# OuroMamba: A Data-Free Quantization Framework for Vision Mamba Models

**Akshat Ramachandran**[g]†, **Mingyu Lee**[g]†, **Huan Xu**[g], **Souvik Kundu**[i], **Tushar Krishna**[g]

[g]Georgia Institute of Technology, Atlanta, USA
[i]Intel Labs, USA
†Equal Contribution Authors
[g]{akshat.r, mlee864, hxu398}@gatech.edu, tushar@ece.gatech.edu
[i]souvikk.kundu@intel.com

## Abstract

We present **OuroMamba**, the first data-free post-training quantization (DFQ) method for vision Mamba-based models (VMMs). We identify two key challenges in enabling DFQ for VMMs, (1) VMM's recurrent state transitions restricts capturing of long-range interactions and leads to semantically weak synthetic data, (2) VMM activations exhibit dynamic outlier variations across time-steps, rendering existing static PTQ techniques ineffective. To address these challenges, OuroMamba presents a two-stage framework: (1) OuroMamba-`Gen` to generate semantically rich and meaningful synthetic data. It applies contrastive learning on patch level VMM features generated through neighborhood interactions in the latent state space, (2) OuroMamba-`Quant` to employ mixed-precision quantization with lightweight dynamic outlier detection during inference. In specific, we present a thresholding based outlier channel selection strategy for activations that gets updated every time-step. Extensive experiments across vision and generative tasks show that our **data-free** OuroMamba surpasses existing data-driven PTQ techniques, achieving state-of-the-art performance across diverse quantization settings. Additionally, we implement efficient GPU kernels to achieve practical latency speedup of up to $\mathbf{2.36\times}$. Code will be released soon.
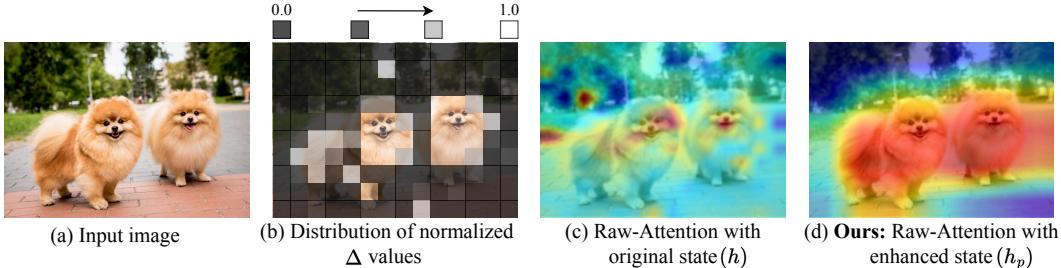
Figure 1: Qualitative comparison of VMM implicit attention. (a) Sample input image. (b) Normalized distribution of the input gate values ($\Delta$). Visualization of implicit attention Ali et al. (2024) using (c) the original state ($h$) and (d) the **proposed** patched state ($h_p$), which incorporates spatial dependencies through patched neighborhood interactions.

## 1 Introduction

The sub-quadratic compute alternative of the vision State space models (SSMs) Gu & Dao (2023), particularly vision mamba models (VMMs) Zhu et al. (2024) has made them a promising alternative to ViTs Dosovitskiy et al. (2020). Like ViTs, larger VMM variants face deployment challenges due to high memory and latency constraints Li et al. (2025b). Quantization Ramachandran et al. (2024a;b;c) is a popular method to tackle the high memory and latency demand by mapping full-precision (FP) weights and/or activations to lower-bit representations. In particular, post-training quantization (PTQ) Ramachandran et al. (2024a); Kundu et al. (2022) converts a pre-trained FP model to low-precision and requires calibration on a small dataset Ramachandran et al. (2024a) to

recover model performance. However, PTQ's calibration data is typically drawn from the original training dataset, potentially restricting its applicability in situations demanding privacy, security of training set Kundu et al. (2021); Zhang et al. (2023).

