Mixture-of-Mamba: Enhancing Multi-Modal State-Space Models with Modality-Aware Sparsity

Firstname1 Lastname1^{*1} Firstname2 Lastname2^{*12} Firstname3 Lastname3² Firstname4 Lastname4³ Firstname5 Lastname5¹

Abstract

State Space Models (SSMs) have emerged as efficient alternatives to Transformers for sequential modeling, but their inability to leverage modalityspecific features limits their performance in multimodal pretraining. Here, we propose Mixtureof-Mamba, a novel SSM architecture that introduces modality-aware sparsity through modalityspecific parameterization of the Mamba block. Building on Mixture-of-Transformers (W. Liang et al. arXiv:2411.04996; 2024), we extend the benefits of modality-aware sparsity to SSMs while preserving their computational efficiency. We evaluate Mixture-of-Mamba across three multimodal pretraining settings: Transfusion (interleaved text and continuous image tokens with diffusion loss), Chameleon (interleaved text and discrete image tokens), and an extended three-modality framework incorporating speech. Mixture-of-Mamba consistently reaches the same loss values at earlier training steps with significantly reduced computational costs. In the Transfusion setting, Mixture-of-Mamba achieves equivalent image loss using only 34.76% of the training FLOPs at the 1.4B scale. In the Chameleon setting, Mixture-of-Mamba reaches similar image loss with just 42.50% of the FLOPs at the 1.4B scale, and similar text loss with just 65.40% of the FLOPs. In the three-modality setting, MoM matches speech loss at 24.80% of the FLOPs at the **1.4B** scale. Our ablation study highlights the synergistic effects of decoupling projection components, where joint decoupling yields greater gains than individual modifications. These results establish modality-aware sparsity as a versatile and effective design principle, extending its impact from Transformers to SSMs and setting new benchmarks in multi-modal pretraining.

1. Introduction

State Space Models (SSMs) (Gu et al., 2021; Gu & Dao, 2023) have emerged as efficient alternatives to Transformers for sequential modeling, offering linear scaling in sequence length and strong performance in single-modality tasks. Mamba, a recent SSM variant, has demonstrated exceptional efficiency and scalability across diverse tasks by leveraging advanced gating mechanisms and selective state-space scanning (Gu & Dao, 2023). Despite these advantages, SSMs, including Mamba, remain inherently dense, applying the same set of parameters across all input tokens, regardless of modality. This uniform parameterization limits their ability to capture modality-specific features, leading to suboptimal performance in multi-modal pretraining.

Recent efforts have extended SSMs to multi-modal tasks. Works like VLMamba (Qiao et al., 2024) and Cobra (Zhao et al., 2024) augment Mamba for vision-language modeling by adding LLaVA-style projection modules that map image features into the token space of Mamba. In the vision domain, Vision Mamba (Zhu et al., 2024) and VMamba (Liu et al., 2024c) incorporate bidirectional scanning schemes and selective 2D scanning paths for image patch modeling. Similarly, Mamba has been explored for diffusion-based image and video generation, as seen in DiffuSSM (Yan et al., 2024) and Zigma (Hu et al., 2024), which employ unique state-space scanning patterns. While these approaches demonstrate the adaptability of Mamba, they are orthogonal to our focus, which introduces **modality-aware sparsity** directly into the Mamba block itself.

A promising approach to address such limitations is **model sparsity**, exemplified by Mixture-of-Experts (MoE) (Jacobs et al., 1991; Eigen et al., 2013; Shazeer et al., 2017; Lepikhin et al., 2020; Fedus et al., 2022; Jiang et al., 2024; Sukhbaatar et al., 2024). MoE reduces computational load by activating only a subset of model components for each input token, allowing experts to specialize in specific aspects of the data. Despite its potential, MoE-based architec-

053 054 Preprint

 ¹Department of XXX, University of YYY, Location, Country
 ²Company Name, Location, Country ³School of ZZZ, Institute
 of WWW, Location, Country. Correspondence to: Anonymous
 Author <anon.email@domain.com>.

tures face challenges such as imbalanced expert utilization,
bi-level optimization instability, and inefficient load balancing (Shazeer et al., 2017; Lepikhin et al., 2020; Fedus et al.,
2022). These issues motivate the need for alternative sparse
architectures that are computationally efficient and easier to
optimize.

061 In multi-modal contexts, prior work (Bao et al., 2022b; 062 Wang et al., 2022; Shen et al., 2023; Lin et al., 2024) has 063 introduced modality-aware sparsity in Transformer-based 064 MoE architectures. These approaches activate specific ex-065 perts or parameters based on modality, enabling models to 066 specialize in handling diverse data types. Other methods 067 fine-tune modality-specific modules atop dense LLM back-068 bones (Wang et al., 2023; He et al., 2024). Such methods 069 show that simple rule-based modality routing often out-070 performs learned routing, likely due to improved training stability and reduced optimization challenges. 072

The closest work to our approach is MoE-Mamba (Pióro 074 et al., 2024) and the related Blackmamba architecture (An-075 thony et al., 2024), which interleave Mamba blocks with 076 MoE-augmented MLP layers. While effective, these hy-077 brid designs apply sparsity only to the MLP layers, leav-078 ing the dense Mamba blocks unmodified. In contrast, we 079 present Mixture-of-Mamba, a novel architecture that directly introduces modality-aware sparsity into the Mamba 081 block itself. Inspired by Mixture-of-Transformers (Liang 082 et al., 2024), our approach dynamically selects modality-083 specific weights in every input processing component of Mamba, enabling stable and efficient multi-modal pretraining. Furthermore, prior work (Liang et al., 2024) shows 086 that MoE techniques can complement sparse architectures 087 like Mixture-of-Transformers, suggesting that Mixture-of-088 Mamba and MoE-based MLP sparsification can be com-089 bined to achieve further gains.

To rigorously evaluate Mixture-of-Mamba, we conduct experiments across three multi-modal pretraining settings:

093

094

095

096

097

098

099

100

104

105

106

109

- **Transfusion:** Interleaved text and continuous image tokens with distinct autoregressive and diffusion-based objectives. Mixture-of-Mamba achieves equivalent image loss using only **34.76%** of the training FLOPs at the **1.4B** scale.
- Chameleon: Interleaved text and discrete image tokens. Mixture-of-Mamba reaches similar image loss with just 42.50% of the FLOPs and similar text loss with only 65.40% of the FLOPs at the 1.4B scale.
- **Three-Modality:** Extension of the Chameleon setting to include speech. Mixture-of-Mamba matches speech loss using only **24.80%** of the FLOPs at the **1.4B** scale, while maintaining strong performance across image and text modalities.



Figure 1. **Multi-modal pretraining on interleaved text and image data.** Training loss on the image modality is shown for models with 1.4B parameters: Mamba Dense (cyan), Flex-Attention Transformer (dark gray), and Mixture-of-Mamba (orange). The Mixture-of-Mamba achieves significantly lower training loss and requires **2.5x fewer training steps** (indicated by the green arrow) to reach the same loss level as the other baselines.

Additionally, we perform an ablation study to analyze the contribution of modality-specific parameterization. Our findings reveal a synergistic effect: jointly decoupling all components yields greater gains than individual modifications, underscoring the importance of modality-aware sparsity as a holistic design principle.

In summary, Mixture-of-Mamba establishes a versatile and efficient architecture for SSMs by extending *modality-aware sparsity* into the Mamba block. This approach delivers robust performance gains and substantial computational savings across diverse multi-modal settings, setting new benchmarks in scalable multi-modal pretraining.

2. Method

2.1. The Mixture-of-Mamba Block

Our hypothesis is that explicitly parametrizing the selection in SSMs with the modality can improve the data efficiency of multi-modality training (Liang et al., 2024).

Following the setting of other SSMs (Gu et al., 2021), Mixture-of-Mamba is composed of homogeneous Mixtureof-Mamba blocks (line 1-13 of Algorithm 1).

