
Mechanistic Insights into Grokking from the Embedding Layer

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Grokking, a delayed generalization in neural networks after perfect training per-
2 formance, has been observed in Transformers and MLPs, but the components
3 driving it remain underexplored. We show that embeddings are central to grokking:
4 introducing them into MLPs induces delayed generalization in modular arithmetic
5 tasks, whereas MLPs without embeddings can generalize immediately. Our analy-
6 sis identifies two key mechanisms: (1) Embedding update dynamics, where rare
7 tokens stagnate due to sparse gradient updates and weight decay, and (2) Bilinear
8 coupling, where the interaction between embeddings and downstream weights
9 introduces saddle points and increases sensitivity to initialization. To confirm
10 these mechanisms, we investigate frequency-aware sampling, which balances token
11 updates by minimizing gradient variance, and embedding-specific learning rates,
12 derived from the asymmetric curvature of the bilinear loss landscape. We prove
13 that an adaptive learning rate ratio, $\frac{\eta_E}{\eta_W} \propto \frac{\sigma_{\max}(E)}{\sigma_{\max}(W)} \cdot \frac{f_W}{f_E}$, mitigates bilinear cou-
14 pling effects, accelerating convergence. Our methods not only improve grokking
15 dynamics but also extend to broader challenges in Transformer optimization, where
16 bilinear interactions hinder efficient training.

17 1 Introduction

18 The phenomenon of grokking, in which a neural network exhibits delayed generalization after
19 achieving close to or perfect training performance, has emerged as a compelling topic in deep learning.
20 Initially observed in Transformer architectures by [19], grokking presents a puzzling challenge
21 where models that seem to overfit to training data eventually demonstrate remarkable generalization
22 capabilities after extensive training. Subsequent research has identified this phenomenon across
23 various architectures, including convolutional neural networks (CNNs) and multi-layer perceptrons
24 (MLPs) [13, 12]. Despite growing interest, the underlying mechanisms of grokking remain elusive.

25 Existing studies have sought to unravel grokking by exploring its connection to delayed robustness,
26 local complexity, and model architecture [3, 6]. For instance, [6] suggest that grokking coincides with
27 a phase transition in the linear regions of a model’s input space, leading to robust partitions that enable
28 generalization after extended training. Others have attributed grokking to emergent circuit behaviors
29 or optimization dynamics [17, 21]. However, these studies often focus on high-level phenomena,
30 overlooking the role of specific components, such as embedding layers, in shaping the dynamics of
31 grokking.

32 In this work, we argue that embedding layers are central to understanding the grokking phenomenon.
33 By introducing embedding layers into MLP architectures, we observe clear grokking patterns even in
34 simple modular arithmetic tasks, such as modular addition. Interestingly, MLPs without embedding
35 layers can often generalize without grokking, suggesting that embeddings introduce unique dynamics
36 that delay generalization. Our analysis identifies two critical factors that influence these dynamics:

- 37 1. **Embedding update dynamics:** Embedding parameters are updated through gradient de-
38 scent and weight decay. However, embeddings corresponding to tokens not present in a
39 given batch are updated solely via weight decay or residual effects from previous gradi-
40 ents in optimizers like Adam. This imbalance delays stabilization and can hinder training,
41 particularly for low-probability tokens.
- 42 2. **Coupling with the first-layer weights:** When embeddings are multiplied with the weights
43 of the first layer, they form a bilinear interaction. This coupling introduces structural
44 complexity into the optimization landscape, making the process more susceptible to saddle
45 points and increasing the sensitivity to initialization.

46 Building on these insights, we propose two strategies to address and prove the hypotheses introduced
47 for embedding layers. **First:** A refined sampling methodology that ensures more uniform updates
48 across all embeddings, mitigating frequency imbalance. **Second:** A learning rate adjustment for
49 embeddings, setting it higher than that of the rest of the model. This adjustment counteracts the
50 coupling effect with the first-layer weights, enabling faster stabilization and reducing the risk of
51 optimization stagnation. Our experiments demonstrate that these strategies not only accelerate the
52 grokking process but also enable generalization in scenarios where traditional approaches fail.

53 Additionally, the bilinear coupling observed in embedding-based MLPs highlights broader challenges
54 in optimizing Transformer architectures. Transformers, which rely on multiplicative interactions in
55 attention mechanisms, exhibit similar issues due to the bilinearity of query, key, and value projections.
56 While softmax attention and scaling by the dimensionality d help smooth the optimization landscape,
57 these mechanisms may still struggle with increased saddle points in certain layers [5]. In summary,
58 this work contributes to the understanding of grokking and its broader implications for deep learning
59 by:

- 60 • Highlighting the unique role of embedding layers in delaying generalization and their
61 coupling with the first layer in MLPs.
- 62 • Proposing strategies to accelerate grokking, including refined sampling and embedding-
63 specific learning rates.
- 64 • Connecting the challenges in embedding-based optimization to broader issues in Transformer
65 training, such as bilinearity, saddle points, and the effectiveness of adaptive optimizers like
66 Adam.

67 By bridging insights from grokking and Transformer optimization, we provide a unified perspective
68 on the interplay between embedding dynamics, optimization challenges, and generalization.

69 2 Related Work

70 The phenomenon of grokking, where generalization emerges abruptly after prolonged overfitting, was
71 first observed in transformers [19] and later extended to CNNs and ResNets [13, 12], indicating it is
72 architecture-agnostic. Various explanations have been proposed. [7] attribute it to phase transitions in
73 local complexity (“delayed robustness”), while others link it to circuit efficiency [17, 21, 11]. Though
74 insightful, these perspectives don’t fully explain the delayed generalization. Connections to double
75 descent have also been explored [1, 16], but grokking’s dynamics remain distinct.

76 The closest work to ours studies modular addition using permutation-equivariant models [15], where
77 one-hot inputs interact with the first layer as a fixed embedding. Their analysis, however, is limited to
78 modular tasks and specific activations. In contrast, we generalize across datasets and highlight how
79 embedding layers, especially when trainable, interact bilinearly with downstream weights, affecting
80 optimization dynamics.

81 Related studies like Tensor Programs IV [24] prescribe per-layer scaling based on width, assuming
82 independent layer evolution. Our setup differs: the embedding layer’s updates depend on both its
83 own width and the spectrum of the coupled layer. Prieto et al. [20] connect delayed generalization to
84 numerical instability (Softmax Collapse), proposing solutions that complement our focus on structural
85 coupling and gradient imbalance.

86 Unlike works that focus on final representations [4], we analyze the embedding layer’s evolving role
87 during training. Even with one-hot inputs, its interaction with the first linear layer forms a learnable

88 embedding mechanism. Concurrent work shows that transferring embeddings from small to large
 89 models can accelerate grokking [23]; while we share this motivation, we also observe in preliminary
 90 trials that transferring other MLP layers may offer similar benefits.

91 Finally, the bilinear coupling we analyze in MLPs parallels challenges in Transformer architectures,
 92 where attention mechanisms introduce similar multiplicative dynamics. Prior work highlights how
 93 adaptive optimizers like Adam outperform SGD due to gradient noise and curvature heterogeneity
 94 [25, 10, 26]. Our findings help bridge these perspectives by showing how embedding-layer coupling
 95 shapes optimization and generalization.

96 3 Preliminaries

97 3.1 Embedding Layers

98 The Transformer model [22] utilizes a self-attention
 99 mechanism to capture dependencies between tokens. In
 100 this framework, embeddings map input tokens to high-
 101 dimensional vectors, which are processed through atten-
 102 tion layers. These embeddings help the model capture
 103 contextualized representations. In contrast, MLPs rely on
 104 fully connected layers without attention mechanisms. We
 105 investigate the role of embeddings in MLPs, specifically
 106 how they improve model generalization. The core contri-
 107 bution of this work is to examine the role of embedding
 108 layers in MLPs. These layers map discrete tokens to dense,
 109 high-dimensional vectors, enabling models to handle non-
 110 linear tasks like modular arithmetic. Even with one-hot
 111 inputs—as studied in theoretical settings [2, 15]—the first
 112 weight matrix effectively functions as a learned embed-
 113 ding. Thus, embeddings, whether explicit or implicit, play
 114 a central role in shaping model dynamics. While com-
 115 monly associated with Transformers, we focus on MLPs
 116 as a simpler and more interpretable setting. MLPs avoid
 117 the added complexity of self-attention while still exhibit-
 118 ing phenomena like grokking. Importantly, the bilinear
 119 coupling between embeddings and downstream weights,
 120 central to our analysis, also arises in Transformers but
 121 is further complicated by attention. Studying MLPs al-
 122 lows us to isolate and understand this coupling in a clean,
 123 controlled environment.

124 3.2 Algorithmic Datasets and Modular Arithmetic

125 Algorithmic datasets are synthetic datasets carefully con-
 126 structed with controlled mathematical properties, typically
 127 involving operations over finite sets such as modular addi-
 128 tion or multiplication. One well-known example is the
 129 modular arithmetic dataset studied by [19], where the goal is to uncover relationships between binary
 130 inputs and produce consistent outputs based on these operations. For instance, given inputs a and
 131 b , the model is tasked to compute $(a + b) \bmod P$ or $(a \times b) \bmod P$, where P is a prime number, and
 132 both inputs and outputs are constrained within $\{0, 1, \dots, P - 1\}$ (refer to Figure 1).

133 This dataset highlights the challenging nature of generalization in grokking: the relationship between
 134 inputs is defined purely by a deterministic operation, not by a probabilistic distribution. Unlike
 135 typical machine learning datasets, where examples are drawn from an underlying (often unknown)
 136 data distribution, algorithmic datasets consist of a finite and complete set of all possible input-output
 137 combinations. In such cases, there is no statistical "distribution" in the conventional sense; instead, the
 138 generalization task relies on uncovering the underlying relationship between inputs, which demands a
 139 model to internalize the algorithm itself. Moreover, any hypothesis consistent with training examples
 140 can initially seem plausible from a statistical perspective, as no known distribution governs the data.