Consequently, data-free quantization (DFQ), a subset of post-training quantization (PTQ) has emerged as a promising alternative. DFQ allows calibration via generating synthetic data from Gaussian noise Ramachandran et al. (2024a); Kim et al. (2025); Li et al. (2023b; 2022) mimicking the distribution of the original train-set. In particular, DFQ for ViTs primarily rely on the distinctive response of self-attention to noise and real data to generate synthetic data Li et al. (2023b); Ramachandran et al. (2024a). However, despite their significant progress for ViTs, *DFQ techniques for VMM remains largely unexplored*. Interestingly, VMM's recurrent state transition in the S6 layer Gu & Dao (2023) lacks an explicit self-attention mechanism, limiting the usage of any attention based patch similarity improvement. While the S6 layer can be formulated as a variant of linear attention Han et al. (2025), our observations reveal that this implicit attention struggles to distinguish foreground from the background of an image. As observed in Figure 1(c), Figure 4(a), despite VMM's "theoretical ability" for global interactions Azizi et al. (2025), it fails to capture long-range dependencies in feature representation.

Notably, VMM's recurrent state transitions introduce dynamic variations in activation characteristics Li et al. (2025b). These variations occur across time steps and tokens[1], potentially demanding online adaptive outlier selection. This is in sharp contrast to ViTs, which largely demonstrate a static outlier activation pattern Dong et al. (2023). A contemporary work, QMamba Li et al. (2025b), rightly identified the dynamic variations in SSM activations. However, it employs statically determined temporal grouped quantization, unable to adapt to the potential changes in outlier locations, yielding large quantization error at low bit precision. Moreover, QMamba lacks a data-free calibration strategy and limited practical speedup, limiting its use in privacy and latency-sensitive scenarios.

**Our Contributions.** To address the above challenges in data generation and quantization stages, we present **OuroMamba**[2], a *first-of-its-kind* method to enable DFQ for VMMs. For data generation, we investigate on VMM's implicit attention struggles to differentiate foreground with background. In specific, we hypothesize it to be dependent on scanning direction of the compressed hidden state $(h(t))$, which limits explicit spatial token interactions. To mitigate this, we propose **OuroMamba**-Gen, that essentially enhances the implicit attention through patched spatial interactions in the latent state space, resulting in a refined hidden state namely, patched state $h_p(t)$. We then leverage an enhanced implicit attention representation based on $h_p(t)$ (Figure 1(d)) to generate synthetic data by employing patch-level contrastive learning Ramachandran et al. (2024a).

For quantization, we first empirically validate that VMM operation leads to dynamic inter-time-step variations of the outlier channel positions. We then present **OuroMamba**-Quant, a mixed-precision quantization scheme that performs channel-wise static weight quantization and dynamic activation quantization per-time-step during inference. OuroMamba-Quant reduces the quantization error associated to dynamic outliers via three key strategies, (1) a lightweight channel-wise activation outlier detection, (2) mixed-precision quantization, where outlier channels are quantized at higher precision, while remaining inlier channels use lower bit quantization and (3) efficient outlier management using an adaptive outlier list.

To validate the efficacy of OuroMamba, we conduct extensive experimental evaluations on VMMs Zhu et al. (2024); Liu et al. (2025); Huang et al. (2024) and hybrid Transformer-Mamba models Hatamizadeh & Kautz (2025) for classification, detection, segmentation, and generative modeling tasks Hu et al. (2024). OuroMamba achieves SoTA performance over existing data-driven VMM PTQ alternatives, with an accuracy improvement of up to 39%. Additionally, through efficient kernel implementation we demonstrate practical end-to-end latency speedup of up to **2.36× over FP16 baseline**.

## 2 RELATED WORKS

**Data-Driven PTQ.** Existing PTQ techniques for ViTs Lin et al. (2021); Jiang et al. (2025); Li et al. (2023c); Yang et al. (2024); Ma et al. (2024); Yuan et al. (2022) address long-tailed distributions and

---

[1]We use token and patch interchangeably based on context.

[2]Inspired by **Ouroboros**—the self-consuming snake— the name OuroMamba reflects its ability to self-generate data for VMM quantization.

static activation outliers but fail to handle VMMs' dynamically changing outliers Cho et al. (2025); Li et al. (2025b). While PTQ4VM Cho et al. (2025) adapts SmoothQuant Xiao et al. (2023) to shift activation outliers into weights, it does not quantize SSM activations, limiting its effectiveness. QMamba Li et al. (2025b) employs temporal grouped quantization but relies on static groupings, making it unsuitable for ultra-low bit-widths. *OuroMamba is the first to identify dynamic outlier variations in VMMs and introduce a mixed-precision quantization framework that achieves SoTA accuracy at low precision.*