In Mixture-of-Mamba, modality-specific parameterization is applied to all projections that explicitly process input features belonging to a single modality, including input projection (\mathbf{O} W_{in_proj}), intermediate projections (\mathbf{O} W_{x_proj} and \mathbf{O} W_{dt_proj}), and output projection (\mathbf{O} W_{out_proj}). Conv1D and state transitions A remain shared because they operate across multiple features or on aggregated RNN-like states, where the notion of modality is not well-defined. After parametrized by modality M, the linear transformation XW + b becomes $\mathcal{M}(X, W, b; M)$. \mathcal{M} applies the weight of modality m (W_m) to tokens of modality m (X_m) in par-

Mixture-of-Mamba: Enhancing Multi-Modal State-Space Models with Modality-Aware Sparsity

Algorithm 1 Mixture-of-Mamba block	
input $F_{in}, A, W_{in_proj}, W_{x_proj}, W_{dt_p}$	$_{proj}, W_{out_proj}, b, M$
output F _{out}	
1: $x, z \leftarrow \mathcal{M}(F_{in}, W_{in_proj}; M)$	▷ Block starts
2: $u \leftarrow \text{SiLU}(\text{Conv1D}(x))$	⊳ [b,ℓ,d]
3: $\delta, B, C \leftarrow \mathcal{M}(u, W_{x_proj}; M)$	⊳ [b,ℓ,(r,n,n)]
4: $\Delta \leftarrow \log(1 + \exp((\mathcal{M}(\delta, W_{dt_proj}))))$	(b; M))))
5: $\overline{A} \leftarrow \Delta * A$	⊳ [b,ℓ,d,n]
6: $\overline{B} \leftarrow \Delta * (u \times B)$	⊳ [b,ℓ,d,n]
7: $h = 0$	⊳ [b,d,n]
8: for $i = 0N - 1$ do	
9: $h = h * \overline{A_i} + \overline{B_i}$	⊳ [b,d,n]
10: $y_i = h \cdot C_i$	⊳[b,d]
11: end for	
12: $o \leftarrow (y+u) * \operatorname{SiLU}(z)$	
13: $F_{out} \leftarrow \mathcal{M}(o, W_{out, proj}; M)$	▷ Block ends
14:	
15: function $\mathcal{M}(X, W, b = \text{None}; M)$	
16: for each modality $m \in M$ do	
17: $I_m \leftarrow \{i : m_i = m\}$	
18: $X_m \leftarrow \{x_i : i \in I_m\}$	
19: $Y_m \leftarrow X_m W_m + b_m$	
20: end for	
21: return $Y \leftarrow \bigcup_{m \in M} Y_m$	
22. end function	

allel based on the modality mask. The output shape of \mathcal{M} is the same as the corresponding linear transformation.

136

138

139

140 The shape of W_{in_proj} is [f,(d,d)] where f is the feature 141 dimension of input F_{in} and d is the expanded feature dimen-142 sion. These two projections are fused together for efficiency 143 and W_{x_proj} uses the same technique. Line 1, 12 and 13 144 can be viewed as a SwiGLU (Shazeer, 2020) around the 145 conv+SSM (Line 2-12). x is passed to conv+SSM and z 146 will be transformed to the gate in SwiGLU.

The Conv1D in Line 2 can help collect local information across time as observed in (Sun et al., 2024). Similarly, Conv1D can also gather local information across modalities and we keep the weight-sharing property of convolution without separating the convolution kernel into different modalities.

Line 3-12 is multi-modality selective SSM. It is composedof parameter preparation (line 3-6), RNN update (line 7-11),and residual connection (line 12).

¹⁵⁷ Δ is the discretization time step. It is derived from u through a low-rank approximation $u \rightarrow \delta \rightarrow \Delta$ followed by a softplus as shown in Line 3 and 4. *A* is of shape [d,n] and Δ is of shape [b, ℓ ,d] where b is batch size, ℓ is sequence length, and n is the state dimension. Line 5 is a broadcast element-wise multiplication where Δ is unsqueezed to [b, ℓ ,d,1] and repeated to [b, ℓ ,d,n]. Line 6 first applies



Figure 2. Comparison of (a) the original Mamba block and (b) the proposed Mixture-of-Mamba block. In Mixture-of-Mamba, modality-specific parameterization is applied to all projections that explicitly process input features belonging to a single modality, including input projection (\mathbf{O} W_{in_proj}), intermediate projections (\mathbf{O} W_{x_proj} and \mathbf{O} W_{dt_proj}), and output projection (\mathbf{O} W_{out_proj}). Conv1D and state transitions A remain shared because they operate across multiple features or on aggregated RNN-like states, where the notion of modality is not well-defined. By selectively decoupling these projections, Mixture-of-Mamba enables modalityaware sparsity without compromising computational efficiency.

a batched outer product between u [b, ℓ ,d] and B [b, ℓ ,n] whose result is element-wise multiplied with Δ . Line 5 and 6 apply the selection to A,B and get $\overline{A},\overline{B}$, respectively. \overline{B} can be viewed as a gated input u and \overline{A} can be viewed as a selection gate on the state h.

Line 7-10 is a typical RNN operator with state h and output y_i . The y_i 's are concatenated together as output y. The gate application on input u is fused with gate parameter preparation at line 6 for efficiency.

Line 12 first adds the input u to the output y as residual, which is the final output of SSM. Then, Line 12 applies the gate of "SwiGLU" to the output of SSM. Finally, line 13 projects o back to the feature dimension.

2.2. Multi-objective Training with Diffusion

Following Transfusion (Zhou et al., 2024), Mixture-of-Mamba is trained on interleaved multi-modal sequences of discrete text tokens and continuous image tokens using a combined objective that incorporates both language modeling and diffusion-based image generation. Each image isencoded as a sequence of latent patches using a Variational

Autoencoder (VAE), where each patch is represented as a continuous vector. The patches are sequenced left-to-right,

169 top-to-bottom, and inserted into the discrete text sequence.

The diffusion process follows the Denoising Diffusion Probabilistic Models (DDPM) (Ho et al., 2020), where Gaussian noise is progressively added to the latent image patches during the forward process. Given a clean latent patch \mathbf{x}_0 , a noised version \mathbf{x}_t at timestep t is created as:

176

188

189

190

199

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \qquad (1)$$

where $\bar{\alpha}_t$ is determined by a cosine noise schedule (Nichol & Dhariwal, 2021), approximated as $\sqrt{\bar{\alpha}_t} \approx \cos(\frac{t}{T} \cdot \frac{\pi}{2})$ with adjustments. During training, noise is added to the latent patches at a randomly selected timestep t, and the model is optimized to predict the noise ϵ .

The overall training objective combines the autoregressive language modeling loss \mathcal{L}_{LM} , applied to the discrete text tokens, with the diffusion loss \mathcal{L}_{DDPM} , applied to the latent image patches:

$$\mathcal{L} = \mathcal{L}_{\rm LM} + \lambda \cdot \mathcal{L}_{\rm DDPM},\tag{2}$$

where λ balances the contributions of the two losses.

Critically, the conditioning for image generation is naturally
embedded within the interleaved sequence. When denoising
image patches, the preceding tokens—including both text
describing the image and prior images—serve as context
for conditional generation. This unified approach enables
Mixture-of-Mamba to leverage the modality-aware sparsity
to efficiently model both local intra-image dependencies and
long-range inter-modal relationships across the sequence.

200 2.3. Training with Uniform Representations

201 As an alternative to the multi-objective training paradigm, 202 we explore a unified representation strategy in which both text and image modalities are represented as discrete tokens. 204 Following the Chameleon framework (Chameleon Team, 205 2024), we treat the image data as sequences of discrete to-206 kens obtained through a pre-trained VQ-VAE model (Gafni 207 et al., 2022). Specifically, each image is encoded into a 208 fixed number of tokens (e.g., 1,024) by quantizing its latent 209 features into a learned codebook. These tokens are then ar-210 ranged sequentially, similar to the processing of text tokens, 211 resulting in a uniform discrete representation across both modalities.

During training, both text and image tokens are processed
using the same autoregressive objective, where the model
learns to predict the next token in the sequence given all
previous tokens. Formally, the training objective is:

218
219
$$\mathcal{L}_{\text{uniform}} = \mathbb{E}_{\mathbf{x}_{1:T}} \left[-\log P(\mathbf{x}_t \mid \mathbf{x}_{1:t-1}) \right], \quad (3)$$

where $\mathbf{x}_{1:T}$ represents the interleaved sequence of text and image tokens. This objective allows the model to treat text and image data equivalently, unifying the training process across modalities while relying solely on an autoregressive loss. The use of discrete tokens for images simplifies the training procedure by removing the need for separate loss formulations, as in the diffusion-based approach. It also aligns with the inherent sequence-to-sequence nature of Mixture-of-Mamba, where the same modality-aware sparsity design can be applied seamlessly across the discrete text and image tokens.