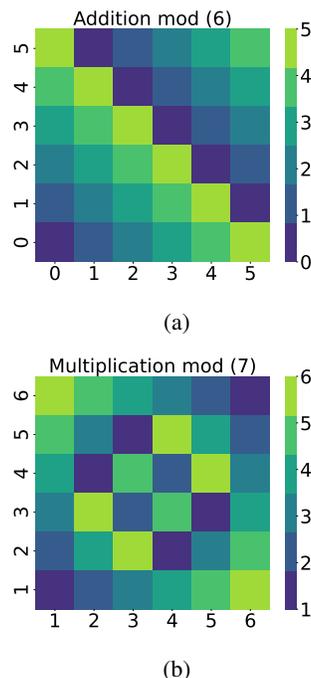


Figure 1: Heatmaps for (a) additive group (mod 6) and (b) multiplicative group (mod 7). The two groups are isomorphic despite differing appearances.

141 The difficulty of generalization thus lies not in interpolating unseen samples but in discovering the
 142 underlying relation, making it a fundamentally different task.

143 We note that there is an equivalence between modular addition and modular multiplication in certain
 144 settings. Namely, given a prime number p , the groups (in mathematical sense) of modular addition
 145 $(\{0, 1, \dots, p - 2\}, +)$ (where addition is performed modulo $p - 1$), and of modular multiplication
 146 $(\{1, \dots, p - 1\}, *)$ (where multiplication is performed modulo p) are isomorphic. Both groups have
 147 the same number of elements (which is $p - 1$), and are simple (meaning, there is an element g , called
 148 generator, such that every other element is of the form $g * \dots * g$, where $*$ is the group operation and
 149 the number of operations used is less than p . In the first group, any element different from 0 is the
 150 group generator while in the second group, any element different from 1 is the generator (see Figure
 151 1).

152 The embedding layer strips the input group elements of their numerical meanings, and assigns a
 153 general, abstract vector to each element. In this way, training on modular addition or multiplication
 154 presents no difference for MLP (or other architectures) with the embedding layer. In contrast to
 155 this, the MLP without the embedding layer is able to fit and generalize on modular addition, while it
 156 completely fails on modular multiplication.

157 3.3 Problem Setup and Motivations

158 Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ represent an algorithmic dataset, where each x_i is an input token sequence
 159 (e.g., $a, b, \text{operation}, =$), and y_i is the output derived from an operation modulo a positive integer P .
 160 The task is to learn a mapping $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ parameterized by θ , capable of generalizing to unseen
 161 samples from $\mathcal{D}_{\text{test}}$.

162 To process inputs effectively, we tokenize them as sequences of their digit representations, as the
 163 model does not inherently interpret numerical values. Each operand a and b is assigned a token in the
 164 range 0 to $P - 1$, while the operation and equality symbols are represented by tokens P and $P + 1$,
 165 respectively. For instance, the modular arithmetic expression $(3 + 2)(\text{mod } 5) = 0$ is tokenized as
 166 $[3, 5, 2, 6, 0]$.

167 Embedding layers in models provide a dense representation of tokens. However, delayed updates to
 168 embeddings for infrequent tokens can significantly impact convergence and generalization. Our work
 169 explores these dynamics, with a focus on the impact of p_i , the i^{th} -token sampling probability, and
 170 proposes adjustments to improve convergence. We investigate the use of embeddings in MLPs for
 171 algorithmic tasks. We started by training a MLP on modular addition and multiplication datasets,
 172 comparing setups with and without embedding layers.

173 **MLP Without Embeddings.** In this setup, input tokens ($a, b, \text{operation } (P)$, and equality sign
 174 $(P + 1)$) are encoded directly into a 4-dimensional input vector. The MLP processes these inputs as:

$$\begin{aligned} \mathbf{h}_1 &= \sigma(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1), & \mathbf{h}_2 &= \mathbf{W}_2 \mathbf{h}_1 + \mathbf{b}_2, \\ \hat{\mathbf{y}} &= \text{Softmax}(\mathbf{h}_2). \end{aligned} \tag{1}$$

175 where $\mathbf{x} \in \mathbb{R}^4$ is the encoded input vector (with first and third entry a and b , respectively), $\mathbf{W}_1, \mathbf{W}_2$
 176 are weight matrices, $\mathbf{b}_1, \mathbf{b}_2$ are biases, σ is the ReLU activation function, and $\hat{\mathbf{y}}$ represents the
 177 predicted output.

178 This configuration demonstrates that the MLP can fit the addition task with ease, but struggles to
 179 generalize multiplication. This difficulty arises because multiplication modulo P is not linearly
 180 separable, as evident in the non-trivial patterns in Figure 1.

181 **MLP With Embeddings.** To overcome the challenges of non-linear separability, we introduced
 182 an embedding layer. Each token x is mapped to a dense vector e_x through an embedding matrix
 183 $\mathbf{E} \in \mathbb{R}^{V \times d}$, where d is the embedding dimension. Our input consists of 4 token embeddings of the
 184 form $\hat{\mathbf{e}} = [e_i, e_{*'}, e_k, e_{='}]^\top$, and the modified forward pass is:

$$\begin{aligned} \mathbf{h}_1 &= \sigma(\mathbf{W} \hat{\mathbf{e}} + \mathbf{b}_1), \\ \mathbf{h}_2 &= \mathbf{W}_2 \mathbf{h}_1 + \mathbf{b}_2, & \hat{\mathbf{y}} &= \text{Softmax}(\mathbf{h}_2), \end{aligned} \tag{2}$$

185 Adding embeddings allows the model to capture more expressive input representations. With this
 186 setup, we observed that the model generalized well to both addition and multiplication tasks, but with
 187 a delayed generalization for multiplication. This delay corresponds to the grokking phenomenon,
 188 which appears as a "trapezoid pattern" in performance plots: a phase of memorization followed by a
 189 sudden leap in test accuracy, as illustrated in figure 2 .

190 These observations motivate a deeper analysis of embedding dynamics during training. In particular,
 191 we investigated the gradient heatmaps to understand the role of embeddings in delaying generalization.
 192 By visualizing gradient magnitudes across training epochs, we point out that embeddings receive
 193 smaller updates compared to other weights of the model, potentially causing grokking. This investi-
 194 gation will help establish a connection between embedding behavior and the observed generalization
 195 delays.

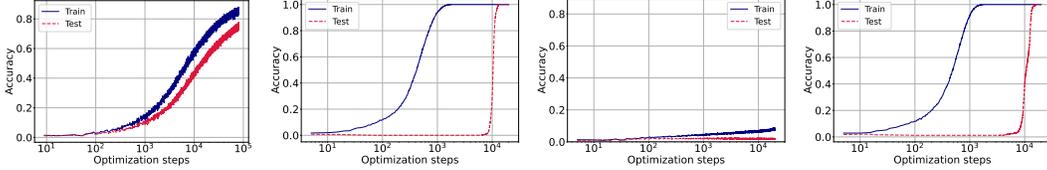


Figure 2: Training and validation accuracies of the MLP model on modular arithmetic tasks, trained with Adam. *Left two:* Addition task, without (first) and with (second) embeddings. *Right two:* Multiplication task, without (third) and with (fourth) embeddings. In the embedding-free cases, training and validation accuracies increase together only for addition; multiplication fails to generalize. In contrast, models with embeddings reach 100% training accuracy in both tasks, but only begin generalizing after a delay exhibiting the grokking phenomenon.

196 4 Main Results

197 Our methodology investigates the dynamics of embedding layers within MLPs to address challenges
 198 in generalization, particularly in the context of algorithmic tasks. The key contributions include:
 199 (1) exploring the novel role of embedding layers attached to MLP architectures, (2) examining the
 200 impact of embedding sampling probability p_i on training dynamics, and (3) understanding how
 201 initialization and the coupling of embedding and weight matrices affect learning efficiency. These
 202 factors contribute to the grokking phenomenon, where generalization is delayed during training.

203 4.1 Embedding Dynamics

204 Let the loss function of the model be $\mathcal{L}(\theta, \mathbf{E})$, where θ is model parameters other than embedding
 205 weights. Let $\mathbf{e}_{i,t}$ denote the embedding vector for token i at step t . Under stochastic gradient descent
 206 (SGD) with weight decay λ , the embedding update rule is:

$$\mathbf{e}_{i,t+1} - \mathbf{e}_{i,t} = -\eta\lambda\mathbf{e}_{i,t} - \eta\nabla_{\mathbf{e}_{i,t}}\mathcal{L}, \quad (3)$$

207 where η is the learning rate, and $\nabla_{\mathbf{e}_i}\mathcal{L}$ is the gradient¹. Token embeddings are updated using
 208 corresponding gradients only when the associated tokens appear in a batch. Assume that token i
 209 being sampled in a batch with a probability p_i . Consequently, taking into account the randomness of
 210 batch sampling, the expected update can be expressed as:

$$\mathbb{E}[\mathbf{e}_{i,t+1} - \mathbf{e}_{i,t}] = -\eta\lambda\mathbf{e}_{i,t} - \eta p_i \nabla_{\mathbf{e}_{i,t}}\mathcal{L}. \quad (4)$$

211 To summarize, the sampling probability p_i directly influences the gradient dynamics of the embedding
 212 layer. While gradients contribute to tokens only probabilistically, weight decay affects all embeddings
 213 uniformly, leading to imbalances in parameter updates. This dynamic, visualized in Figure 3,
 214 highlights the need for a deeper understanding of how p_i affects convergence.