**Date-Free PTQ.** Existing DFQ techniques for ViTs, PSAQ-ViT v1 Li et al. (2022) and v2 Li et al. (2023b) optimize noise into synthetic data by maximizing global patch similarity entropy, but overlook spatial sensitivity and semantic inter-patch relationships, leading to simplistic data. CLAMP-ViT Ramachandran et al. (2024a) addresses this by introducing patch-level contrastive learning on self-attention, aligning foreground patches while separating background patches. Following data generation, these methods apply uniform symmetric quantization to weights and activations, which is suboptimal for VMM. *DFQ for VMM thus far remains unexplored, and this work makes the first effort in this direction.*

## 3 PRELIMINARIES

### 3.1 MAMBA BLOCK IN VMM

**The S6 Layer.** Each Mamba block in VMM utilizes the selective SSM (S6) mechanism Gu & Dao (2023) that maps an input signal $u(t) \in \mathbb{R}^E$ to an intermediate hidden state $h(t) \in \mathbb{R}^{E \times N}$ through a **linear recurrent transition**. Here, $E, N, t$ correspond to the number of channels, state dimension and time-step, respectively. $h(t)$ then produces the output $o(t) \in \mathbb{R}^E$ via a linear transformation, as shown below.

$$h(t) = \bar{A}(t) \odot h(t-1) + \bar{B}(t) \odot u(t); \ o(t) = C(t)h(t) \tag{1}$$

where $\odot$ is element-wise product, $\bar{A}(t), \bar{B}(t) \in \mathbb{R}^{E \times N}$ and $C(t) \in \mathbb{R}^{N \times 1}$ are the discrete time-variant system, input, and output matrices, respectively. The S6 layer generates these per-time-step discrete matrices from the input sequence and their corresponding continuous parameter counterparts as follows,

$$\bar{A}(t) = e^{(A\Delta(t))}; B(t) = W_B(u(t)); \bar{B}(t) = B\Delta(t)$$
$$\Delta(t) = S^+(u(t)\Delta(t)_{\text{proj}}); C(t) = (W_C(z(t)))^T \tag{2}$$

where, $S^+$ is softplus, $\Delta(t)_{\text{proj}}, W_B$, and $W_C$ are linear projection layers. $\Delta(t) \in \mathbb{R}^{E \times 1}$ is the discretization tensor. $\bar{A}, \bar{B}, C, \Delta$ are input-dependent discrete S6 parameters that exhibit dynamic activation variations at each time step. Equation 1 deals with scalar inputs. To operate over an input sequence $u$ of batch size $B$, sequence length $M$ with $E$ channels, $u \in \mathbb{R}^{B \times M \times E}$, Equation 2, Equation 1 are applied independently to each $u(t) \in \mathbb{R}^{B \times E}$ Gu & Dao (2023) across $M$ time-steps.

**Mamba Block.** The VMM model is composed by stacking multiple Mamba blocks that maps its input sequence of tokens $X$ to its output sequence $Y$ as follows,

$$G = \sigma(W_{gate\_proj}X); U = Conv1D(W_{\text{in\_proj}}X)$$
$$O = S6(U); \ Y = O \odot G \tag{3}$$

here, $G$ is a gating function obtained from a linear transformation of X, followed by a SiLU activation $\sigma$. $\odot$ between $G$ and $O$ enables the model to selectively emphasize or suppress input features. The S6 input, $U$, is a linearly transformed version of $X$, followed by a 1D convolution.

**Selective Scan for Vision.** While selective scan in S6 is well-suited for 1D NLP tasks, vision data is 2D, requiring VMMs to adopt specialized scanning strategies to flatten images into 1D sequences before processing them through Mamba's selective scan. Examples include 2-way scanning Zhu et al. (2024) (Figure 2(a)), 4-way scanning Liu et al. (2025). However, scan direction choice can

introduce directional dependency, impacting the model's ability to capture spatial relationships Xiao et al. (2024).

## 3.2 IMPLICIT ATTENTION IN VMM

Prior works Han et al. (2025); Ali et al. (2024) have shown that S6 layer operation can be reformulated into an implicit self-attention mechanism. By setting $h(0) = 0$ and unrolling Equation 1, the hidden state can be derived as, $h(t) = \sum_{j=1}^{t}(\prod_{k=j+1}^{t} \bar{A}(k))\bar{B}(j)u(j)$ and $o(t)$ can be computed by using this derived value of $h(t)$ in Equation 1. Additionally, in computing $o(t)$, the contribution of the $j^{th}$ token on the $i^{th}$ token is captured as follows,

$$\tilde{\alpha}[i,j] = \sum_{m=1}^{N} \alpha^m[i,j] = \sum_{m=1}^{N} C(i)(\prod_{k=j+1}^{i} \bar{A}(k)_m)\bar{B}(j)[m] \tag{4}$$

where, $\bar{A}(k)_m$ is the $m^{\text{th}}$ diagonal element. The output derivation can be simplified as $o = \tilde{\alpha}u$. Thus, $\tilde{\alpha} \in \mathbb{R}$ is the implicit attention score matrix and $\alpha \in \mathbb{R}^N$ is the $N$-dimensional implicit attention matrix of VMM.