Motivation and Robustness Testing. We include this alternative strategy to evaluate the robustness of our Mixtureof-Mamba architecture under different choices of training objectives and data representations. By experimenting with uniform discrete representations, we demonstrate that Mixture-of-Mamba consistently outperforms Mamba Dense models across various settings, including both continuous (multi-objective) and discrete (uniform) representations. This highlights the versatility of Mixture-of-Mamba and its ability to deliver performance gains regardless of the underlying choice of modality representations or training objectives.

3. Results

3.1. Results in Multi-objective Training (Transfusion)

We evaluate Mixture-of-Mamba (MoM) against Mamba Dense and Flex-Attention Transformer in the **Transfusion** setting, where pretraining is performed on interleaved text and image data across three model scales: 163M, 760M, and 1.4B. See our training configuration in Appendix Table 5. Performance gain is quantified as:

Performance Gain (%) =
$$\frac{\text{Loss}_{\text{Dense}} - \text{Loss}_{\text{Mixture}}}{\text{Loss}_{\text{Dense}}} \times 100,$$

where $Loss_{Dense}$ and $Loss_{Mixture}$ are the final losses of Mamba Dense and Mixture-of-Mamba, respectively. Relative training FLOPs reflect the computational cost required for MoM to match the training dynamics (similar loss) of Mamba Dense. The detailed results are summarized in Table 1 and Figure 3, with further visualizations provided in Appendix Figures 4, 5, and 6.¹

¹Flex-Attention Transformer (i.e., Transfusion (Zhou et al., 2024)) combines both attention patterns by applying causal attention to every element in the sequence and *bidirectional attention within the elements of each individual image*. This makes Flex-Attention Transformer an overestimated baseline for transformers because both Mamba and Mixture-of-Mamba are strictly causal across all elements, while Flex-Attention Transformer benefits from bidirectional attention within images.





Figure 3. Multi-modal pretraining in the Transfusion setting on interleaved text and image data across model scales. Training loss and loss matching are reported for image and text modalities at three model sizes: 1.4B, 760M, and 163M. (a, e, i) Image training loss shows significant improvements for Mixture-of-Mamba (orange), which consistently achieves lower loss compared to Mamba Dense (cyan) and Flex-Attention Transformer (dark gray) across all scales. (b, f, j) Image loss matching compares the training dynamics and shows that Mixture-of-Mamba and Flex-Attention Transformer reach the same loss values at earlier training steps compared to Mamba Dense. (c, g, k) Text training loss shows competitive results, with Mixture-of-Mamba performing better than Mamba Dense and on par with the Flex-Attention Transformer. (d, h, l) Text loss matching illustrates that Mixture-of-Mamba and Flex-Attention Transformer exhibit more efficient training dynamics than Mamba Dense, requiring fewer steps to achieve comparable loss values, though the primary improvements are observed in the image modality. Overall, in the Transfusion setting, Mixture-of-Mamba demonstrates substantial gains in image loss and training efficiency while maintaining strong performance on text.

253 Image Modality. Mixture-of-Mamba (MoM) consistently 254 demonstrates superior performance in image modality 255 training loss across all model scales. At the 1.4B scale, 256 MoM achieves a training loss of **0.2138**, outperforming 257 Mamba Dense by 2.20% while requiring only 34.76% of 258 the training FLOPs. Similar trends are observed at smaller 259 scales: at the 760M scale, MoM achieves a training loss of 0.2172, a 2.37% improvement over Mamba Dense, while 261 reducing training FLOPs to 37.76%.

263 The validation loss curves on the CC12M dataset ((Table 1, 264 Appendix Figure 5) further illustrate these trends. Mixture-265 of-Mamba consistently achieves lower image validation 266 loss compared to Mamba Dense and Flex-Attention Transformer, with the improvements becoming more pronounced 267 as model size increases. Additionally, loss matching curves 269 demonstrate that MoM reaches equivalent loss values at 270 earlier training steps, highlighting its improved training effi-271 ciency.

Text Modality. In the text modality, Mixture-of-Mamba consistently outperforms Mamba Dense across both training and validation metrics. At the 1.4B scale, MoM achieves lower validation losses on both the C4 (2.2695) and Wikipedia (1.7164) datasets compared to Mamba Dense, despite their similar training losses. This indicates better generalization to unseen text data. Importantly, MoM also performs comparably to or better than Flex-Attention Transformer, particularly on validation losses, as shown in Appendix Figure 4. Similar trends are observed at smaller scales (760M and 163M), where MoM reduces validation losses while maintaining high training efficiency.

Loss matching results in Appendix Figure 4 (b, f, j) confirm that Mixture-of-Mamba aligns closely with or surpasses Mamba Dense, reaching comparable loss values earlier during training. These improvements highlight MoM's strong performance in text tasks while maintaining its computational efficiency.

272

243

244

245

246

247

248

249

250

251 252

Mixture-of-Mamba: Enhancing Multi-Modal State-Space Models with Modality-Aware Sparsity

Model Scale	Metric Category	Metric Name	Mamba Loss (↓)	Flex-AttentionTransformerLoss (\downarrow)	Mixture-of- Mamba Loss (↓)	Performance Gain over Mamba (%) (†)	Relative Training FLOPs to Match Mamba (%) (↓)
	Image Metrics	Training Loss CC12M Val. Loss	0.2262 0.2295	0.2250 0.2293	0.2199 0.2255	2.80% 1.74%	49.21% 50.61%
163M	Text Metrics	Avg Training Loss C4 Val. Loss Wikipedia Val. Loss	2.4702 2.6917 2.1884	2.4424 2.6862 2.1715	2.4690 2.6912 2.1870	$0.05\% \\ 0.02\% \\ 0.06\%$	98.80% 99.88% 99.81%
	Overall	Train Avg Loss	3.6014	3.5674	3.5685	0.91%	86.11%
	Image Metrics	Training Loss CC12M Val. Loss	0.2225 0.2272	0.2213 0.2253	0.2172 0.2201	2.37% 3.13%	37.76% 35.27%
760M	Text Metrics	Avg Training Loss C4 Val. Loss Wikipedia Val. Loss	2.1394 2.3593 1.8191	2.1253 2.3559 1.8143	2.1353 2.3555 1.8149	0.19% 0.16% 0.23%	96.82% 99.01% 99.11%
	Overall	Train Avg Loss	3.2519	3.2318	3.2214	0.94%	82.94%
	Image Metrics	Training Loss CC12M Val. Loss	0.2186 0.2264	0.2221 0.2247	0.2138 0.2190	2.20% 3.29%	34.76% 36.15%
1.4B	Text Metrics	Avg Training Loss C4 Val. Loss Wikipedia Val. Loss	2.0761 2.2726 1.7205	2.0673 2.2728 1.7218	2.0737 2.2695 1.7164	0.12% 0.13% 0.24%	98.27% 99.34% 99.30%
	Overall	Train Avg Loss	3.1693	3.1777	3.1429	0.84%	83.10%

297 Table 1. Training and validation metrics across model scales in the Transfusion setting. Loss values are reported for image and text 298 modalities at three model sizes: 163M, 760M, and 1.4B. Mixture-of-Mamba consistently achieves competitive or superior performance 299 in image metrics and maintains strong text performance compared to Mamba Dense and Flex-Attention Transformer. The table also reports relative training FLOPs required for Mixture-of-Mamba and Flex-Attention Transformer to match Mamba's training dynamics, 300 highlighting improved training efficiency. Best loss values in each row are highlighted. 301

303 Overall Performance and Efficiency. Across both image 304 and text modalities, Mixture-of-Mamba consistently out-305 performs Mamba Dense in terms of loss reduction while 306 requiring significantly fewer training FLOPs to achieve sim-307 ilar learning dynamics. At the 1.4B scale, MoM improves 308 the overall training loss by 0.84% while requiring only 309 83.10% of the training FLOPs. At smaller scales, such as 310 760M and 163M, MoM reduces the overall training loss 311 by up to 0.94%, while requiring just 82.94% and 86.11% 312 of the FLOPs, respectively (Table 1, Appendix Figure 6). 313 These results, summarized in Table 1 and Figure 3, and 314 further supported by Appendix Figures 4, 5, and 6, under-315 scoring MoM's effectiveness, scalability, and efficiency in 316 the Transfusion setting.