215 To analyze the reduction of the loss, we assume that the model’s overall loss function $\mathcal{L}(\theta, \{\mathbf{e}_i\})$ is
 216 β -smooth. This means it satisfies the following inequality for all updates:

$$\mathcal{L}(\theta_{t+1}, \{\mathbf{e}_{i,t+1}\}) \leq \mathcal{L}(\theta_t, \{\mathbf{e}_{i,t}\}) + \langle \nabla\mathcal{L}, \Delta \rangle + \frac{\beta}{2}\|\Delta\|^2.$$

¹Assuming the SGD update rule without momentum.

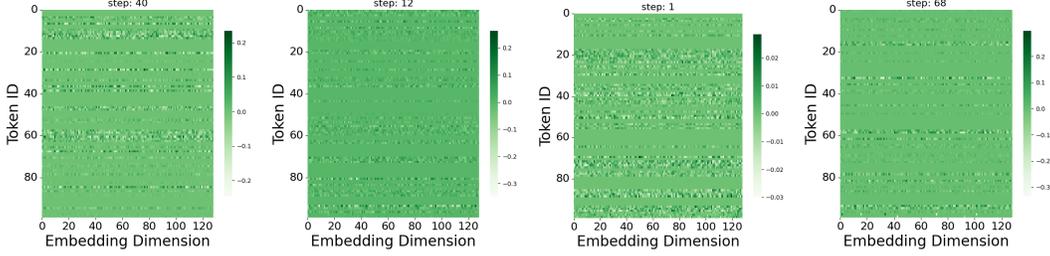


Figure 3: Gradient heat maps of the MLP model at random optimization steps. Sparse columns in the embedding gradients reflect the absence of certain tokens in sampled batches, leading to uneven learning dynamics and contributing to delayed generalization.

217 where $\Delta = (\theta_{t+1} - \theta_t, \mathbf{e}_{i,t+1} - \mathbf{e}_{i,t})$.

218 Denote $\mathcal{L}_t := \mathcal{L}(\theta_t, \{\mathbf{e}_{i,t}\})$ then taking expectations over randomness of batch sampling leads to the
 219 following expected update:

$$\begin{aligned} \mathbb{E}[\mathcal{L}_{t+1} - \mathcal{L}_t] &\leq \nabla_{\theta_t} \mathcal{L}^T(\theta_{t+1} + \theta_t) \\ &\quad - \sum_{i=1}^V \nabla_{\mathbf{e}_{i,t}} \mathcal{L}^T \mathbb{E}(\mathbf{e}_{i,t+1} - \mathbf{e}_{i,t}) + \frac{\beta}{2} \|\Delta\|^2, \end{aligned} \quad (5)$$

220 Substituting the embedding update based on equation 4 into the smoothness inequality,

$$\begin{aligned} \mathbb{E}[\mathcal{L}_{t+1} - \mathcal{L}_t] &\leq \nabla_{\theta_t} \mathcal{L}^T(\theta_{t+1} - \theta_t) \\ &\quad - \eta \sum_{i=1}^V (p_i \|\nabla_{\mathbf{e}_{i,t}} \mathcal{L}\|^2 + \lambda \mathbf{e}_{i,t}^T \nabla_{\mathbf{e}_{i,t}} \mathcal{L}) + \frac{\beta}{2} \|\Delta\|^2, \end{aligned} \quad (6)$$

222 and noting from the right hand side of the inequality above, p_i plays important role in reduction of
 223 the expected loss. However, the dependence on p_i , is coupled with weight decay, which explains why
 224 these two parameters are important to study more deeply to draw a conclusion about grokking.

225 4.2 Dataset Splitting Strategies

226 To further explore the role of p_i , we investigate how train-test splitting strategies affect its value
 227 and, consequently, the grokking process. The train-test split determines the probability of token i
 228 appearing in a batch.

229 We begin by assuming that the weight decay parameter λ is zero and that the learning rate η is uniform
 230 across all parameters. This reduces the optimization problem to focusing on p_i , under the constraints
 231 $\sum_{i=1}^V p_i = 1, p_i \geq 0 \forall i$. Specifically, the optimal p_i can be found by solving for the following:

$$\min_{p_i | p_i \geq 0, \sum p_i = 1} -\eta \sum_{i=1}^V p_i \|\nabla_{\mathbf{e}_{i,t}} \mathcal{L}\|^2. \quad (7)$$

232 However, solving this exactly is challenging in practice due to the need for estimating all embedding
 233 gradient norms. Instead, we adopt approximate strategies for splitting the training data, guided by
 234 various assumptions about the gradient structure (see Appendix A for details).

- 235 1. **Uniform Sampling:** Distribute all combinations of a and b evenly across training and test
 236 sets.
- 237 2. **Skewed Sampling:** Introduce a bias in the combinations of a that are distributed across
 238 training and test sets.
- 239 3. **Random Sampling:** Randomly distribute the examples across training and test sets.

240 These splits enable us to regulate token sampling probabilities, offering a direct assessment of the
 241 impact of p_i on embedding convergence and grokking. Furthermore, Section 5.1 provides a detailed
 242 experiments conducted on two algorithmic datasets.

243 **4.3 Embedding Convergence and Initialization**

244 While the frequency of embedding updates plays
 245 a crucial role in training dynamics, as demon-
 246 strated in our experiments, it alone cannot fully
 247 explain phenomena such as grokking after fit-
 248 ting, its relationship to initialization, weight de-
 249 cay, or the structure of the loss landscape.

250 Stabilization (or convergence) occurs when the
 251 embedding e_i reaches a steady state where the
 252 updates become negligibly small, i.e., when the
 253 change in the embedding $\|e_{i,t+1} - e_{i,t}\|$ is ap-
 254 proximately zero. This condition implies that,
 255 $(\eta\lambda)e_{i,t} \approx \eta p_i \nabla_{e_i} \mathcal{L}$. from equation 4.

256 For small learning rates ($\eta \ll 1$), the embedding
 257 updates behave like a continuous system, and we
 258 can model this as a differential equation (along
 259 every dimension):

$$\frac{de_i}{dt} = -\lambda e_i - p_i \nabla_{e_i} \mathcal{L}, \quad (8)$$

260 where $\nabla_{e_i} \mathcal{L}$ is the gradient of the loss function
 261 with respect to the embedding i . Assuming that
 262 the gradient $\nabla_{e_i} \mathcal{L}$ stabilizes to a constant value
 263 g , the solution to this equation is:

$$e_i(t) = Ce^{-\lambda t} - \frac{\eta p g}{\lambda}, \quad (9)$$

264 where C is an integration constant determined
 265 by the initial conditions. As time t increases,
 266 the embedding $e_i(t)$ converges to the equilib-
 267 rium value $e_i(t) \rightarrow -\frac{\eta p g}{\lambda}$. Thus, conver-
 268 gence is achieved when $e_i(t)$ stabilizes around this
 269 equilibrium point. The time T to reach conver-
 270 gence is bounded as $T \geq \frac{1}{\lambda} \ln\left(\frac{C}{\epsilon}\right)$, where ϵ is a small
 271 threshold. In summary, convergence time is governed by the embedding gradient g , the weight decay
 272 λ , and the initialization magnitude C : stronger gradients and larger λ accelerate convergence, while
 273 larger initial values C slow it down.

274 In bilinear models such as MLPs and Transformers, embedding gradients are tightly coupled with
 275 those of downstream weights (e.g., \mathbf{W}), forming a feedback loop: poor updates to \mathbf{E} degrade \mathbf{W} ,
 276 and vice versa. To study the role of initialization in this dynamic, we tested two setups: frozen
 277 embeddings, which led to slow convergence due to limited representational flexibility; and small
 278 initial embeddings, which improved convergence by allowing stronger early gradients—an effect also
 observed in prior work [26, 12], though without analyzing embedding-weight coupling.

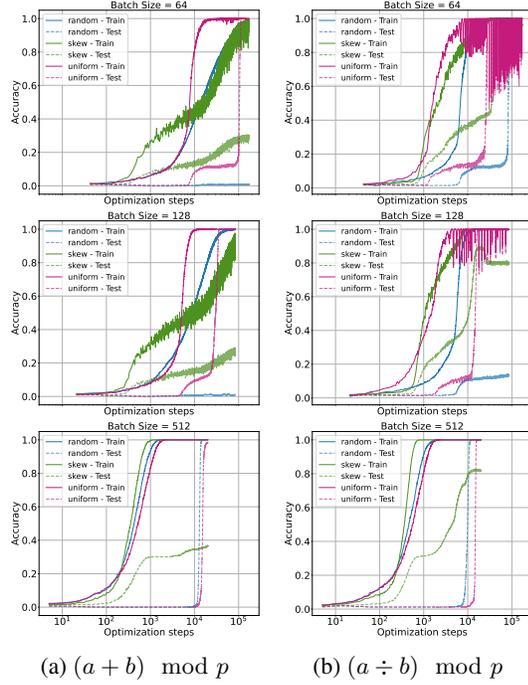
279 Motivated by these observations, we propose the **Adam-LR Optimizer**, which adjusts the embedding
 280 learning rate to balance update magnitudes between \mathbf{E} and \mathbf{W} . This coupling-aware scaling is
 281 formalized below:

282 **Proposition 4.1.** *Let \mathbf{E} and \mathbf{W} be the embedding matrix and first-layer weights. To equalize update*
 283 *scales under cross-entropy loss, the learning rate ratio $c = \frac{\eta_E}{\eta_W}$ should satisfy:*

$$c \propto \frac{\sigma_{\max}(\mathbf{E})}{\sigma_{\max}(\mathbf{W})} \cdot \frac{f_W}{f_E},$$

284 where $\sigma_{\max}(\cdot)$ denotes the largest singular value and f_E, f_W are the respective update frequen-
 285 cies, (see appendix B for details).

