are from Qi et al. (2025).

	1B and under	4B
<code>cutoff_len</code>	2048	2048
<code>batch_size</code>	64	64
<code>learning_rate</code>	0.00001	0.0000075
<code>lr_scheduler_type</code>	cosine	cosine
<code>warmup_ratio</code>	0.1	0.1
<code>n_epochs</code>	3	3

C.3 TEMPLATE

We use the following template for supervised fine-tuning.

Human: {question}
Assistant: {response}

D EVALUATION

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

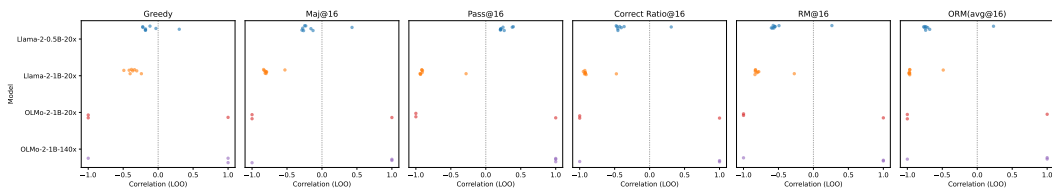


Figure 7: **Stability analysis for Pearson correlation coefficient.** Pearson correlation is computed for each leave-one-out (LOO) subset in Figure 9g. The LOO correlation can change noticeably in magnitude and sign, suggesting that the computed correlation relationship is rather unstable.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

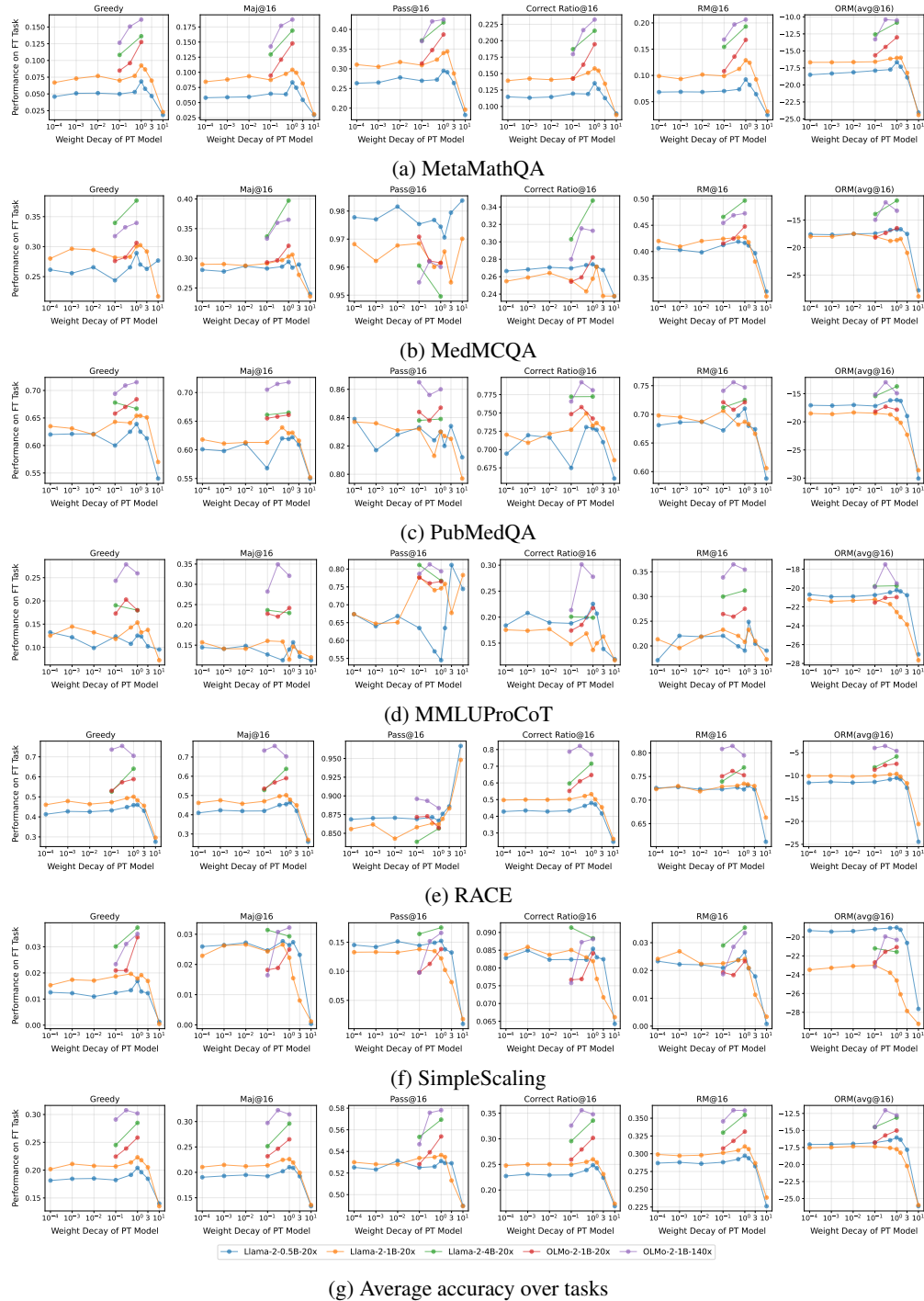


Figure 8: Accuracy of models on each task after fine-tuning measured based on six metrics.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