2 2

302

318

319

320

3.2. Results in Training with Uniform Representations (Chameleon)

We evaluate Mixture-of-Mamba (MoM) in the Chameleon setting, where both image and text modalities are represented as discrete tokens. See our training configuration in 324 Appendix Table 6. Results are summarized in Table 2, with 325 full results across all five scales (37M, 94M, 443M, 880M, and **1.5B**) provided in Appendix Table 7. Training dynamics and validation loss trends are visualized in Appendix 328 Figures 7, 8, and 9. 329

Image Modality. Mixture-of-Mamba (MoM) consistently demonstrates better performance in image modality training loss across all model scales, achieving substantial efficiency gains over Mamba Dense. At the 443M scale, MoM achieves a training loss of 5.1703, a 3.46% improvement over Mamba Dense, while requiring only 33.40% of the training FLOPs. Similar trends are observed at other scales: at the largest 1.5B scale, MoM achieves a training loss of 5.0591, a 2.51% improvement, with only 42.50% of the training FLOPs. At the smallest 37M scale, MoM reduces training loss to 5.9561, outperforming Mamba Dense by 2.85% while requiring just 25.90% of the FLOPs (Appendix Table 7). These results highlight MoM's ability to achieve improved performance and convergence efficiency consistently in the image modality across all model scales.

Text Modality. Mixture-of-Mamba (MoM) demonstrates consistent improvements in text modality training loss across all model scales. At the largest 1.5B scale, MoM reduces training loss to 2.1614, a 3.01% improvement over Mamba Dense, while requiring only 65.40% of the training FLOPs. Validation loss on Obelisc and SSTK datasets exhibits similar trends, with MoM achieving notable improvements in loss values while maintaining significant efficiency gains (Appendix Figures 8 and 9). These results further highlight MoM's ability to deliver strong text perfor-

Mixture-of-Mamba: Enhancing Multi-Modal State-Space Models with Modality-Aware Sparsity

Model Scale	Metric Category	Metric Name	Mamba Loss (↓)	Mixture-of- Mamba Loss (↓)	Performance Gain (%) (†)	Relative Training FLOPs to Match Mamba (%) (↓)
		Training Loss	5.3558	5.1703	3.46%	33.40%
	Image Metrics	Obelisc Val. Loss	4.5258	4.3546	3.78%	35.10%
		SSTK Val. Loss	5.9179	5.7471	2.89%	35.30%
143M		Training Loss	2.4637	2.3864	3.14%	62.00%
	Text Metrics	Obelisc Val. Loss	3.0544	2.9820	2.37%	66.70%
		SSTK Val. Loss	2.7569	2.6250	4.78%	54.70%
	Overall	Avg Training Loss	3.6584	3.5364	3.33%	47.90%
		Training Loss	5.2260	5.1201	2.03%	48.40%
	Image Metrics	Obelisc Val. Loss	4.4127	4.3105	2.32%	49.30%
		SSTK Val. Loss	5.7987	5.6986	1.73%	50.50%
380M		Training Loss	2.3073	2.2438	2.75%	65.60%
	Text Metrics	Obelisc Val. Loss	2.8886	2.8313	1.99%	72.80%
		SSTK Val. Loss	2.5483	2.4548	3.67%	67.90%
	Overall	Avg Training Loss	3.5130	3.4320	2.31%	58.30%
		Training Loss	5.1892	5.0591	2.51%	42.50%
	Image Metrics	Obelisc Val. Loss	4.3692	4.2510	2.71%	44.50%
1.5B		SSTK Val. Loss	5.7546	5.6335	2.10%	44.60%
		Training Loss	2.2284	2.1614	3.01%	65.40%
	Text Metrics	Obelisc Val. Loss	2.8020	2.7393	2.24%	71.60%
		SSTK Val. Loss	2.4614	2.3455	4.71%	62.10%
	Overall	Avg Training Loss	3.4602	3.3670	2.69%	54.70%

Table 2. Training and validation metrics across model scales in the Chameleon setting. In this setting, both image and text modalities are represented as discrete tokens. Mixture-of-Mamba achieves substantial performance improvements over Mamba Dense, with the **image modality** showing the largest gains. The **text modality** also exhibits significant improvements, in contrast to the Transfusion setting where text gains were more modest. The current table shows results for three model scales: **443M**, **880M**, and **1.5B**, due to space constraints. See Appendix Table 7 for the full results across all five model scales: **37M**, **94M**, **443M**, **880M**, and **1.5B**. These results further highlight the effectiveness and efficiency of Mixture-of-Mamba, which consistently achieves strong performance with reduced relative training FLOPs.

mance with improved convergence efficiency. These results highlight Mixture-of-Mamba's robust and efficient improvements in the Chameleon setting across both image and text modalities, with substantial computational savings.

3.3. Results in Training with Three Modalities (Chameleon+Speech)

355

357

358

359

360

361

362 363

364

367

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

To evaluate the robustness and scalability of Mixture-of-Mamba (MoM), we extend the Chameleon framework to include a third modality: **speech**, alongside image and text, with all modalities represented as discrete tokens. Speech data is tokenized using an in-house tokenizer, a variant of DinoSR (Liu et al., 2024a), which extracts semantic tokens with a vocabulary size of 500, where each token corresponds to 40ms of audio content. Results are summarized in Table 3, with additional training dynamics and evaluation loss trends visualized in Appendix Figures 11, 12, 13, and 14.

Speech Modality. Mixture-of-Mamba (MoM) achieves substantial improvements in **speech modality training loss**

across all model scales. At the **443M** scale, MoM improves speech training loss by **7.14%** compared to Mamba Dense. To match the training loss achieved by Mamba Dense, MoM requires only **19.20%** of the training FLOPs, demonstrating significant efficiency gains. Similar trends hold at the largest **1.5B** scale, where MoM achieves a **5.75%** improvement in speech training loss and matches Mamba Dense's loss with just **24.80%** of the training FLOPs.

Overall training loss is consistently reduced across scales. At the **1.5B** scale, MoM lowers the overall training loss by **3.05%**. When targeting the same loss as Mamba Dense, MoM achieves this with a **56.20%** reduction in relative training FLOPs, highlighting its improved computational efficiency.

Performance in the **image** and **text** modalities similarly shows consistent improvements in training and validation losses relative to Mamba Dense. Full results and trends are presented in Appendix Figures 13 and 14, where MoM's robust performance across all three modalities is further

Model Scale	Metric Category	Metric Name	Mamba Loss (↓)	Mixture-of- Mamba Loss (↓)	Performance Gain (%) (†)	Relative Training FLOPs to Match Mamba (%) (↓)
		Training Loss	1.8159	1.6909	6.88%	10.30%
2714	Speech Metrics	LL60K Val. Loss	1.6756	1.5217	9.18%	13.60%
57111		PPL30K Val. Loss	1.8147	1.6845	7.17%	13.60%
	Overall Metrics	Avg Training Loss	4.2299	4.0759	3.64%	45.00%
		Training Loss	1.6911	1.5662	7.38%	11.90%
0414	Speech Metrics	LL60K Val. Loss	1.5235	1.3747	9.76%	14.80%
Over		PPL30K Val. Loss	1.6951	1.6152	4.71%	12.60%
	Overall Metrics	Avg Training Loss	3.7756	3.6371	3.67%	43.10%
		Training Loss	1.5414	1.4313	7.14%	19.20%
44214	Speech Metrics	LL60K Val. Loss	1.3466	1.2113	10.05%	24.70%
443M		PPL30K Val. Loss	1.5634	1.4790	5.40%	22.00%
	Overall Metrics	Avg Training Loss	3.3317	3.2096	3.66%	44.00%
		Training Loss	1.4902	1.4054	5.69%	22.40%
00014	Speech Metrics	LL60K Val. Loss	1.2939	1.1757	9.13%	30.10%
880M		PPL30K Val. Loss	1.5400	1.4619	5.07%	24.30%
	Overall Metrics	Avg Training Loss	3.2289	3.1571	2.22%	54.30%
1.5B		Training Loss	1.4790	1.3940	5.75%	24.80%
	Speech Metrics	LL60K Val. Loss	1.2592	1.1552	8.26%	32.10%
		PPL30K Val. Loss	1.5200	1.4387	5.35%	27.60%
	Overall Metrics	Avg Training Loss	3.1507	3.0545	3.05%	56.20%

Table 3. Training and validation metrics across model scales with three modalities: image, text, and speech. This setting extends the Chameleon framework by incorporating speech alongside image and text, with all modalities represented as discrete tokens. Mixture-of-Mamba achieves consistent improvements over Mamba Dense across all scales (37M, 94M, 443M, 880M, and 1.5B), particularly in the speech modality, where performance gains reach up to 9.18%. These gains are achieved with substantial reductions in training FLOPs, ranging from 10.30% to 56.20% relative to Mamba Dense. The results demonstrate that Mixture-of-Mamba generalizes effectively to a multi-modal setting with three modalities while delivering significant computational efficiency.

validated.