286 In practice, we set $c = 10$, guided by empirical singular value trends and supported by sensitivity
 287 analysis (see Fig. 7, §5.2). This adjustment improves convergence and stability, especially under
 288 sparse embedding updates common in skewed token distributions.



(a) $(a + b) \bmod p$ (b) $(a \div b) \bmod p$
 Figure 4: Sampling strategy comparison for two modular tasks—addition and division—across all batch sizes. Uniform sampling generalizes faster; skewed sampling fails to generalize due to token imbalance.

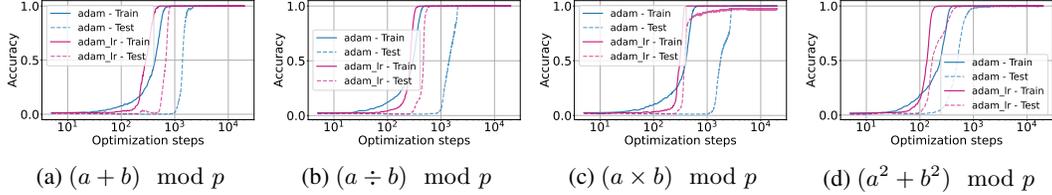


Figure 5: Performance comparison of Adam-LR and Adam optimizers on four algorithmic datasets. Adam-LR scales the embedding learning rate based on the singular values of the embedding matrix. This adaptive adjustment accelerates convergence and enhances generalization across all datasets. The results demonstrate that Adam-LR significantly speeds up the grokking process compared to the standard Adam optimizer under identical training settings ($lr = 0.01$, batch size = 512).

289 5 Experiments and Discussions

290 We begin our exploration with a MLP model. The architecture consists of two layers, where the
 291 hidden dimension of the first layer is set to four times the embedding dimension (where four is the
 292 sequence length), and embedding dimension is set to 128, as per prior work on grokking. The second
 293 layer has a dimension of $P = 97$. The activation function used throughout is ReLU, and optimization
 294 is performed using the Adam optimizer with a weight decay of 0.001.

295 5.1 The Effect of Embedding Probability

296 The first set of experiments investigates various strategies for splitting the training and testing datasets.
 297 Specifically, we explore three approaches, namely; uniform sampling, skewed sampling, and random
 298 sampling.

299 The expression $(a + b) \bmod p$ represents the sum of a and b modulo p . For our experiments, we
 300 randomly set aside 20% of the data as a test set, ensuring that evaluation is performed on unseen
 301 samples. From the remaining data, 30/80% (i.e. 30% from total set) is sampled as the training set
 302 according to each sampling strategy.

303 Figure 4 compare the performance of the sampling methods (random, uniform, skew) across different
 304 splits of the dataset (see appendix D.1 for further datasets and settings). Each represents a specific
 305 datasets, while the rows compare batch sizes, and columns compare datasets. The x-axis is logarithmic
 306 to emphasize the convergence trends.

307 Uniform sampling generally promotes faster generalization and convergence compared to random
 308 sampling. However, its benefits diminish at larger batch sizes (e.g., beyond 512), where random
 309 sampling becomes nearly as effective due to broader token coverage. Crucially, our results show
 310 that skewed sampling—despite fitting the training data and preserving the overall train-test ratio—
 311 consistently leads to suboptimal generalization. This suggests that models can converge to lower
 312 subaccuracy plateaus when token probabilities are heavily imbalanced. Importantly, even uniform
 313 sampling does not guarantee optimality: unless the batch size is sufficiently large, some tokens may
 314 be consistently omitted from updates. These findings underscore that token probability, both in
 315 expectation and in per-batch coverage, plays a central role in embedding dynamics and grokking
 316 behavior.

317 5.2 Comparison of Optimizers

318 To evaluate the effectiveness of our proposed optimizer, Adam-LR, which incorporates a simple yet
 319 effective strategy for treating the embedding layer differently to avoid stagnation or saddle points,
 320 we conducted experiments on four datasets. The results are shown in Figure 5, where we compare
 321 the performance of the two optimizers, Adam-LR and the standard Adam optimizer, under identical
 322 training settings ($lr = 0.01$, batch size = 512).

323 Using our proposed optimizer, Adam-LR, which scales the embedding learning rate by a factor of 10,
 324 the results demonstrate a significant acceleration in the grokking process compared to the baseline
 325 Adam optimizer across all datasets.

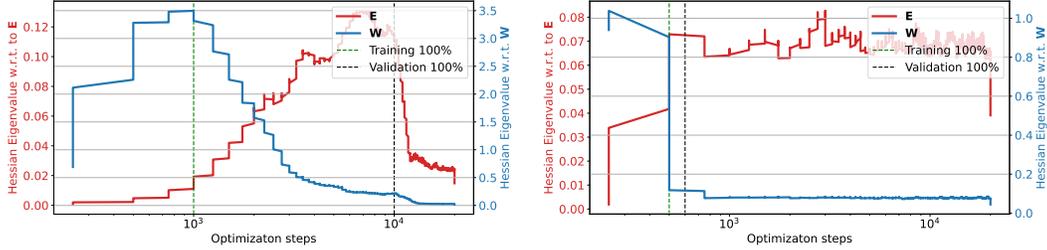


Figure 6: Maximum eigenvalues of the Hessian with respect to embedding weights (**E**) and downstream weights (**W**) during training. The left plot corresponds to the Adam optimizer, while the right plot uses Adam_lr optimizer (ours). With Adam (left), the eigenvalues for **E** are significantly smaller than those for **W**, reflecting differences in dimensionality and update frequency. In contrast, with Adam_lr (right), the eigenvalues of **W** are notably reduced and become closer to those of **E**, suggesting a more balanced optimization dynamic. Training accuracy reaches 100% when the eigenvalues of **W** begin to decrease, while validation accuracy improves as the eigenvalues of **E** decrease. This suggests that **W** drives early optimization progress, while **E** fine-tunes generalization. The Adam_lr optimizer (ours) appears to regularize **W**, leading to a more stable training process.

326 5.3 Analysis of singular values of embedding layer

327 Prior work attributes Adam’s superiority over SGD in Transformers to factors like gradient noise,
 328 descent direction, and Hessian block heterogeneity [25, 10, 18, 26]. However, these studies largely
 329 overlook the role of embeddings and their bilinear interactions. Our analysis supports the view that
 330 such bilinear structure, especially in embeddings, contributes significantly to the observed curvature
 331 differences (see appendix C.1 for more discussion).

332 To analyze the curvature of the loss landscape, we compute the maximum eigenvalue
 333 of the Hessian matrix using the power method with Hessian-vector products (HVPs).

334 Figure 6 shows the maximum eigenvalues of the Hessian with respect to **E** and **W** during
 335 training. The results highlight distinct curvature
 336 properties for **E** and **W**, reflecting their roles in
 337 the bilinear interaction.
 338

339 6 Discussions

340 In this study, we explored the interplay between
 341 embedding layers and downstream weights in
 342 neural networks, highlighting how their bilinear
 343 coupling influences optimization and drives
 344 the grokking phenomenon. We demonstrated
 345 that embedding layers play a central role in
 346 delayed generalization and introduced the Adam-
 347 LR optimizer to address the imbalance in update
 348 dynamics, scaling the embedding learning rate
 349 based on singular values and update frequencies.

350 A key limitation of this work is its focus on
 351 MLPs, which provide a simplified setting for
 352 analyzing embedding-weight coupling. While
 353 this enables controlled analysis, it leaves open
 354 how these insights transfer to more complex
 355 architectures such as Transformers, where similar
 356 bilinear interactions appear in attention mechanisms but with added structural complexity. Extending
 357 our framework to the Transformer setting is a promising direction for future work.

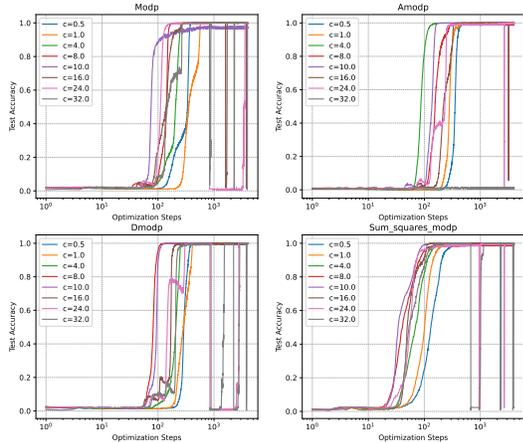


Figure 7: Sensitivity of test accuracy to the learning rate ratio $c = \eta_E/\eta_W$ across four tasks. Small c leads to under-updating, large c causes instability, and $c = 10$ consistently balances convergence and stability.