## 3.3 QUANTIZATION

In this paper, we perform symmetric group quantization Ramachandran et al. (2024a); Li et al. (2022) of both weights and activations for VMMs. The quantization process for an input FP tensor $X$ with a target bit-precision $b$ is given by,

$$Q(X, S, b) = clip(\lfloor \frac{X}{S} \rceil, -2^{b-1} + 1, 2^{b-1} - 1) \tag{5}$$

where, $S = \max(|X|)$ is the scale factor.

## 3.4 CONTRASTIVE OBJECTIVE

Contrastive learning based on the infoNCE loss Wang et al. (2023) helps learn an anchor patch from both similar (positive) and dissimilar (negative) patches. Following Ramachandran et al. (2024a), the formulation of patch-level contrastive loss ($\mathcal{L}_{l,(i,j)}^{C}$) we employ for data generation is as follows,

$$\mathcal{L}_{l,(i,j)}^{C} = -\log \frac{\sum_{p+} \exp(\lambda^* \cdot \lambda^+/\tau)}{\sum_{p+} \exp(\lambda^* \cdot \lambda^+/\tau) + \sum_{p-} \exp(\lambda^* \cdot \lambda^-/\tau)} \tag{6}$$

where, $*$, $+$, and $-$ correspond to the anchor patch, positive, and negative patches, respectively. $\tau$ controls the concentration level Wang et al. (2023). $\lambda_{l,(i,j)}$ (shown without subscript above for brevity) represents the embedding employed to calculate $\mathcal{L}_{l,(i,j)}^{C}$. $l$ is the layer and $(i,j)$ is the location of anchor patch. The final loss is given as $\mathcal{L}^C = \sum_l \sum_i \sum_j \mathcal{L}_{l,(i,j)}^{C}$.

# 4 MOTIVATIONAL ANALYSIS

## 4.1 SYNTHETIC DATA GENERATION

> **Observation 1:** *Synthetic data generated for ViT fails to transfer effectively to VMM.*

We investigate the W4A8 quantization performance of Vim-S Zhu et al. (2024) on an existing VMM PTQ framework, QMamba Li et al. (2025b), using different calibration data sources. In specific, we perform calibration with both real data and ViT generated synthetic data of same sample size of 1024 images. We use synthetic calibration data produced by ViT DFQ technique, namely, CLAMP-ViT Ramachandran et al. (2024a). Figure 2 reveals that across tasks, ViT synthetic data results in significant accuracy degradation (7-14%) as compared to using real calibration data.
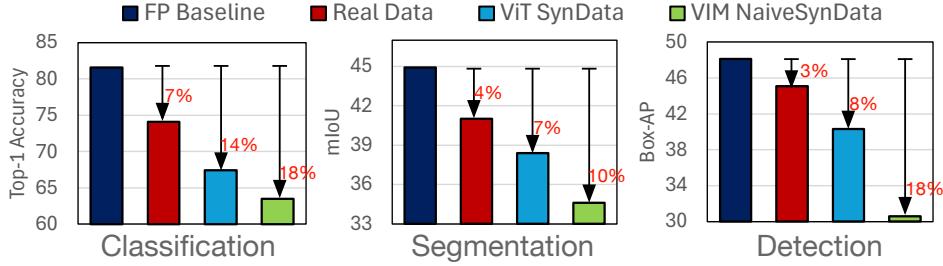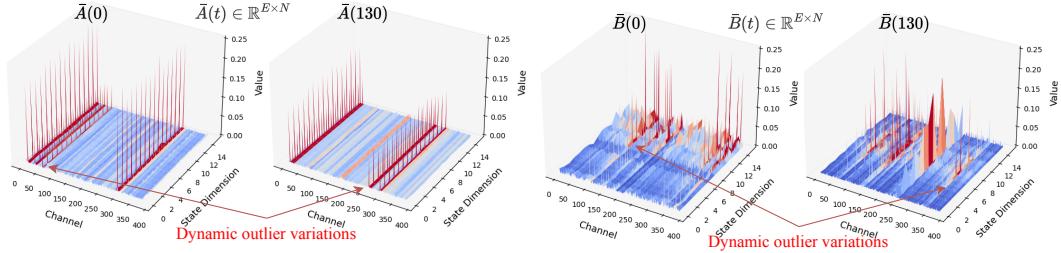
Figure 2: W4A8 quantization performance comparison with different calibration data sources.