408

409

410

411

412

413

414 415

416

417

418

419

420

421

422

423

424

425

426

427

3.4. Ablation Study on Decoupling Components

To better understand the design choices underpinning Mixture-of-Mamba, we conduct an ablation study on the Chameleon + Speech setting at the 443M scale. We evaluate the impact of decoupling four key components— $W_{\text{in-proj}}$ (**0**), $W_{\text{x-proj}}$ (**2**), $W_{\text{dt-proj}}$ (**3**), and $W_{\text{out-proj}}$ (O)—individually and in various combinations. This analysis enables us to test both individual and combined contributions to the model's overall performance.

The results show that decoupling components individually 428 yields varying degrees of improvement, with performance 429 gains ranging from 0.63% ($W_{\text{out-proj}}$) to 1.22% ($W_{\text{in-proj}}$). 430 Interestingly, some components $(W_{x-proj} \text{ and } W_{dt-proj})$ ex-431 hibit minimal or even slightly negative impact when de-432 coupled alone. However, decoupling multiple components 433 in combination leads to significantly larger gains. For ex-434 ample, decoupling $W_{\text{in-proj}}$ and $W_{\text{out-proj}}$ ($\mathbf{0+0}$) achieves a 435 2.20% improvement, while decoupling three components 436 $(\mathbf{0}+\mathbf{2}+\mathbf{4})$ further increases the gain to 3.11%. 437

438 Most importantly, decoupling all four components simul-439

taneously (**0**+**0**+**3**+**4**, Mixture-of-Mamba) achieves the largest improvement, with a performance gain of 3.80% over the Mamba baseline. This result highlights a key observation: the gain from decoupling all components together exceeds the sum of individual gains, demonstrating a synergistic effect. The combination of all decoupled projections enables better parameter allocation across modalities, leading to more efficient and effective learning. In summary, the ablation study confirms that the design of Mixture-of-Mamba is both effective and interdependent. Decoupling all key components simultaneously is critical to achieving the observed substantial performance gains.

4. Related Work

4.1. State-Space Models and Multi-Modal Extensions

State-space models (SSMs) (Gu et al., 2021; Gu & Dao, 2023) have recently gained traction as computationally efficient alternatives to Transformers for sequential modeling. Mamba (Gu & Dao, 2023), in particular, demonstrates strong performance on single-modality tasks by leveraging linear time complexity and advanced gating mechanisms. Extending Mamba to multi-modal tasks remains an active

Ablation Study	Avg Training Loss (\downarrow)	Performance Gain (%) (†)	
443M Mamba (without 0239)	3.3317	0% (baseline)	
\mathbf{O} (decouple W_{in_proj})	3.2916	1.22%	
\mathbf{Q} (decouple W_{x_proj})	3.3580	-0.79%	
$(\text{decouple } W_{dt,proj})$	3.3525	-0.62%	
$(\text{decouple } W_{out_proj})$	3.3109	0.63%	
1 + 2 (decouple W_{in_proj}, W_{x_proj})	3.2780	1.64%	
1+3 (decouple W_{in_proj}, W_{dt_proj})	3.2687	1.93%	
0+9 (decouple $W_{in_proj}, W_{out_proj}$)	3.2599	2.20%	
Q+O (decouple W_{x_proj}, W_{dt_proj})	3.3214	0.31%	
2+4 (decouple W_{x_proj}, W_{out_proj})	3.2829	1.49%	
$(\text{decouple } W_{dt_proj}, W_{out_proj})$	3.2509	2.48%	
1 + 2 + 3 (not decoupling W_{out_proj})	3.2593	2.22%	
1 + 2 + 4 (not decoupling W_{dt_proj})	3.2312	3.11%	
1 + 3 + 4 (not decoupling W_{x_proj})	3.2342	3.01%	
$2+3+4$ (not decoupling W_{in_proj})	3.2773	1.66%	
1 + 2 + 3 + 3 (Mixture-of-Mamba)	3.2096	3.80%	

Table 4. Ablation study on the Chameleon + Speech setting. This study evaluates the impact of decoupling individual components (1, 2, 3, 4) and their combinations on model performance. The results demonstrate that decoupling all components (1+2+3+4, Mixture-of-Mamba) achieves the best performance with a 3.80%gain over the Mamba baseline. Notably, the performance gain achieved by decoupling all components together exceeds the sum of gains from decoupling each component individually, highlighting the synergistic effect of combined decoupling. Green shading indicates positive performance gains, with the darkest green highlighting the best configuration.

research area.

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468 In vision-language modeling, VLMamba (Qiao et al., 2024) 469 and Cobra (Zhao et al., 2024) augment Mamba by incor-470 porating LLaVA-style projection modules, enabling image 471 features to be mapped into the token space of the Mamba 472 model for sequence modeling. In the vision domain, Vision 473 Mamba (Zhu et al., 2024) introduces bidirectional scan-474 ning by chaining forward and backward SSM blocks, while 475 VMamba (Liu et al., 2024c) further enhances image patch 476 processing with a 2D Selective Scan (SS2D) module that 477 traverses patches across multiple scanning paths. 478

For diffusion-based models, works such as DiffuSSM (Yan 479 et al., 2024) and Zigma (Hu et al., 2024) replace attention 480 mechanisms with SSMs for image and video generation. 481 Zigma introduces a zigzag scanning scheme to improve 482 efficiency for sequential diffusion tasks, while other ap-483 proaches (Mo & Tian, 2024; Fei et al., 2024) explore bi-484 directional SSM architectures. While these works highlight 485 the flexibility of Mamba in generative tasks, they focus pri-486 marily on architectural modifications for specific domains 487 rather than general multi-modal pretraining. 488

The most related work to ours is MoE-Mamba (Pióro et al., 2024) and Blackmamba (Anthony et al., 2024), which interleave Mamba blocks with MoE-augmented MLPs to introduce sparsity. However, these hybrid designs apply sparsity only to the MLP layers, leaving the dense Mamba block

unmodified. In contrast, our proposed Mixture-of-Mamba integrates modality-aware sparsity directly into the Mamba block by decoupling its projection components, enabling specialized computations for different modalities. This general design complements existing methods and offers new opportunities for computationally efficient multi-modal pretraining.

4.2. Sparse Architectures for Multi-Modal Pretraining

Model sparsity, particularly Mixture-of-Experts (MoE), has been extensively explored in Transformers to reduce computational cost (Jacobs et al., 1991; Eigen et al., 2013; Shazeer et al., 2017; Lepikhin et al., 2020; Fedus et al., 2022; Jiang et al., 2024). MoE selectively activates subsets of parameters for each input token, allowing the model to specialize in different aspects of the data. However, challenges such as expert imbalance, bi-level optimization, and load balancing remain prevalent (Shazeer et al., 2017; Lepikhin et al., 2020).

In multi-modal tasks, modality-aware sparsity has emerged as an effective strategy. Works such as VLMo (Shen et al., 2023), MoMA (Lin et al., 2024), and related approaches (Wang et al., 2022; Bao et al., 2022a; Long et al., 2023) assign modality-specific experts to handle the unique statistical properties of text, images, and other data types. This improves specialization while avoiding the complexities of learned routing mechanisms (Liang et al., 2022).

Transformer-based architectures have further extended sparsity into attention mechanisms (Wang et al., 2023; Shen et al., 2024; Liu et al., 2024b). CogVLM (Wang et al., 2023) applies sparse techniques on top of a pre-trained Vicuna-7B model but remains limited to generating text outputs. Concurrently, Playground v3 (PGv3) (Liu et al., 2024b) integrates DiT-style image transformers with a frozen LLaMA-3 backbone to achieve state-of-the-art performance in text-toimage generation.

Our work differs fundamentally in two key aspects. First, Mixture-of-Mamba introduces *modality-aware sparsity* into the Mamba block itself, generalizing sparse architectures beyond Transformers to SSMs. Unlike prior works that sparsify only the MLP or attention components, we decouple projection components of the Mamba block, enabling efficient and specialized computations across modalities. Second, Mixture-of-Mamba is trained from scratch for multimodal generation tasks, unlike approaches like CogVLM and PGv3 that fine-tune pre-trained backbones.