References

- 358
- 359 [1] X. Davies, L. Langosco, and D. Krueger. Unifying grokking and double descent. *arXiv preprint*
360 *arXiv:2303.06173*, 2023.
- 361 [2] D. Doshi, T. He, A. Das, and A. Gromov. Grokking modular polynomials. *arXiv preprint*
362 *arXiv:2406.03495*, 2024.
- 363 [3] S. Fan, R. Pascanu, and M. Jaggi. Deep grokking: Would deep neural networks generalize
364 better? *arXiv preprint arXiv:2405.19454*, 2024.
- 365 [4] A. Gromov. Grokking modular arithmetic. *arXiv preprint arXiv:2301.02679*, 2023.
- 366 [5] X. S. Huang, F. Perez, J. Ba, and M. Volkovs. Improving transformer optimization through
367 better initialization. In *International Conference on Machine Learning*, pages 4475–4483.
368 PMLR, 2020.
- 369 [6] A. I. Humayun, R. Balestriero, and R. Baraniuk. Deep networks always grok and here is why.
370 *arXiv preprint arXiv:2402.15555*, 2024.
- 371 [7] A. Jeffares, A. Curth, and M. van der Schaar. Deep learning through a telescoping lens: A
372 simple model provides empirical insights on grokking, gradient boosting & beyond. *Advances*
373 *in Neural Information Processing Systems*, 37:123498–123533, 2024.
- 374 [8] S. Kobayashi, Y. Akram, and J. Von Oswald. Weight decay induces low-rank attention layers.
375 *Advances in Neural Information Processing Systems*, 37:4481–4510, 2024.
- 376 [9] T. Kumar. *Grokking as the transition from lazy to rich training dynamics*. PhD thesis, none,
377 2024.
- 378 [10] F. Kunstner, J. Chen, J. W. Lavington, and M. Schmidt. Noise is not the main factor behind
379 the gap between sgd and adam on transformers, but sign descent might be. *arXiv preprint*
380 *arXiv:2304.13960*, 2023.
- 381 [11] J. Lee, B. G. Kang, K. Kim, and K. M. Lee. Grokfast: Accelerated grokking by amplifying
382 slow gradients. *arXiv preprint arXiv:2405.20233*, 2024.
- 383 [12] Z. Liu, O. Kitouni, N. S. Nolte, E. Michaud, M. Tegmark, and M. Williams. Towards un-
384 derstanding grokking: An effective theory of representation learning. *Advances in Neural*
385 *Information Processing Systems*, 35:34651–34663, 2022.
- 386 [13] Z. Liu, E. J. Michaud, and M. Tegmark. Omnigrok: Grokking beyond algorithmic data. In *The*
387 *Eleventh International Conference on Learning Representations*, 2022.
- 388 [14] K. Lyu, J. Jin, Z. Li, S. S. Du, J. D. Lee, and W. Hu. Dichotomy of early and late phase implicit
389 biases can provably induce grokking. *arXiv preprint arXiv:2311.18817*, 2023.
- 390 [15] M. A. Mohamadi, Z. Li, L. Wu, and D. J. Sutherland. Why do you grok? a theoretical analysis
391 of grokking modular addition. *arXiv preprint arXiv:2407.12332*, 2024.
- 392 [16] P. Nakkiran, G. Kaplun, Y. Bansal, T. Yang, B. Barak, and I. Sutskever. Deep double descent:
393 Where bigger models and more data hurt. *Journal of Statistical Mechanics: Theory and*
394 *Experiment*, 2021(12):124003, 2021.
- 395 [17] N. Nanda, L. Chan, T. Lieberum, J. Smith, and J. Steinhardt. Progress measures for grokking
396 via mechanistic interpretability. *arXiv preprint arXiv:2301.05217*, 2023.
- 397 [18] Y. Pan and Y. Li. Toward understanding why adam converges faster than sgd for transformers.
398 *arXiv preprint arXiv:2306.00204*, 2023.
- 399 [19] A. Power, Y. Burda, H. Edwards, I. Babuschkin, and V. Misra. Grokking: Generalization beyond
400 overfitting on small algorithmic datasets. *arXiv preprint arXiv:2201.02177*, 2022.
- 401 [20] L. Prieto, M. Barsbey, P. A. Mediano, and T. Birdal. Grokking at the edge of numerical stability.
402 *arXiv preprint arXiv:2501.04697*, 2025.

- 403 [21] V. Varma, R. Shah, Z. Kenton, J. Kramár, and R. Kumar. Explaining grokking through circuit
404 efficiency. *arXiv preprint arXiv:2309.02390*, 2023.
- 405 [22] A. Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*,
406 2017.
- 407 [23] Z. Xu, Z. Ni, Y. Wang, and W. Hu. Let me grok for you: Accelerating grokking via embedding
408 transfer from a weaker model. *arXiv preprint arXiv:2504.13292*, 2025.
- 409 [24] G. Yang and E. J. Hu. Tensor programs iv: Feature learning in infinite-width neural networks. In
410 M. Meila and T. Zhang, editors, *Proceedings of the 38th International Conference on Machine*
411 *Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 11727–11737.
412 PMLR, 18–24 Jul 2021.
- 413 [25] J. Zhang, S. P. Karimireddy, A. Veit, S. Kim, S. Reddi, S. Kumar, and S. Sra. Why are adaptive
414 methods good for attention models? *Advances in Neural Information Processing Systems*,
415 33:15383–15393, 2020.
- 416 [26] Y. Zhang, C. Chen, T. Ding, Z. Li, R. Sun, and Z.-Q. Luo. Why transformers need adam: A
417 hessian perspective. *arXiv preprint arXiv:2402.16788*, 2024.

418 **Appendix**

419 **A Optimizing for Sampling Probability**

420 **Uniform Importance Assumption**

421 If we assume that all gradients are equally important, i.e., $\|\nabla_{\mathbf{E}_{i,t}} \mathcal{L}\|^2$ is uniform across all embed-
422 dings:

$$\|\nabla_{\mathbf{E}_{i,t}} \mathcal{L}\|^2 = c, \quad \forall i,$$

423 where c is a constant.

424 In this case, the optimization of $-\sum_{i=1}^V p_i \|\nabla_{\mathbf{E}_{i,t}} \mathcal{L}\|^2$ becomes independent of p_i . To satisfy the
425 normalization constraint $\sum_{i=1}^V p_i = 1$, the optimal solution is:

$$p_i = \frac{1}{V}, \quad \forall i. \quad (10)$$

426 This corresponds to a uniform distribution, where all embeddings are treated equally (see Figure
427 8). While computationally efficient, this approach may lead to suboptimal convergence if some
428 embeddings contribute disproportionately to the loss reduction.

429 **Gradient Norm Bounded by L_i**

430 Now, let us assume that the gradient norm for each embedding is bounded,

$$\|\nabla_{\mathbf{E}_{i,t}} \mathcal{L}\| \leq L_i, \quad \forall i, \quad (11)$$

431 where L_i is a known upper bound for embedding i . Using this bound, we approximate,

$$-\sum_{i=1}^V p_i \|\nabla_{\mathbf{E}_{i,t}} \mathcal{L}\|^2 \geq -\sum_{i=1}^V p_i L_i^2. \quad (12)$$

432 To maximize $\sum_{i=1}^V p_i L_i^2$ subject to the constraint $\sum_{i=1}^V p_i = 1$, we note that the objective function
433 is linear in \mathbf{p} . Therefore, the maximum is attained at a vertex of the probability simplex, meaning the
434 optimal solution is:

$$p_k = 1, \quad \text{where } k = \arg \max_i L_i^2, \quad \text{and } p_i = 0, \quad \forall i \neq k. \quad (13)$$

435 This result indicates that the optimal probability distribution assigns all weight to the embedding with
436 the highest gradient bound, ignoring all others. Therefore, to obtain a smooth probability distribution,
437 we introduce an entropy regularization term as follow,

$$H(\mathbf{p}) = -\sum_{i=1}^V p_i \log p_i. \quad (14)$$

438 We now optimize the modified objective,

$$\sum_{i=1}^V p_i L_i^2 + \gamma H(\mathbf{p}), \quad (15)$$

439 subject to the constraint $\sum_{i=1}^V p_i = 1$, where $\gamma > 0$ controls the strength of the regularization.

440 The corresponding Lagrangian is as follow,

$$\mathcal{L}_p = \sum_{i=1}^V p_i L_i^2 + \gamma \left(-\sum_{i=1}^V p_i \log p_i \right) + \mu \left(\sum_{i=1}^V p_i - 1 \right). \quad (16)$$

441 Taking the derivative with respect to p_i and setting it to zero, we get,

$$L_i^2 - \gamma(1 + \log p_i) + \mu = 0. \quad (17)$$

442 Solving for p_i gives:

$$\log p_i = \frac{L_i^2 + \mu - \gamma}{\gamma} \implies p_i = \exp\left(\frac{L_i^2 + \mu - \gamma}{\gamma}\right). \quad (18)$$

443 Applying the constraint $\sum_{i=1}^V p_i = 1$, would results in the following solution,

$$p_i^* = \frac{\exp(L_i^2/\gamma)}{\sum_{j=1}^V \exp(L_j^2/\gamma)}. \quad (19)$$

444 This result smoothly distributes probabilities based on the gradient bounds, assigning higher probabili-
445 ty to embeddings with larger L_i^2 while ensuring a non-degenerate distribution.

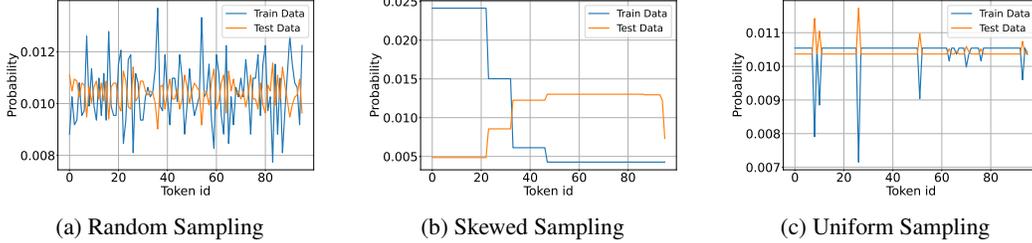


Figure 8: Token probabilities in the training and test sets under different sampling strategies. Imbalanced sampling leads to uneven token occurrences in mini-batches, causing some tokens to be absent in multiple updates while others appear frequently. This results in highly variable gradient updates, where frequently seen tokens converge faster, while rare tokens stagnate due to sparse updates, affecting overall model generalization.