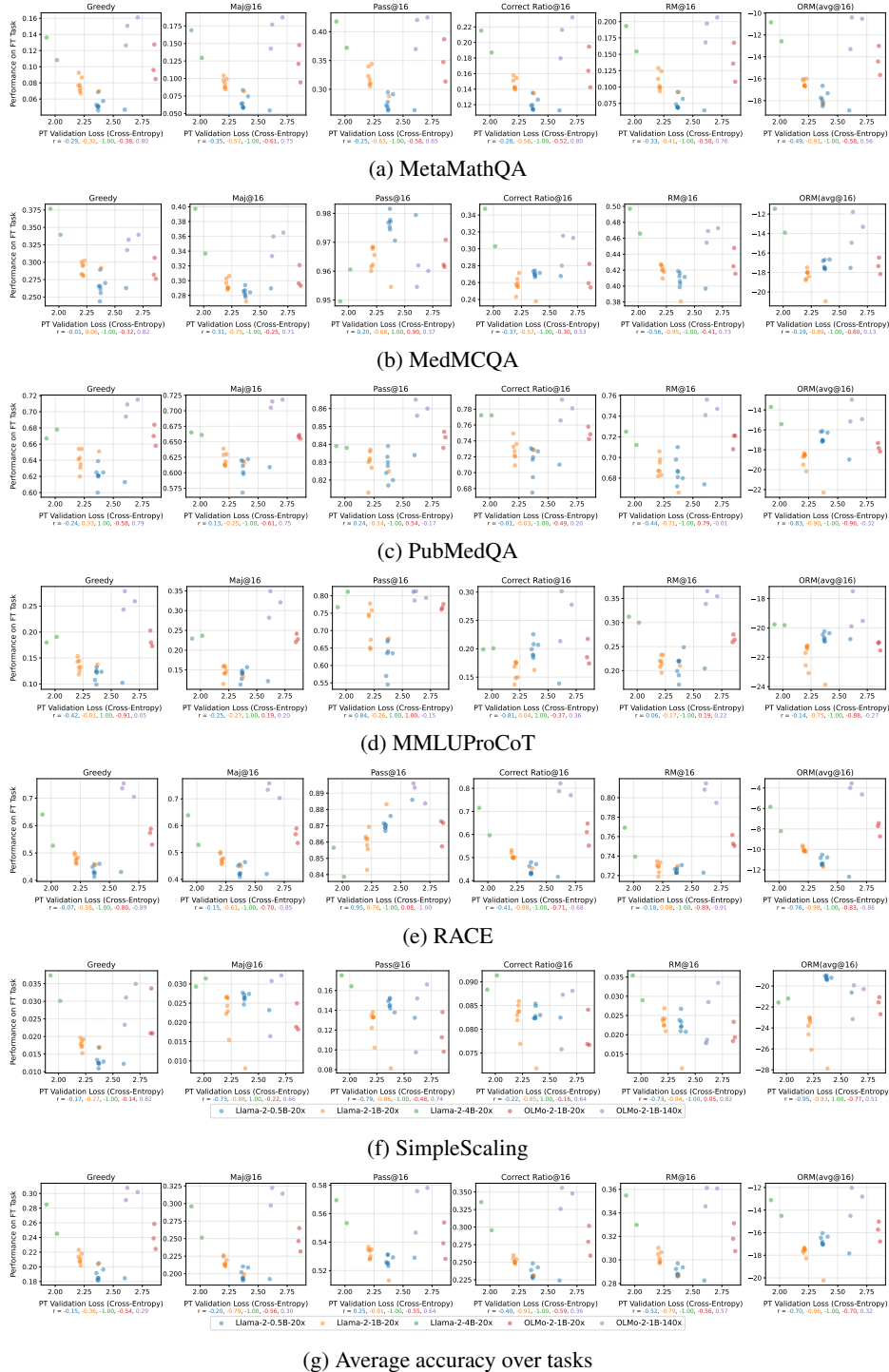
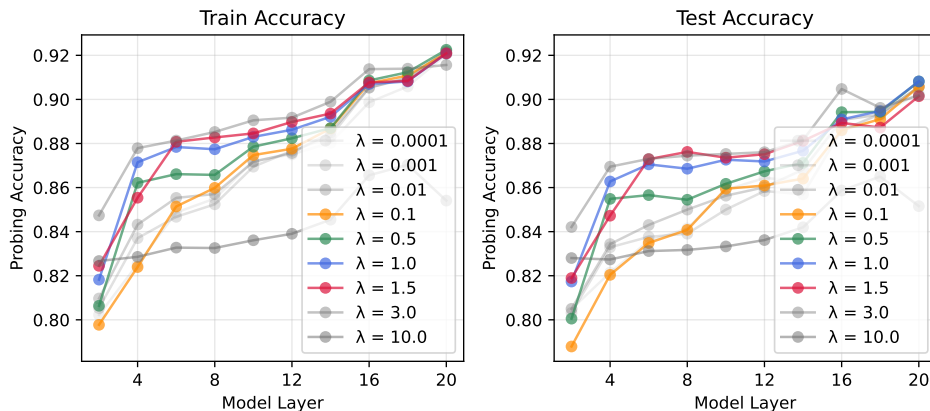