Furthermore, our design is complementary to existing MoE techniques. Prior work (Liang et al., 2024) has demonstrated that MoE-based sparsification can be combined with sparse architectures like Mixture-of-Transformers to achieve additional gains. Similarly, Mixture-of-Mamba can serve as

495 a versatile and computationally efficient solution, offering

496 new pathways for scalable multi-modal pretraining.497

⁴⁹⁸ 5. Conclusion

In this work, we introduced Mixture-of-Mamba, a novel extension of state-space models (SSMs) that incorporates modality-aware sparsity through modality-specific param-eterization. By enabling modality-specific specialization while preserving the computational efficiency of SSMs, Mixture-of-Mamba consistently outperforms dense base-lines across three multi-modal settings: Transfusion (in-terleaved text and continuous image tokens), Chameleon (interleaved text and discrete image tokens), and an extended Chameleon+Speech framework. Our results demonstrate substantial improvements in loss reduction, with training efficiency gains reaching more than double the computa-tional efficiency compared to dense SSMs. Ablation studies further reveal a synergistic effect from jointly decoupling key projection components, highlighting the effectiveness of modality-aware sparsity. These findings establish Mixture-of-Mamba as a scalable and efficient architecture for multi-modal pretraining, paving the way for future exploration in dynamic sparsity and broader multi-modal applications.

550 References

567

568

569

573

574

575

576

577

581

582

583

584

- 551 Anthony, Q., Tokpanov, Y., Glorioso, P., and Millidge, B. 552 Blackmamba: Mixture of experts for state-space models. 553 arXiv preprint arXiv:2402.01771, 2024. 554
- 555 Bao, H., Wang, W., Dong, L., Liu, O., Mohammed, O. K., 556 Aggarwal, K., Som, S., Piao, S., and Wei, F. Vlmo: 557 Unified vision-language pre-training with mixture-of-558 modality-experts. Advances in Neural Information Pro-559 cessing Systems, 35:32897-32912, 2022a.
- 560 Bao, H., Wang, W., Dong, L., Liu, Q., Mohammed, O. K., 561 Aggarwal, K., Som, S., and Wei, F. Vlmo: Unified 562 vision-language pre-training with mixture-of-modality-563 experts, 2022b. URL https://arxiv.org/abs/ 564 2111.02358. 565
- 566 Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models, 2024. URL https://arxiv. org/abs/2405.09818.
- Eigen, D., Ranzato, M., and Sutskever, I. Learning fac-570 tored representations in a deep mixture of experts. arXiv 571 preprint arXiv:1312.4314, 2013. 572
 - Fedus, W., Zoph, B., and Shazeer, N. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity, 2022. URL https://arxiv.org/ abs/2101.03961.
- Fei, Z., Fan, M., Yu, C., and Huang, J. Scalable diffu-578 579 sion models with state space backbone. arXiv preprint arXiv:2402.05608, 2024. 580
 - Gafni, O., Polyak, A., Ashual, O., Sheynin, S., Parikh, D., and Taigman, Y. Make-a-scene: Scene-based textto-image generation with human priors. arXiv preprint arXiv:2203.13131, 2022.
- 586 Gu, A. and Dao, T. Mamba: Linear-time sequence 587 modeling with selective state spaces. arXiv preprint 588 arXiv:2312.00752, 2023.
- 589 Gu, A., Goel, K., and Ré, C. Efficiently modeling long 590 sequences with structured state spaces. arXiv preprint 591 arXiv:2111.00396, 2021. 592
- 593 He, W., Fu, S., Liu, M., Wang, X., Xiao, W., Shu, F., Wang, 594 Y., Zhang, L., Yu, Z., Li, H., et al. Mars: Mixture of 595 auto-regressive models for fine-grained text-to-image syn-596 thesis. arXiv preprint arXiv:2407.07614, 2024. 597
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion proba-598 bilistic models. Advances in neural information process-599 ing systems, 33:6840-6851, 2020. 600
- 601 Hu, V. T., Baumann, S. A., Gui, M., Grebenkova, O., Ma, 602 P., Schusterbauer, J., and Ommer, B. Zigma: A dit-style 603 zigzag mamba diffusion model. In ECCV, 2024. 604

- Jacobs, R. A., Jordan, M. I., Nowlan, S. J., and Hinton, G. E. Adaptive mixtures of local experts. Neural computation, 3(1):79-87, 1991.
- Jiang, A. Q., Sablayrolles, A., Roux, A., Mensch, A., Savary, B., Bamford, C., Chaplot, D. S., de las Casas, D., Hanna, E. B., Bressand, F., Lengyel, G., Bour, G., Lample, G., Lavaud, L. R., Saulnier, L., Lachaux, M.-A., Stock, P., Subramanian, S., Yang, S., Antoniak, S., Scao, T. L., Gervet, T., Lavril, T., Wang, T., Lacroix, T., and Sayed, W. E. Mixtral of experts, 2024. URL https://arxiv. org/abs/2401.04088.
- Langley, P. Crafting papers on machine learning. In Langley, P. (ed.), Proceedings of the 17th International Conference on Machine Learning (ICML 2000), pp. 1207-1216, Stanford, CA, 2000. Morgan Kaufmann.
- Lepikhin, D., Lee, H., Xu, Y., Chen, D., Firat, O., Huang, Y., Krikun, M., Shazeer, N., and Chen, Z. Gshard: Scaling giant models with conditional computation and automatic sharding, 2020. URL https://arxiv.org/abs/ 2006.16668.
- Liang, V. W., Zhang, Y., Kwon, Y., Yeung, S., and Zou, J. Y. Mind the gap: Understanding the modality gap in multi-modal contrastive representation learning. Advances in Neural Information Processing Systems, 35: 17612-17625, 2022.
- Liang, W., Yu, L., Luo, L., Iyer, S., Dong, N., Zhou, C., Ghosh, G., Lewis, M., Yih, W.-t., Zettlemoyer, L., et al. Mixture-of-transformers: A sparse and scalable architecture for multi-modal foundation models. arXiv preprint arXiv:2411.04996, 2024.
- Lin, X. V., Shrivastava, A., Luo, L., Iyer, S., Lewis, M., Gosh, G., Zettlemoyer, L., and Aghajanyan, A. Moma: Efficient early-fusion pre-training with mixture of modality-aware experts. arXiv preprint arXiv:2407.21770, 2024.
- Liu, A. H., Chang, H.-J., Auli, M., Hsu, W.-N., and Glass, J. R. Dinosr: Self-distillation and online clustering for self-supervised speech representation learning, 2024a. URL https://arxiv.org/abs/2305.10005.
- Liu, B., Akhgari, E., Visheratin, A., Kamko, A., Xu, L., Shrirao, S., Lambert, C., Souza, J., Doshi, S., and Li, D. Playground v3: Improving text-to-image alignment with deep-fusion large language models, 2024b. URL https://arxiv.org/abs/2409.10695.
- Liu, Y., Tian, Y., Zhao, Y., Yu, H., Xie, L., Wang, Y., Ye, Q., Jiao, J., and Liu, Y. VMamba: Visual state space model. In The Thirty-eighth Annual Conference on Neural Information Processing Systems, 2024c. URL https: //openreview.net/forum?id=ZgtLQQR1K7.

659

- Long, Z., Killick, G., McCreadie, R., and Camarasa, G. A.
 Multiway-adapater: Adapting large-scale multi-modal
 models for scalable image-text retrieval. *arXiv preprint arXiv:2309.01516*, 2023.
- Mo, S. and Tian, Y. Scaling diffusion mamba with bidirectional ssms for efficient image and video generation. *arXiv preprint arXiv:2405.15881*, 2024.
- Nichol, A. Q. and Dhariwal, P. Improved denoising diffusion probabilistic models. In *International conference on machine learning*, pp. 8162–8171. PMLR, 2021.
- Pióro, M., Ciebiera, K., Król, K., Ludziejewski, J., and Jaszczur, S. Moe-mamba: Efficient selective state space models with mixture of experts, 2024.
- Qiao, Y., Yu, Z., Guo, L., Chen, S., Zhao, Z., Sun, M., Wu, Q., and Liu, J. Vl-mamba: Exploring state space models for multimodal learning. *arXiv preprint arXiv:2403.13600*, 2024.
- Shazeer, N. Glu variants improve transformer. *arXiv* preprint arXiv:2002.05202, 2020.
- Shazeer, N., Mirhoseini, A., Maziarz, K., Davis, A., Le, Q. V., Hinton, G. E., and Dean, J. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *CoRR*, abs/1701.06538, 2017. URL http:// arxiv.org/abs/1701.06538.
- Shen, S., Yao, Z., Li, C., Darrell, T., Keutzer, K., and He, Y. Scaling vision-language models with sparse mixture of experts. *arXiv preprint arXiv:2303.07226*, 2023.
- Shen, Y., Guo, Z., Cai, T., and Qin, Z. Jetmoe: Reaching llama2 performance with 0.1 m dollars. *arXiv preprint arXiv:2404.07413*, 2024.
- Sukhbaatar, S., Golovneva, O., Sharma, V., Xu, H., Lin, X. V., Rozière, B., Kahn, J., Li, D., tau Yih, W., Weston, J., and Li, X. Branch-train-mix: Mixing expert llms into a mixture-of-experts llm, 2024. URL https://arxiv. org/abs/2403.07816.
- Sun, Y., Li, X., Dalal, K., Xu, J., Vikram, A., Zhang, G.,
 Dubois, Y., Chen, X., Wang, X., Koyejo, S., et al. Learning to (learn at test time): Rnns with expressive hidden states. *arXiv preprint arXiv:2407.04620*, 2024.
- Wang, W., Bao, H., Dong, L., Bjorck, J., Peng, Z., Liu, Q.,
 Aggarwal, K., Mohammed, O. K., Singhal, S., Som, S.,
 and Wei, F. Image as a foreign language: Beit pretraining
 for all vision and vision-language tasks, 2022. URL
 https://arxiv.org/abs/2208.10442.
- Wang, W., Lv, Q., Yu, W., Hong, W., Qi, J., Wang, Y., Ji, J., Yang, Z., Zhao, L., Song, X., et al. Cogvlm: Visual