446 B Dynamics of Updates in Bilinear Systems with Initialization Effects

447 We analyze the interaction between embeddings $\mathbf{E} \in \mathbb{R}^{p \times d}$ and weight matrix $\mathbf{W} \in \mathbb{R}^{4d \times d}$ in a
448 bilinear term:

$$z(\mathbf{E}\mathbf{W}), \quad (20)$$

449 where z is an activation function applied elementwise. The gradients of \mathbf{E} and \mathbf{W} are given as:

$$\nabla_{\mathbf{E}} \propto \mathbf{W}^\top \nabla_{\text{loss}}, \quad \nabla_{\mathbf{W}} \propto \mathbf{E}^\top \nabla_{\text{loss}}. \quad (21)$$

450 The gradient norms are influenced by the dominant singular values of \mathbf{W} and \mathbf{E} . Specifically:

$$\|\nabla_{\mathbf{E}}\| \propto \sigma_{\max}(\mathbf{W}), \quad \|\nabla_{\mathbf{W}}\| \propto \sigma_{\max}(\mathbf{E}). \quad (22)$$

451 At initialization, \mathbf{E} and \mathbf{W} are often drawn from distributions with variances that depend on their
452 dimensions (e.g., PyTorch initializes weights with $\mathcal{N}(0, \sqrt{2/d})$ scaling). This initialization typically
453 ensures $\sigma_{\max}(\mathbf{E}) \gg \sigma_{\max}(\mathbf{W})$, as \mathbf{W} is higher-dimensional, amplifying the difference in gradient
454 magnitudes.

455 The embedding matrix \mathbf{E} is updated less frequently than \mathbf{W} because not all tokens appear in every
456 batch. Let f_E and f_W represent the update frequencies of \mathbf{E} and \mathbf{W} , respectively. Typically,
457 $f_W > f_E$, exacerbating the update disparity.

458 To balance the effective updates of \mathbf{E} and \mathbf{W} , the learning rates η_E and η_W must be scaled to account
459 for both their singular values and update frequencies. The effective update ratio is:

$$\frac{\|\Delta \mathbf{E}\|}{\|\Delta \mathbf{W}\|} \propto \frac{\eta_E \cdot \sigma_{\max}(\mathbf{W}) \cdot f_E}{\eta_W \cdot \sigma_{\max}(\mathbf{E}) \cdot f_W}. \quad (23)$$

460 For proportional updates ($\|\Delta \mathbf{E}\| \sim \|\Delta \mathbf{W}\|$), the ratio $c = \frac{\eta_E}{\eta_W}$ must satisfy:

$$c \propto \frac{\sigma_{\max}(\mathbf{E})}{\sigma_{\max}(\mathbf{W})} \cdot \frac{f_W}{f_E}. \quad (24)$$

461 The term $\frac{\sigma_{\max}(\mathbf{E})}{\sigma_{\max}(\mathbf{W})}$ reflects the imbalance in singular values due to initialization and structural
 462 properties. The term $\frac{f_W}{f_E}$ accounts for the frequency imbalance in updates between \mathbf{E} and \mathbf{W} , driven
 463 by sparse token appearances in batches.

464 PyTorch initialization, which scales weights by $\mathcal{O}(\sqrt{2/d})$, ensures that $\sigma_{\max}(\mathbf{W})$ and $\sigma_{\max}(\mathbf{E})$ are
 465 initially proportional to the dimensions d . This contributes to the observed imbalance in their singular
 466 values at the start of training.

467 C More experiments

468 C.1 Analysis of singular values of embedding layer

469 Previous studies (e.g., [25], [10], [18], [26]) have explored the gap between SGD and Adam in
 470 optimizing Transformer models, but the specific role of embeddings and their bilinearity with down-
 471 stream weights remains underexplored. For example, [25] attributes SGD’s suboptimal performance
 472 to the heavy-tailed distribution of stochastic gradient noise. This observation aligns with our findings
 473 regarding the randomness in embedding updates for low- p tokens.

474 On the other hand, [10] argues that gradient noise alone cannot explain Adam’s superiority. Their
 475 experiments demonstrate that, even with full-batch training to eliminate stochastic noise, SGD
 476 underperforms compared to Adam. They suggest that the sign of the gradient might be a more reliable
 477 descent direction than its magnitude, and since Adam optimally balances both, it outperforms SGD,
 478 particularly in small-batch settings.

479 Furthermore, [26] provides a novel explanation for Adam’s advantage over SGD in Transformers
 480 by analyzing the blockwise Hessian spectrum, introducing the concept of “block heterogeneity.”
 481 This refers to significant variations in the Hessian spectra across parameter blocks, a phenomenon
 482 observed in Transformers but not in CNNs. However, the underlying source of this heterogeneity
 483 is not explicitly discussed. We hypothesize that this stems from the bilinear nature of weights,
 484 particularly in the embedding and attention mechanisms. To support this hypothesis, we analyze the
 485 Hessian of embedding weights compared to other weight below.

486 To analyze the curvature of the loss landscape, we compute the maximum eigenvalue of the Hessian
 487 matrix using the power method with Hessian-vector products (HVPs). This approach avoids explicitly
 488 constructing the Hessian, making it computationally efficient for large-scale systems.

489 The power method iteratively approximates the maximum eigenvalue of the Hessian \mathbf{H} as follows:

- 490 1. Initialize a random vector \mathbf{v}_0 with the same dimensionality as the parameters $[\mathbf{E}, \mathbf{W}]$.
- 491 2. Compute the Hessian-vector product $\mathbf{H}\mathbf{v}_k$ using automatic differentiation:

$$\mathbf{H}\mathbf{v}_k = \nabla_{\theta} (\nabla_{\theta} \mathcal{L} \cdot \mathbf{v}_k),$$

492 where $\theta = [\mathbf{E}, \mathbf{W}]$.

- 493 3. Normalize the vector and update the eigenvalue estimate:

$$\mathbf{v}_{k+1} = \frac{\mathbf{H}\mathbf{v}_k}{\|\mathbf{H}\mathbf{v}_k\|}, \quad \sigma_{\max} \approx \mathbf{v}_k^{\top} \mathbf{H}\mathbf{v}_k.$$

494 Figure 9 shows the maximum eigenvalues of the Hessian with respect to \mathbf{E} and \mathbf{W} during training.
 495 The results highlight distinct curvature properties for \mathbf{E} and \mathbf{W} , reflecting their roles in the bilinear
 496 interaction.

497 Extending these insights to attention mechanisms highlights further challenges in bilinear optimization
 498 and demonstrates how adaptive learning rates (e.g., Adam) help escape saddle points. This suggests
 499 a deeper connection between the bilinearity of weight interactions and the optimization challenges
 500 unique to Transformer models.

501 C.2 Rank Evolution and Implicit Regularization

502 Recent work has shown that weight decay in bilinear models (e.g., $\mathbf{Z} = \mathbf{E}\mathbf{W}$) implicitly regularizes
 503 the nuclear norm of the product matrix, promoting low-rank solutions and improved generalization

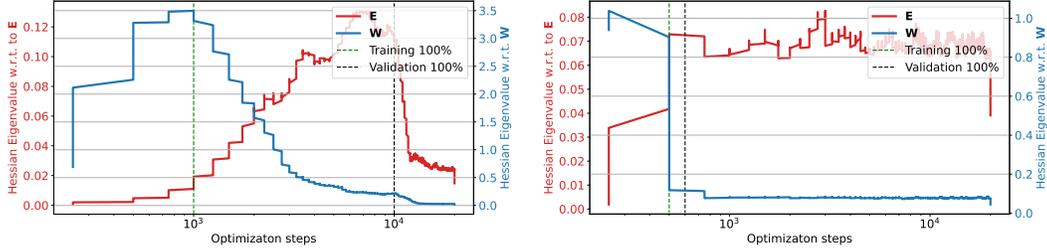


Figure 9: Maximum eigenvalues of the Hessian with respect to embedding weights (\mathbf{E}) and downstream weights (\mathbf{W}) during training. The left plot corresponds to the Adam optimizer, while the right plot uses Adam_lr optimizer (ours). With Adam (left), the eigenvalues for \mathbf{E} are significantly smaller than those for \mathbf{W} , reflecting differences in dimensionality and update frequency. In contrast, with Adam_lr (right), the eigenvalues of \mathbf{W} are notably reduced and become closer to those of \mathbf{E} , suggesting a more balanced optimization dynamic. Training accuracy reaches 100% when the eigenvalues of \mathbf{W} begin to decrease, while validation accuracy improves as the eigenvalues of \mathbf{E} decrease. This suggests that \mathbf{W} drives early optimization progress, while \mathbf{E} fine-tunes generalization. The Adam_lr optimizer (ours) appears to regularize \mathbf{W} , leading to a more stable training process.

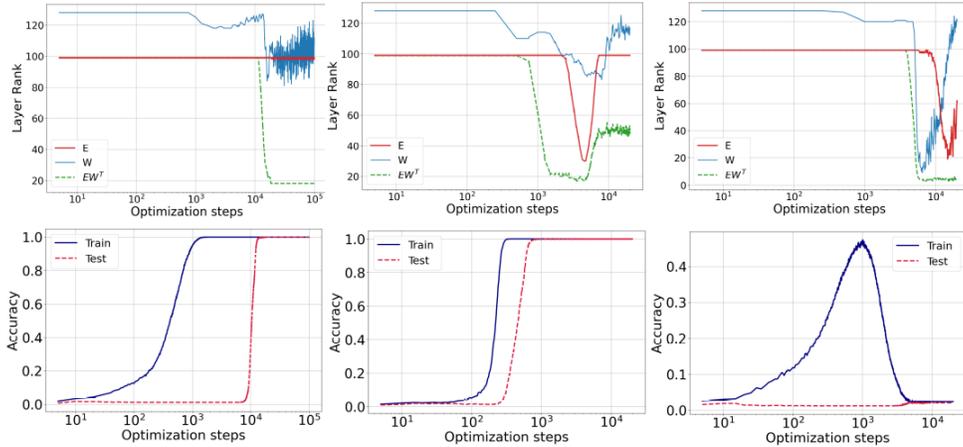


Figure 10: Rank evolution during training for three optimization setups: Adam ($wd=0.001$), Adam-LR ($wd=0.001$ with learning rate ratio), and Adam with stronger weight decay ($wd=0.005$). While all runs show decreasing $\text{rank}(\mathbf{E}\mathbf{W})$, only Adam-LR continues to adjust rank after generalization. This suggests that rank behavior alone does not fully explain grokking, and supports the need to analyze embedding-weight coupling dynamics.