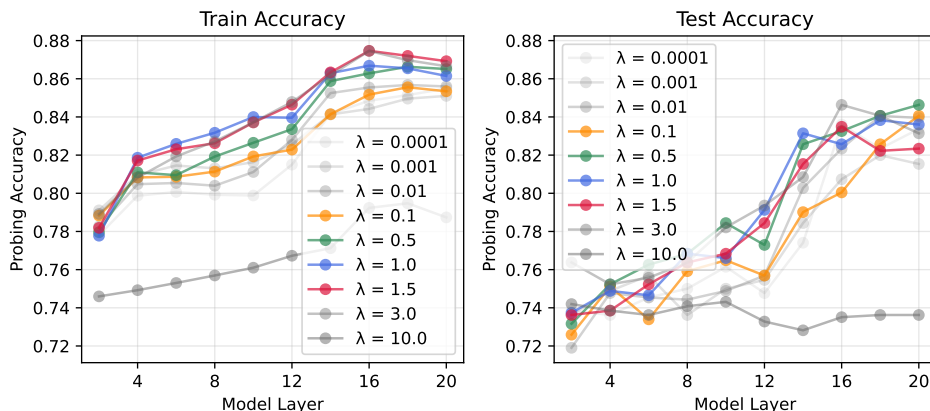
Figure 9: **Loss after model pre-training is not predictive of model performance after fine-tuning.** Figures show model performance on individual datasets after fine-tuning measured based on six datasets.

E ANALYSES ON WEIGHT DECAY’S MECHANISTIC EFFECTS ON MODEL BEHAVIOR

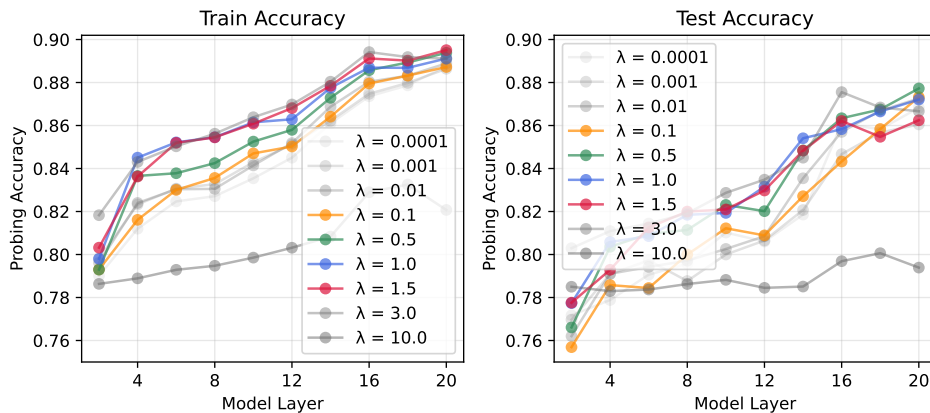
E.1 WEIGHT DECAY’S EFFECT ON MODEL REPRESENTATIONS