expert for pretrained language models. *arXiv preprint arXiv:2311.03079*, 2023.

- Yan, J. N., Gu, J., and Rush, A. M. Diffusion models without attention. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pp. 8239– 8249, 2024.
- Zhao, H., Zhang, M., Zhao, W., Ding, P., Huang, S., and Wang, D. Cobra: Extending mamba to multi-modal large language model for efficient inference. *arXiv preprint arXiv:2403.14520*, 2024.
- Zhou, C., Yu, L., Babu, A., Tirumala, K., Yasunaga, M., Shamis, L., Kahn, J., Ma, X., Zettlemoyer, L., and Levy, O. Transfusion: Predict the next token and diffuse images with one multi-modal model. *arXiv preprint arXiv:2408.11039*, 2024.
- Zhu, L., Liao, B., Zhang, Q., Wang, X., Liu, W., and Wang, X. Vision mamba: Efficient visual representation learning with bidirectional state space model. In *Proceedings* of the 41st International Conference on Machine Learning, Proceedings of Machine Learning Research. PMLR, 2024.

A. Appendix

Model Size	Hidden Dim.	Layers	Heads	Seq. Length	Batch Size/GPU	GPUs	Tokens/Batch	Steps
163M	768	16	12	4,096	4	56	1,048,576	250,000
760M	1,536	24	24	4,096	4	56	1,048,576	250,000
1.4B	2,048	24	16	4,096	2	128	1,048,576	250,000

Table 5. Architectural specifications and training configurations of models across different parameter scales (Transfusion setting).

Model Size	Hidden Dim.	Layers	Heads	Seq. Length	Batch Size/GPU	GPUs	Tokens/Batch	Steps
37M	256	4	8	4,096	2	64	524,288	160,000
94M	512	8	8	4,096	2	64	524,288	160,000
443M	1,024	24	16	4,096	2	64	524,288	160,000
880M	1,536	24	24	4,096	2	64	524,288	120,000
1.5B	2,048	24	16	4,096	1	128	524,288	120,000

Table 6. Architectural specifications and training configurations of models across different parameter scales (Chameleon setting and Chameleon+Speech setting).

Model Scale	Metric Category	Metric Name	Mamba Loss (↓)	Mixture-of- Mamba Loss (↓)	Performance Gain (%) (↑)	Relative Training FLOPs to Match Mamba (%) (↓)
		Training Loss	6.1308	5.9561	2.85%	25.90%
	Image Metrics	Obelisc Val. Loss	5.2866	5.1124	3.29%	26.60%
37M ⁻		SSTK Val. Loss	6.6694	6.5023	2.51%	27.50%
		Training Loss	3.6262	3.5175	3.00%	60.90%
	Text Metrics	Obelisc Val. Loss	4.1244	4.0469	1.88%	64.80%
		SSTK Val. Loss	4.0417	3.9533	2.19%	57.50%
	Overall	Avg Training Loss	4.6607	4.5247	2.92%	50.70%
		Training Loss	5.7609	5.6057	2.69%	35.70%
	Image Metrics	Obelisc Val. Loss	4.9231	4.7683	3.14%	35.30%
		SSTK Val. Loss	6.3130	6.1652	2.34%	37.00%
94M		Training Loss	3.0294	2.9414	2.90%	58.40%
	Text Metrics	Obelisc Val. Loss	3.6016	3.5270	2.07%	62.60%
		SSTK Val. Loss	3.4109	3.2901	3.54%	61.40%
	Overall	Avg Training Loss	4.1577	4.0419	2.78%	49.80%
		Training Loss	5.3558	5.1703	3.46%	33.40%
	Image Metrics	Obelisc Val. Loss	4.5258	4.3546	3.78%	35.10%
		SSTK Val. Loss	5.9179	5.7471	2.89%	35.30%
443M		Training Loss	2.4637	2.3864	3.14%	62.00%
	Text Metrics	Obelisc Val. Loss	3.0544	2.9820	2.37%	66.70%
		SSTK Val. Loss	2.7569	2.6250	4.78%	54.70%
	Overall	Avg Training Loss	3.6584	3.5364	3.33%	47.90%
		Training Loss	5.2260	5.1201	2.03%	48.40%
	Image Metrics	Obelisc Val. Loss	4.4127	4.3105	2.32%	49.30%
		SSTK Val. Loss	5.7987	5.6986	1.73%	50.50%
880M		Training Loss	2.3073	2.2438	2.75%	65.60%
	Text Metrics	Obelisc Val. Loss	2.8886	2.8313	1.99%	72.80%
		SSTK Val. Loss	2.5483	2.4548	3.67%	67.90%
	Overall	Avg Training Loss	3.5130	3.4320	2.31%	58.30%
		Training Loss	5.1892	5.0591	2.51%	42.50%
	Image Metrics	Obelisc Val. Loss	4.3692	4.2510	2.71%	44.50%
1.5B ⁻		SSTK Val. Loss	5.7546	5.6335	2.10%	44.60%
		Training Loss	2.2284	2.1614	3.01%	65.40%
	Text Metrics	Obelisc Val. Loss	2.8020	2.7393	2.24%	71.60%
		SSTK Val. Loss	2.4614	2.3455	4.71%	62.10%
	Overall	Avg Training Loss	3.4602	3.3670	2.69%	54.70%

Table 7. Training and validation metrics across model scales in the Chameleon setting. In this setting, both image and text modalities are represented as discrete tokens. Mixture-of-Mamba achieves substantial performance improvements over Mamba Dense, with the image modality showing the largest gains across all five model scales: 37M, 94M, 443M, 880M, and 1.5B. Notably, the text modality also exhibits significant improvements, in contrast to the Transfusion setting where text gains were more modest. These results further highlight the effectiveness and efficiency of Mixture-of-Mamba, which consistently achieves strong performance with reduced relative training FLOPs.



Figure 4. Validation loss and loss matching for text modality across model scales (C4 and Wikipedia datasets) during multi-modal pretraining in the Transfusion setting. Results are shown for Mixture-of-Mamba, Mamba Dense, and Flex-Attention Transformer at three model scales: 163M, 760M, and 1.4B. (a, e, i) Validation loss on the C4 dataset shows that Mixture-of-Mamba achieves comparable performance at 163M and performs marginally better than Mamba Dense and Flex-Attention Transformer at the 760M and 1.4B scales. (b, f, j) Loss matching for C4 demonstrates that Mixture-of-Mamba reaches similar or slightly lower loss values at earlier training steps compared to Mamba Dense. (c, g, k) Validation loss on the Wikipedia dataset follows a similar trend, with Mixture-of-Mamba showing marginal improvements at the 760M and 1.4B scales. (d, h, l) Loss matching for Wikipedia illustrates efficient training dynamics, with Mixture-of-Mamba aligning closely with Flex-Attention Transformer while reaching comparable or slightly lower loss values than Mamba Dense. Overall, Mixture-of-Mamba demonstrates moderate improvements over both baselines at the larger scales (760M and 1.4B).