504 [8]. This complements our focus on embedding dynamics, as both highlight the impact of bilinear
 505 coupling on optimization.

506 To explore this in our setup, we track the rank evolution of \mathbf{E} , \mathbf{W} , and the product $\mathbf{E}\mathbf{W}$. As shown in
 507 Figure 10, \mathbf{W} exhibits three distinct phases: an early drop during training loss reduction, a plateau,
 508 and a final decline aligned with grokking. In contrast, \mathbf{E} 's rank remains largely stable throughout.

509 Figure 10 compares three optimization setups: Adam (with weight decay 0.001), Adam-LR (our
 510 proposed variant with a learning rate ratio), and Adam with stronger weight decay (0.005). All
 511 configurations lead to a reduction in $\text{rank}(\mathbf{E}\mathbf{W})$, consistent with implicit nuclear norm regularization.
 512 However, only Adam-LR shows continued rank changes after generalization, suggesting that rank
 513 evolution alone does not capture the onset of grokking.

514 These findings reinforce that implicit regularization in bilinear systems depends not just on decay
 515 strength, but also on the interplay between initialization, update frequency, and curvature.

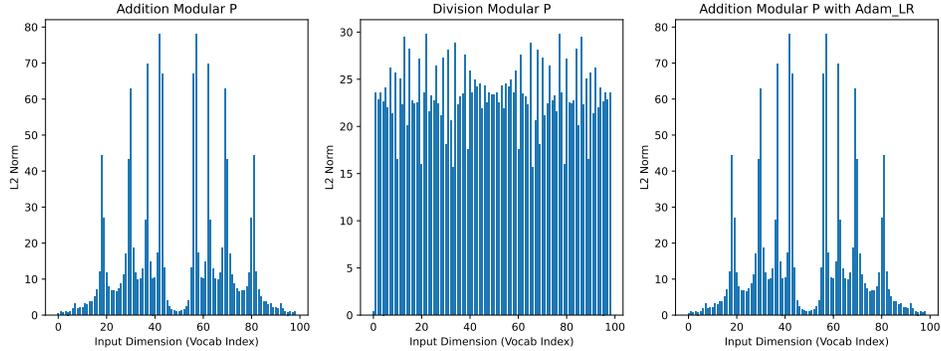


Figure 11: Discrete Fourier analysis of learned embedding representations across tasks. For each embedding matrix, we compute the DFT across the input dimension and the ℓ_2 -norm across the embedding dimension. Peaks indicate frequency localization that naturally aligns with the periodic structure of the task (e.g., modular addition), while tasks like modular division show more diffuse spectra.

516 D Fourier Analysis of Embedding Representations

517 Fourier features offer a structured way to encode modular arithmetic directly into the input space. By
 518 encoding periodicity into the representation, such features can bypass the need for learned embeddings
 519 and mitigate challenges like sparse updates for rare tokens. However, this approach requires prior
 520 knowledge of the task’s structure—e.g., periodicity—which may not apply in more complex tasks
 521 such as modular division or nonlinear compositions.

522 To investigate whether embedding layers naturally learn such structure, we analyze their frequency
 523 characteristics. Following the approach in [12], we apply the Discrete Fourier Transform (DFT)
 524 along the input dimension of the embedding matrix and compute the ℓ_2 -norm across the embedding
 525 dimension. We then plot the first $P/2$ components, leveraging the symmetry of the DFT.

526 The results for different tasks are shown in Figure 11. Clear frequency peaks indicate that the model
 527 internally captures task-specific periodic structure. Notably, such structure emerges even without
 528 explicit Fourier features, especially for modular addition and multiplication. However, in more
 529 complex tasks, such as modular division, this frequency localization diminishes—suggesting the
 530 limits of periodic encoding and the growing need for learned representations.

531 D.1 Additional Datasets and Learning Rate Sensitivity

532 In addition to modular addition and division, we evaluate our methods on two further tasks: modular
 533 multiplication $(a \div b) \bmod p$ and sum of squares $(a^2 + b^2) \bmod p$. These tasks share the same
 534 architecture and tokenization as described in Section 5.

535 We emphasize that our experimental design is not centered on hyperparameter optimization. While
 536 aggressive tuning of learning rates and batch sizes can suppress or delay grokking, our goal is to
 537 study it where it naturally occurs. To that end, we identify configurations where grokking persists
 538 and focus our analysis there. This approach aligns with prior work on mechanistic understanding
 539 of grokking [9, 14], which likewise prioritize clarity of dynamics over benchmark performance.
 540 For illustration, Figures 13 and 14 show learning rate sensitivity on four datasets, confirming the
 541 robustness of our findings across reasonable settings (skewed distribution of embedding update delay
 542 the generalization).

543 Compute Resources

544 All experiments were conducted using an NVIDIA A6000 GPU. Training runs were performed
 545 using PyTorch, with each configuration fitting comfortably within the GPU’s 48 GB memory. No
 546 distributed training or multi-GPU setups were used.

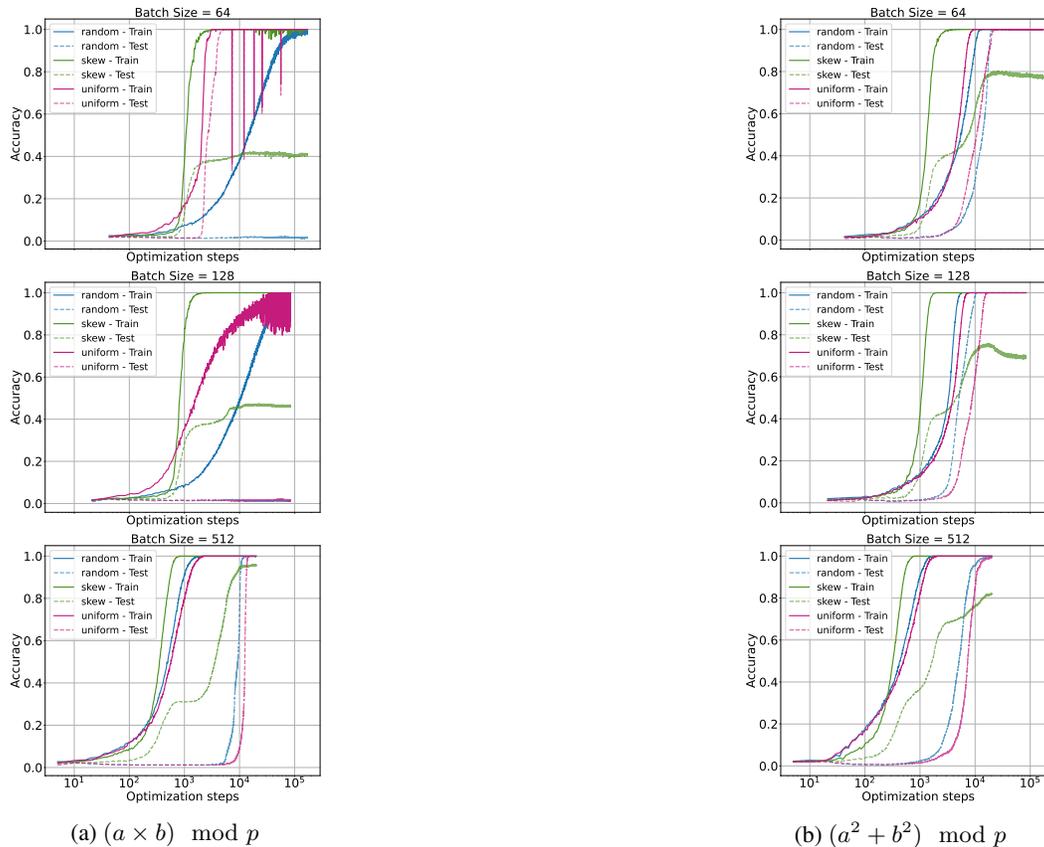


Figure 12: Sampling strategy comparison for multiplication and sum-of-squares tasks ($lr = 0.001$). Larger batch sizes narrow the performance gap, but skewed sampling still harms generalization.