(a) AG News



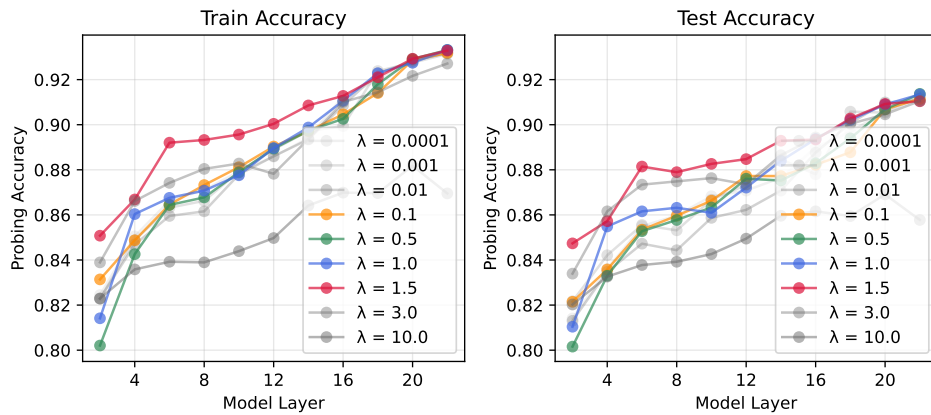
(b) SST



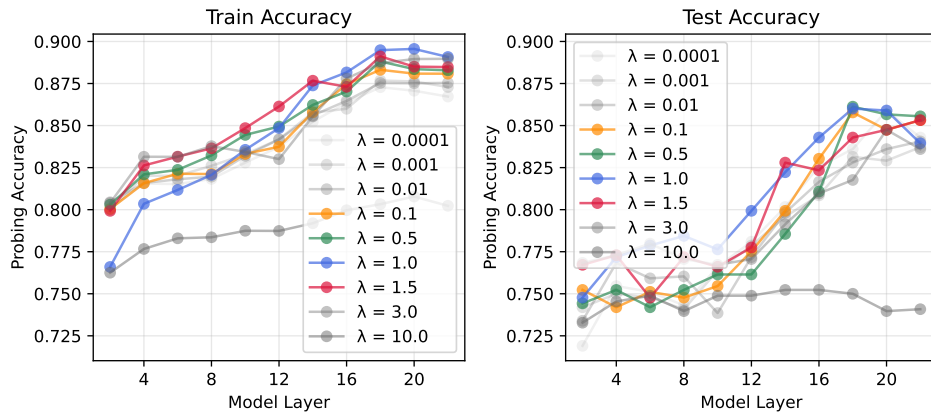
(c) Averaged over tasks

Figure 10: **Linear probing experiments for Llama-2-0.5B-20x.** The train and test accuracies of the linear probes for the SST and AG News datasets and the average train and test accuracy over the two datasets.

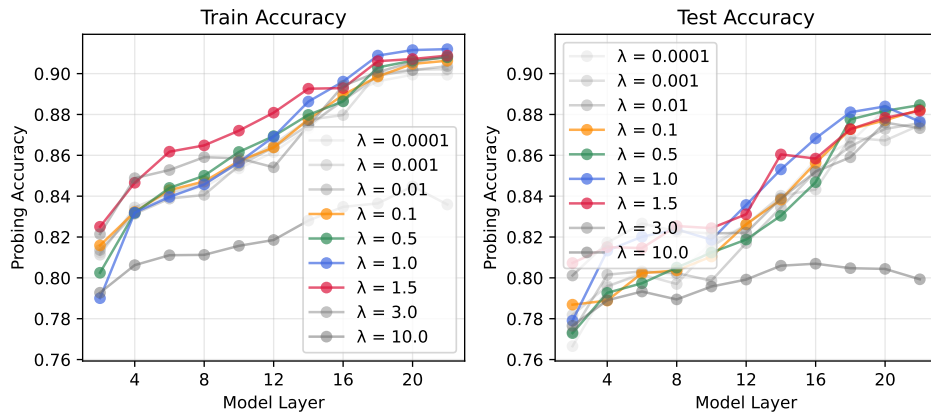
1134
 1135
 1136
 1137
 1138
 1139
 1140
 1141
 1142
 1143
 1144
 1145
 1146
 1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172
 1173
 1174
 1175
 1176
 1177
 1178
 1179
 1180
 1181
 1182
 1183
 1184
 1185
 1186
 1187