Figure 5. Image validation loss and loss matching on the CC12M dataset across three model scales: 163M, 760M, and 1.4B during
 multi-modal pretraining in the Transfusion setting. (a, c, e) Validation loss curves show that Mixture-of-Mamba achieves substantially
 lower image validation loss compared to Mamba Dense and Flex-Attention Transformer across all scales, with the improvement becoming
 more pronounced as model size increases. (b, d, f) Loss matching curves demonstrate that Mixture-of-Mamba reaches the same loss
 values at earlier training steps compared to Mamba Dense, highlighting improved training efficiency. Overall, Mixture-of-Mamba achieves
 large improvements in image validation loss on the CC12M dataset, showcasing its effectiveness in the image modality.



across multi-modal pretraining.



modality.

Figure 7. Modality-specific pre-training loss and step matching plots across model scales (Chameleon setting). Training loss and loss matching are reported for image and text modalities across five model scales: 37M, 94M, 443M, 880M, and 1.5B. (a, e, i, m, q) Image

training loss shows significant improvements for Mixture-of-Mamba (orange), which consistently achieves lower loss compared to Mamba

Dense (cvan) across all scales. (b, f, j, n, r) Image loss matching compares the training dynamics and shows that Mixture-of-Mamba

reaches the same loss values at earlier training steps compared to Mamba Dense, highlighting its improved efficiency. (c, g, k, o, s) Text

training loss demonstrates competitive performance, with Mixture-of-Mamba achieving slightly lower loss values compared to Mamba

Dense. (d, h, l, p, t) Text loss matching illustrates that Mixture-of-Mamba reaches the same loss values at earlier training steps compared

to Mamba Dense, reflecting its efficient training dynamics. Overall, in the Chameleon setting, Mixture-of-Mamba achieves consistent

improvements in the image modality, with substantial computational savings, while also demonstrating meaningful gains in the text



Figure 8. Training and evaluation losses for image and text modalities across model scales in the Chameleon setting on the Obelisc dataset. Results are shown for Mixture-of-Mamba and Mamba Dense across five model scales: 37M, 94M, 443M, 880M, and 1.5B. (a, e, i, m, q) Image evaluation loss demonstrates consistent improvements for Mixture-of-Mamba (orange), achieving lower loss compared to Mamba Dense (cyan) across all scales. (b, f, j, n, r) Image loss matching shows that Mixture-of-Mamba reaches the same loss values at earlier training steps compared to Mamba Dense, reflecting its improved training efficiency. (c, g, k, o, s) Text evaluation loss indicates competitive results for Mixture-of-Mamba, achieving lower losses relative to Mamba Dense. (d, h, l, p, t) Text loss matching highlights that Mixture-of-Mamba reaches the same loss values at earlier training steps, further demonstrating its efficiency in the text modality.
Overall, Mixture-of-Mamba achieves strong and consistent improvements in both image and text modalities across all model scales in the Chameleon setting evaluated on the Obelisc dataset.



Figure 9. Training and evaluation losses for image and text modalities across model scales in the Chameleon setting on the 1087 Shutterstock dataset. Results are shown for Mixture-of-Mamba and Mamba Dense across five model scales: 37M, 94M, 443M, 880M, 1088 and 1.5B. (a, e, i, m, q) Image evaluation loss demonstrates consistent improvements for Mixture-of-Mamba (orange), achieving lower 1089 loss compared to Mamba Dense (cyan) across all scales. (b, f, j, n, r) Image loss matching shows that Mixture-of-Mamba reaches the 1090 same loss values at earlier training steps compared to Mamba Dense, reflecting its improved training efficiency. (c, g, k, o, s) Text 1091 evaluation loss indicates competitive results for Mixture-of-Mamba, achieving lower losses relative to Mamba Dense. (d, h, l, p, t) Text 1092 loss matching highlights that Mixture-of-Mamba reaches the same loss values at earlier training steps, further demonstrating its efficiency in the text modality. Overall, Mixture-of-Mamba achieves strong and consistent improvements in both image and text modalities across 1093 all model scales in the Chameleon setting evaluated on the Shutterstock dataset. 1094

- 1095
- 1096
- 1097
- 1098
- 1099



1147 Mixture-of-Mamba consistently reduces average training loss and achieves more efficient convergence across all model scales in the

- 1148 Chameleon setting.



Figure 11. Training and evaluation losses for image, text, and speech modalities (37M and 94M scales) in the Chameleon+Speech setting. Results are reported for Mixture-of-Mamba and Mamba Dense. (a, e, i) Image training loss demonstrates that Mixture-of-Mamba (orange) achieves consistently lower loss compared to Mamba Dense (cyan). ($\mathbf{b}, \mathbf{f}, \mathbf{j}$) Image loss matching highlights Mixture-of-Mamba's ability to reach the same loss values at earlier training steps, showing improved training efficiency. (c, g, k) Text training loss shows competitive results for Mixture-of-Mamba, improving over Mamba Dense. (d, h, l) Text loss matching confirms Mixture-of-Mamba's ability to reach the same loss values at earlier training steps, showing improved training efficiency. (e, m) Speech training loss highlights significant improvements in speech modality performance. (f, n) Speech loss matching shows efficient learning dynamics for Mixture-of-Mamba. (g, o) Speech evaluation loss on LL60K confirms notable performance gains, and (h, p) Speech evaluation loss on PPL30K further highlights the efficiency of Mixture-of-Mamba.



1255 Figure 12. Training and evaluation losses for image, text, and speech modalities (443M, 880M, and 1.5B scales) in the 1256 Chameleon+Speech setting. Results are reported for Mixture-of-Mamba and Mamba Dense. (a, i, q) Image training loss demonstrates that Mixture-of-Mamba (orange) consistently outperforms Mamba Dense (cyan) across larger scales. (b, j, r) Image loss matching highlights improved training efficiency for Mixture-of-Mamba, reaching the same loss values at earlier training steps. (c, k, s) Text 1258 training loss shows Mixture-of-Mamba achieving better performance. $(\mathbf{d}, \mathbf{l}, \mathbf{t})$ Text loss matching further demonstrates efficient learning 1259 dynamics. (e, m, u) Speech training loss confirms substantial gains for Mixture-of-Mamba in the speech modality, consistent across 1260 model scales. (f, n, v) Speech loss matching illustrates the improved efficiency of Mixture-of-Mamba across scales. (g, o, w) Speech evaluation loss on LL60K highlights consistent improvements, while (h, p, x) Speech evaluation loss on PPL30K demonstrates notable 1262 gains and efficient performance across scales. 1263



Figure 13. Training and validation losses for image and text modalities across model scales in the Chameleon+Speech setting evaluated on the Obelisc dataset. Results are shown for Mixture-of-Mamba and Mamba Dense across five model scales: 37M, 94M, 443M, 880M, and 1.5B. (a, e, i, m, q) Image evaluation loss demonstrates consistent gains for Mixture-of-Mamba (orange) over Mamba Dense (cyan), even with the inclusion of the speech modality. (b, f, j, n, r) Image loss matching shows that Mixture-of-Mamba reaches the same loss values at earlier training steps compared to Mamba Dense, highlighting improved efficiency. (c, g, k, o, s) Text evaluation loss indicates consistent reductions for Mixture-of-Mamba relative to Mamba Dense across all scales. (d, h, l, p, t) Text loss matching illustrates that Mixture-of-Mamba reaches the same loss values at earlier training steps compared to Mamba Dense, maintaining its efficiency in the text modality. Overall, Mixture-of-Mamba achieves consistent improvements in both image and text modalities while maintaining its efficiency, even with the addition of the speech modality. These results confirm the robustness of Mixture-of-Mamba in multi-modal settings.



Figure 14. Training and validation losses for image and text modalities across model scales in the Chameleon+Speech setting evaluated on the Shutterstock dataset. Results are shown for Mixture-of-Mamba and Mamba Dense across five model scales: 37M, 94M, 443M, 880M, and 1.5B. (a, e, i, m, q) Image evaluation loss demonstrates consistent gains for Mixture-of-Mamba (orange) over Mamba Dense (cyan), even with the inclusion of the speech modality. (b, f, j, n, r) Image loss matching shows that Mixture-of-Mamba reaches the same loss values at earlier training steps compared to Mamba Dense, highlighting improved efficiency. (c, g, k, o, s) Text evaluation loss indicates consistent reductions for Mixture-of-Mamba relative to Mamba Dense across all scales. (d, h, l, p, t) Text loss matching illustrates that Mixture-of-Mamba reaches the same loss values at earlier training steps compared to Mamba Dense, maintaining its efficiency in the text modality. Overall, Mixture-of-Mamba achieves consistent improvements in both image and text modalities while maintaining its efficiency, even with the addition of the speech modality. These results confirm the robustness of Mixture-of-Mamba in multi-modal settings.