547 NeurIPS Paper Checklist

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

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

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

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

563 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.
 564 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a
 565 proper justification is given (e.g., "error bars are not reported because it would be too computationally
 566 expensive" or "we were unable to find the license for the dataset we used"). In general, answering
 567 "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we

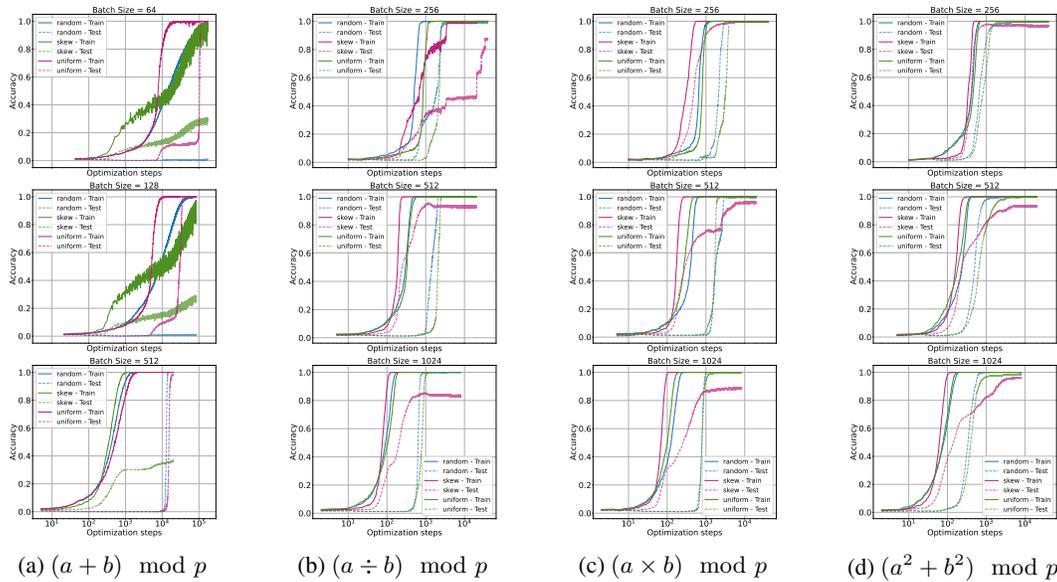


Figure 13: Training and validation accuracies for the modular multiplication dataset for learning rate 0.01 across batch sizes (256, 512, 1024).

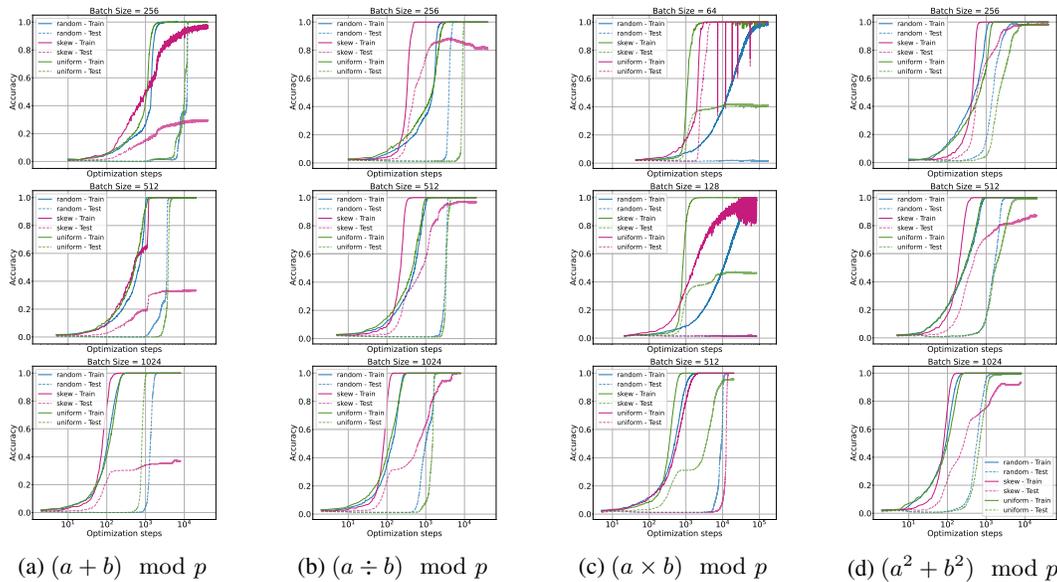


Figure 14: Training and validation accuracies for the modular multiplication dataset for learning rate 0.005 across batch sizes (256, 512, 1024).

568 acknowledge that the true answer is often more nuanced, so please just use your best judgment and
 569 write a justification to elaborate. All supporting evidence can appear either in the main paper or the
 570 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification
 571 please point to the section(s) where related material for the question can be found.

572 IMPORTANT, please:

- 573 • Delete this instruction block, but keep the section heading “NeurIPS Paper Checklist”,
- 574 • Keep the checklist subsection headings, questions/answers and guidelines below.
- 575 • Do not modify the questions and only use the provided macros for your answers.

576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628

1. Claims

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

Answer: [Yes]

Justification: Yes, the main claim made in the abstract and introduction are reflected in the paper Sections 3,4 and 5.

Guidelines:

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

2. Limitations

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

Answer: [Yes]

Justification: The limitations are discussed in section 5.3.

Guidelines:

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

3. Theory assumptions and proofs

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

Answer: [Yes]

Justification: The two main theories are supported by assumptions and proofs in the main body and the appendix.

629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682

Guidelines:

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

4. Experimental result reproducibility

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

Answer: [\[Yes\]](#)

Justification: The details of the experiments are detailed in section 5.

Guidelines:

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

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [\[Yes\]](#)

683 Justification: The data used in the experiments are available online, and the code of experi-
684 ments will be submitted with the supplementary files.

685 Guidelines:

- 686 • The answer NA means that paper does not include experiments requiring code.
- 687 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
688 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 689 • While we encourage the release of code and data, we understand that this might not be
690 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
691 including code, unless this is central to the contribution (e.g., for a new open-source
692 benchmark).
- 693 • The instructions should contain the exact command and environment needed to run to
694 reproduce the results. See the NeurIPS code and data submission guidelines ([https://nips.
695 cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 696 • The authors should provide instructions on data access and preparation, including how
697 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 698 • The authors should provide scripts to reproduce all experimental results for the new
699 proposed method and baselines. If only a subset of experiments are reproducible, they
700 should state which ones are omitted from the script and why.
- 701 • At submission time, to preserve anonymity, the authors should release anonymized
702 versions (if applicable).
- 703 • Providing as much information as possible in supplemental material (appended to the
704 paper) is recommended, but including URLs to data and code is permitted.

705 6. Experimental setting/details

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

709 Answer: [Yes]

710 Justification: The author provides all the details to perform the experiments, and performed
711 sensitivity analysis whenever applicable.

712 Guidelines:

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

718 7. Experiment statistical significance

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

721 Answer: [Yes]

722 Justification: The experiments performed under controlled setting with random seed.

723 Guidelines:

- 724 • The answer NA means that the paper does not include experiments.
- 725 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
726 dence intervals, or statistical significance tests, at least for the experiments that support
727 the main claims of the paper.
- 728 • The factors of variability that the error bars are capturing should be clearly stated (for
729 example, train/test split, initialization, random drawing of some parameter, or overall
730 run with given experimental conditions).
- 731 • The method for calculating the error bars should be explained (closed form formula,
732 call to a library function, bootstrap, etc.)
- 733 • The assumptions made should be given (e.g., Normally distributed errors).

- 734
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
 - 735
 - 736 • It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
 - 737
 - 738
 - 739 • For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
 - 740
 - 741
 - 742 • If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.
 - 743

744 8. Experiments compute resources

745 Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

746 Answer: [No]

747 Justification: The computing resources used in every experiments is detailed in appendix, but we dont prvide memory and time used in training or testing.

748 Guidelines:

- 749 • The answer NA means that the paper does not include experiments.
- 750
- 751 • The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- 752
- 753 • The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- 754
- 755 • The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).
- 756
- 757
- 758
- 759

760 9. Code of ethics

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

762 Answer: [Yes]

763 Justification: We conform to the NeurIPS Code of Ethics.

764 Guidelines:

- 765 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 766
- 767 • If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- 768
- 769 • The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
- 770

771 10. Broader impacts

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

773 Answer: [NA]

774 Justification: The study doesn't has societal element to be discussed.

775 Guidelines:

- 776 • The answer NA means that there is no societal impact of the work performed.
- 777
- 778 • If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- 779
- 780 • Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- 781
- 782
- 783

- 784
- 785
- 786
- 787
- 788
- 789
- 790
- 791
- 792
- 793
- 794
- 795
- 796
- 797
- 798
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

799 **11. Safeguards**

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

803 Answer: [NA]

804 Justification: It's is not applicable.

805 Guidelines:

- 806
- 807
- 808
- 809
- 810
- 811
- 812
- 813
- 814
- 815
- The answer NA means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

816 **12. Licenses for existing assets**

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

820 Answer: [Yes]

821 Justification: The datasets originated by [19] is properly referenced.

822 Guidelines:

- 823
- 824
- 825
- 826
- 827
- 828
- 829
- 830
- 831
- 832
- 833
- 834
- 835
- The answer NA means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

836 • If this information is not available online, the authors are encouraged to reach out to
837 the asset’s creators.

838 **13. New assets**

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

841 Answer: [No]

842 Justification: The study does not provide new assets.

843 Guidelines:

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

852 **14. Crowdsourcing and research with human subjects**

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

856 Answer: [No]

857 Justification: No crowdsourcing and research with human subjects were conducted.

858 Guidelines:

- 859 • The answer NA means that the paper does not involve crowdsourcing nor research with
860 human subjects.
- 861 • Including this information in the supplemental material is fine, but if the main contribu-
862 tion of the paper involves human subjects, then as much detail as possible should be
863 included in the main paper.
- 864 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
865 or other labor should be paid at least the minimum wage in the country of the data
866 collector.

867 **15. Institutional review board (IRB) approvals or equivalent for research with human
868 subjects**

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

873 Answer: [No]

874 Justification: We did not conduct research with human subjects.

875 Guidelines:

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

886 **16. Declaration of LLM usage**

887 Question: Does the paper describe the usage of LLMs if it is an important, original, or
888 non-standard component of the core methods in this research? Note that if the LLM is used
889 only for writing, editing, or formatting purposes and does not impact the core methodology,
890 scientific rigorousness, or originality of the research, declaration is not required.

891 Answer: [No]

892 Justification: LLMs were not used in the core methods of this research.

893 Guidelines:

- 894 • The answer NA means that the core method development in this research does not
895 involve LLMs as any important, original, or non-standard components.
- 896 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)
897 for what should or should not be described.