(a) AG News



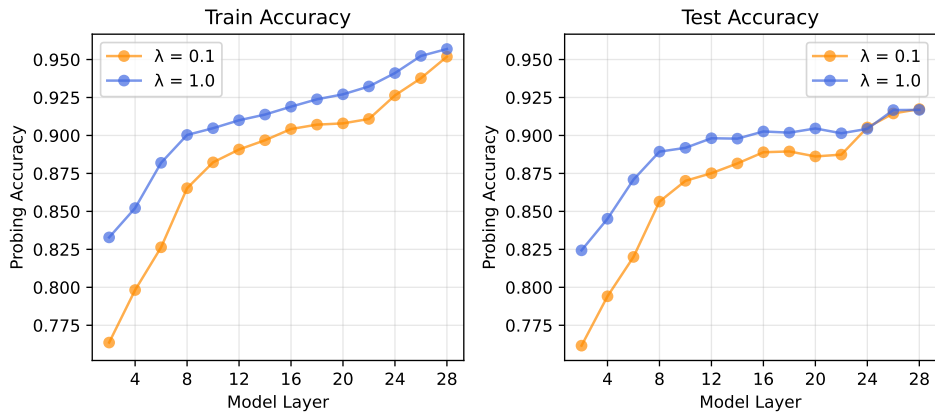
(b) SST



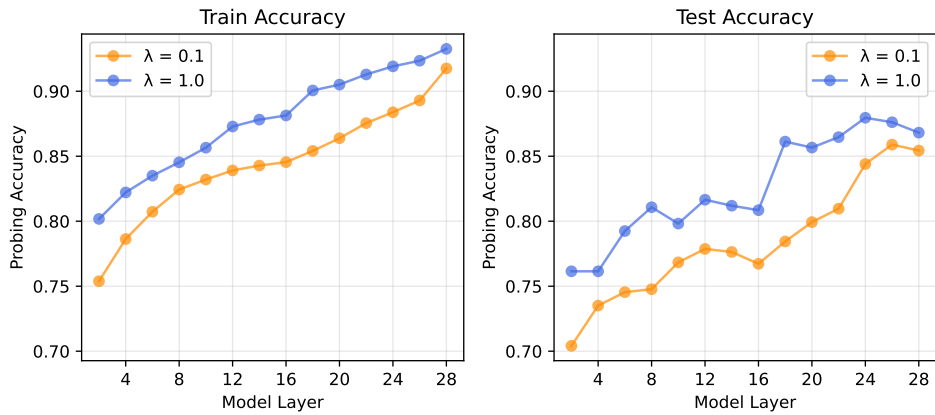
(c) Averaged over tasks

Figure 11: **Linear probing experiments for Llama-2-1B-20x.** The train and test accuracies of the linear probes for the SST and AG News datasets and the average train and test accuracy over the two datasets.

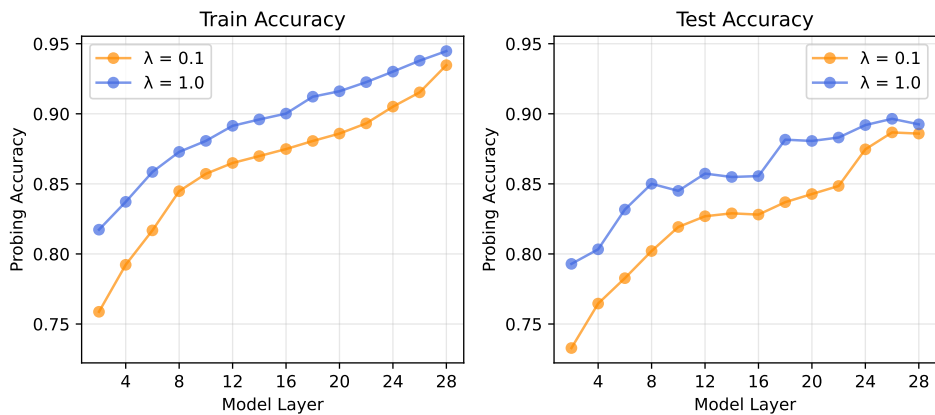
1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241



(a) AG News



(b) SST



(c) Averaged over tasks

Figure 12: **Linear probing experiments for Llama-2-4B-20x.** The train and test accuracies of the linear probes for the SST and AG News datasets and the average train and test accuracy over the two datasets.

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

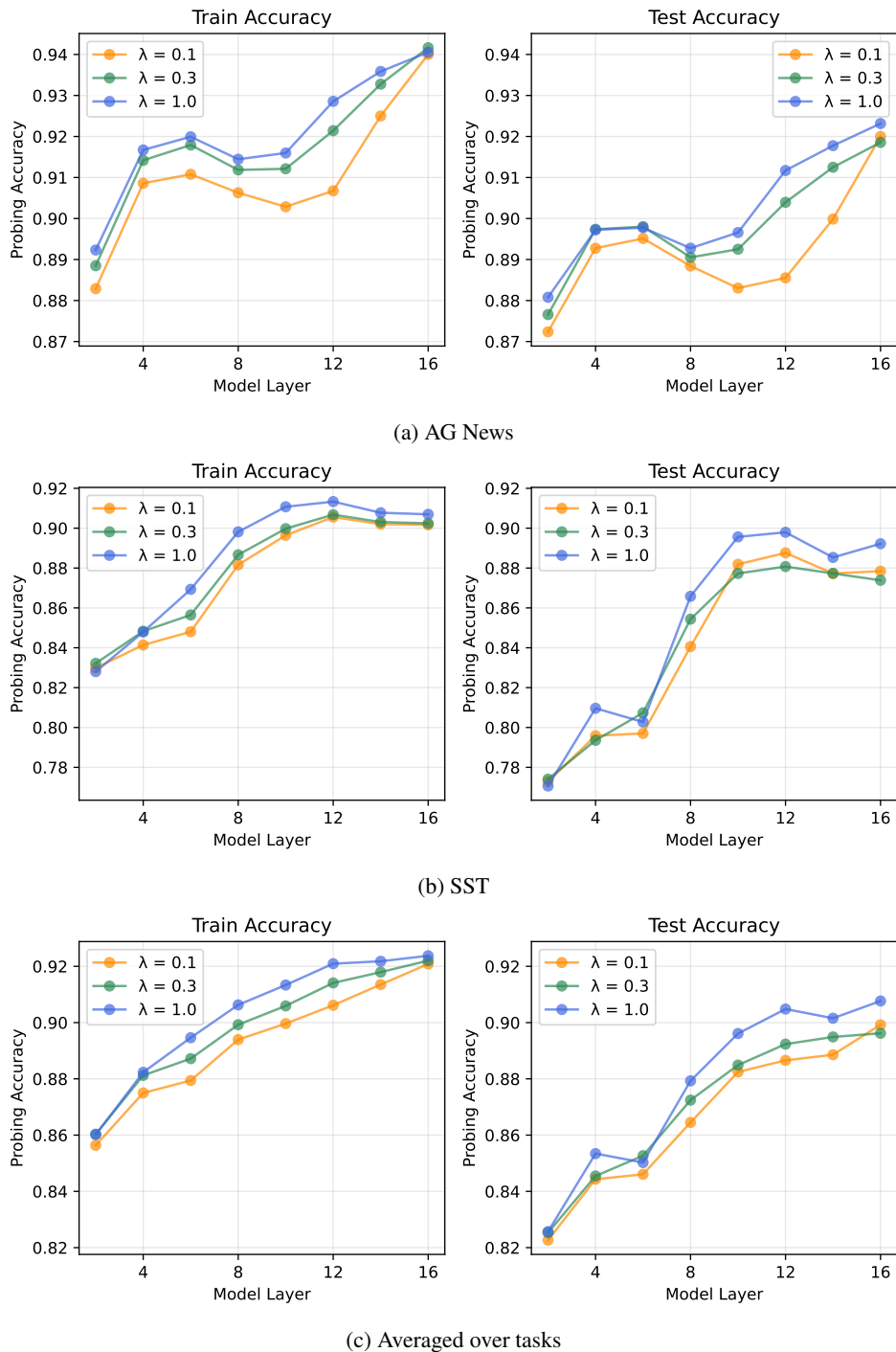
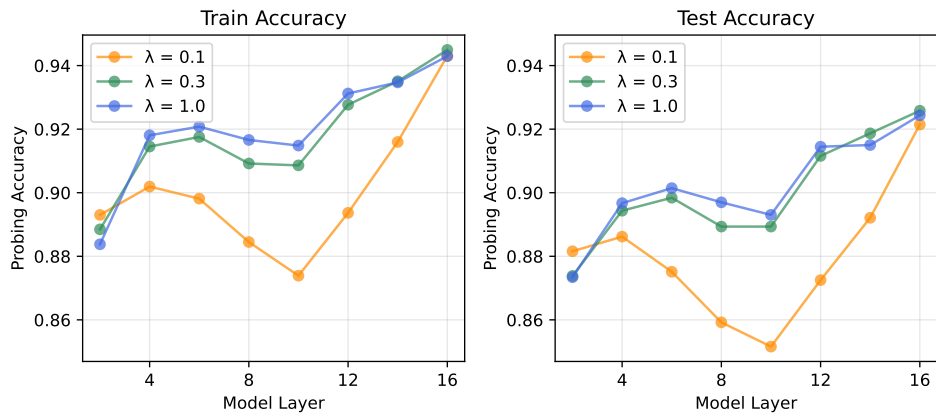
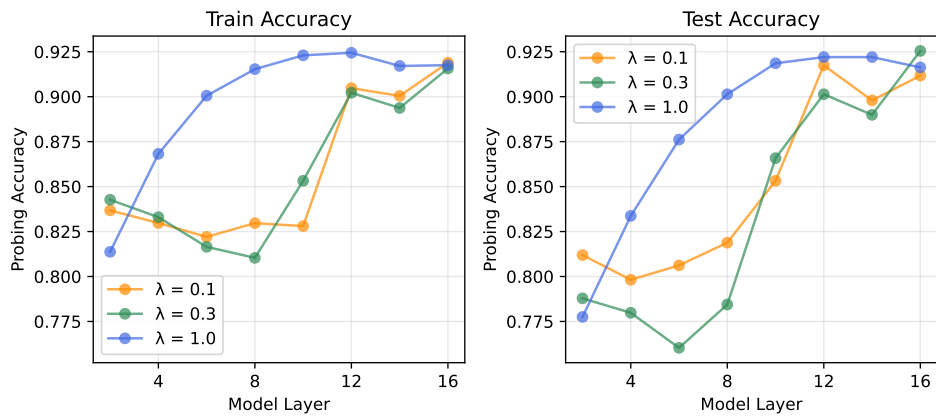


Figure 13: **Linear probing experiments for OLMo-2-1B-20x.** The train and test accuracies of the linear probes for the SST and AG News datasets and the average train and test accuracy over the two datasets.

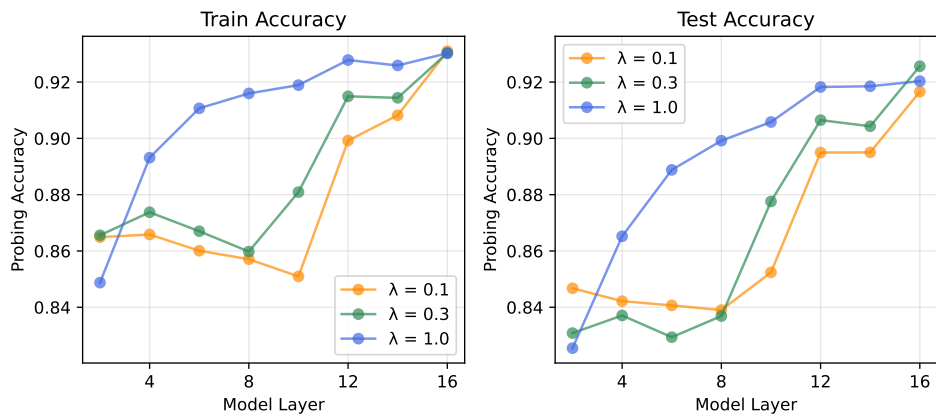
1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349



(a) AG News



(b) SST



(c) Averaged over tasks

Figure 14: **Linear probing experiments for OLMo-2-1B-140x.** The train and test accuracies of the linear probes for the SST and AG News datasets and the average train and test accuracy over the two datasets.

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

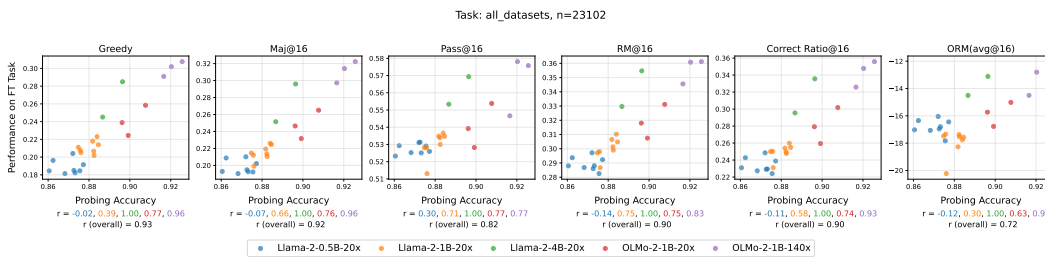


Figure 15: **Probing accuracy is highly predictive of downstream model performance.** The x-axis is the best average probing accuracy of the model (best out of all model layers). The y-axis is the average accuracy of the model over all tasks after fine-tuning. Pretrained models with higher probing accuracies from the linear probing experiments tend to perform better downstream after fine-tuning.

E.2 WEIGHT DECAY’S EFFECT ON ATTENTION MATRIX RANK

E.2.1 ATTENTION PSEUDO-RANK COMPUTATION

To quantify the effective dimensionality of weight matrices, we follow Kobayashi et al. (2024) and compute the pseudo-rank of the matrices. For a matrix W with singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$, the pseudo-rank is defined as the ratio k/n , where k is the smallest integer satisfying:

$$\frac{\sum_{i=1}^k \sigma_i}{\sum_{i=1}^n \sigma_i} \geq 0.95 \tag{4}$$

This metric represents the fraction of the largest singular values required to capture at least 95% of the sum of all singular values. In our analysis, we apply this computation to the product of the key-query matrices ($W_{QK} = W_K^T W_Q$) and the value-projection matrices ($W_{VP} = W_P W_V$) to monitor the emergence of low-rank structures during training.

E.2.2 ADDITIONAL ANALYSES ON ATTENTION MATRIX RANK

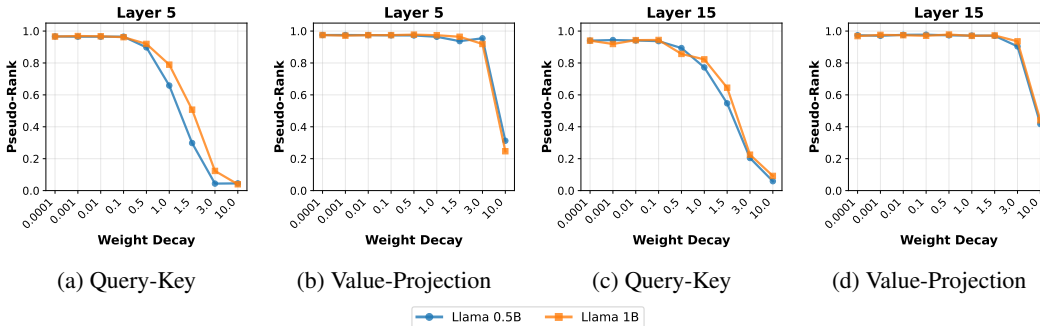


Figure 16: **Weight decay reduces the rank of attention matrices.** The figure depicts the average pseudo-rank (Supplement E.2.1) of the query-key (W_{QK}) and value projection (W_{VP}) matrices in layers 5 and 15 of the fully-trained Llama-2 models at 20 TPP.

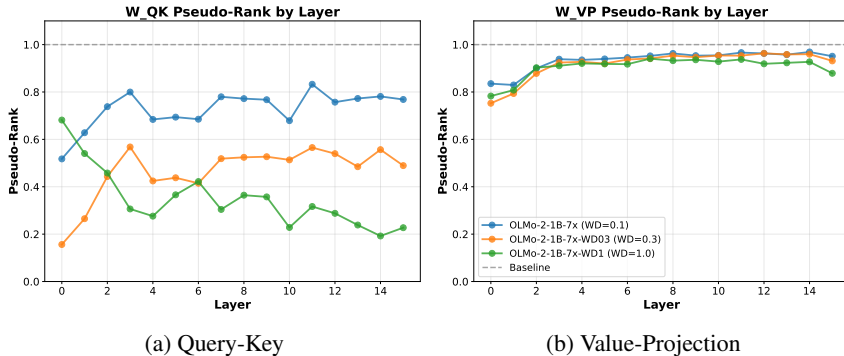


Figure 17: **Weight decay reduces the rank of attention matrices.** This is for the OLMo models trained at 140 TPP. We observe that the rank of attention for weight decay 0.1 is generally smaller than what we observe for both for 20 TPP and for the fully trained OLMo-2-1B-0425 model. Hence, we conjecture that this is because the 140 TPP models were trained with a warmup-stable-decay learning rate schedule, whereas the 1x and 144x models were trained with a cosine learning rate schedule. While it has been shown that WSD leads to a similar validation loss to cosine decay (Hägele et al., 2024), there is emerging evidence that there are important differences between the training dynamics of the two learning rate schedules (Catalan-Tatjer et al., 2025).

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

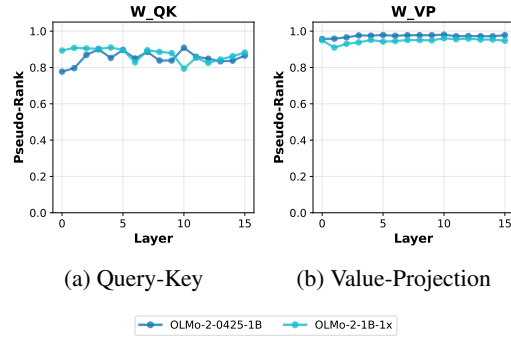


Figure 18: Training time does not reduce the rank of attention matrices.

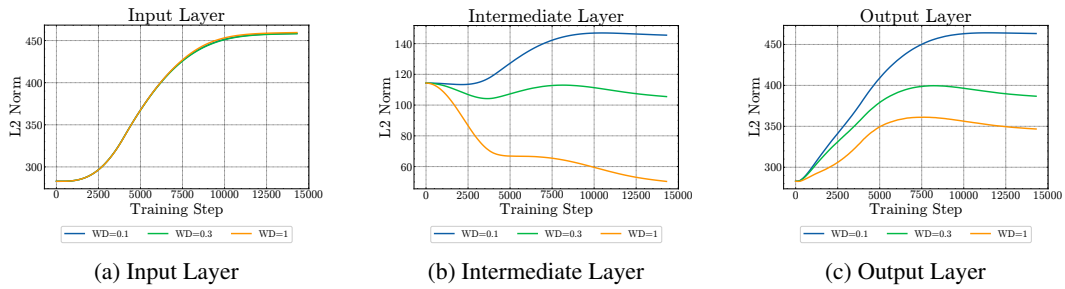


Figure 19: **Weight decay reduces the norm of the weights of the model.** The effect does not occur for the input layer, where the weights are not being decayed. This is for OLMo-2-1B models trained at 20 TPP.