Figure 3: Dynamic inter-time-step outlier channel variations for two representative S6 layer activations: $\bar{A}, \bar{B}$ in layer 3 of Vim-T.

**Observation 2:** *VMM's naive implicit attention reformulation is suboptimal for generating synthetic data.*

As shown in Figure 1(c), VMM's implicit attention struggles to distinguish foreground from background. This is potentially due to the forgetfulness of the states Azizi et al. (2025) as we traverse along the scanning direction (Figure 4(a)) Han et al. (2025). This, unlike self-attention, limits long-range global interactions. Additionally, VMMs Zhu et al. (2024); Liu et al. (2025) employ multiple scanning directions (subsection 3.1) to improve global interactions. However, they process each direction independently with separate S6 blocks. Consequently, each scanning direction's implicit attention exhibits directional bias and are incoherent Han et al. (2025).

Following Ramachandran et al. (2024a), that applies patch-level contrastive loss to ViT self-attention to generate synthetic data, we apply the same method to VMM's forward scan S6 block and present the generated data in Figure 4(b). The "naive synthetic data" appears noisy, lacks semantic structure, and performs worse than ViT-generated synthetic data for calibration, Figure 2. Extending this to other scanning directions yields similar results (see §6). *These observations motivate an optimal VMM data generation technique, capturing spatial dependencies without scanning direction constraints.*



Figure 4: (a) Forward and backward SSM state transition of Vim-S, (b) Naive synthetic data samples generated by applying Ramachandran et al. (2024a) on VMM implicit attention.

## 4.2 QUANTIZATION

**Observation 3:** *VMM activations exhibit dynamic inter-time-step channel variations.*

It has been well established by prior work Dong et al. (2023); Jiang et al. (2025); Yang et al. (2024) that outlier channels in ViT activations follow a fixed pattern, allowing anticipation in advance dur-

ing calibration. However, unlike ViTs, **VMMs exhibit dynamic variations in activation outlier channels** across time-steps. As shown in Figure 3, a comparison of activations $\bar{A}, \bar{B}$ at time-steps 0 and 130 reveals that outlier channels (in red) differ significantly between time-steps. This temporal variance in outlier positions suggests that a static, pre-determined calibration approach, as used in ViTs, is insufficient for VMM quantization. *Therefore, VMM necessitates dynamic identification of outlier channels at inference to reduce quantization error.*

## 5 OUROMAMBA DFQ FRAMEWORK

This section presents OuroMamba, a unified DFQ framework with OuroMamba-`Gen` for data generation and OuroMamba-`Quant` for quantization.

---

**Algorithm 1:** OuroMamba-`Gen`

---

**Input:** A pre-trained FP VMM model $P$ with $L$ layers, Gaussian noise batch $X_\mathcal{B}$, task-specific
    targets $T_{G_\mathcal{B}}$, neighborhood size $\mathcal{N}$, iterations $G$.
**Output:** A set of generated synthetic samples $X_\mathcal{B}^*$.
**for** $g = 1, 2, \ldots, G$ **do**
    Input $X_\mathcal{B}$ into $P$ ;
    **for** $l = 1, 2, \ldots, L$ **do**
        Capture per-time-step original $h^l(t), \Delta^l(t)$ ;
        Set $w_t^l = \texttt{mean}_E(\Delta^l(t))$ ;
        Compute $h_p^l(t)$;
        Extract implicit attention;
        Compute $\mathcal{L}_l^C = \sum_t \mathcal{L}_{l,t}^C$;
    **end**
    Compute $\mathcal{L}^C = \sum_l \mathcal{L}_l^C$;
    Compute output loss $\mathcal{L}^O$;
    Compute $L^{gen} = \mathcal{L}^C + \mathcal{L}^O$;
    Update $X_\mathcal{B}^*$ via backpropagation of $L^{gen}$;
**end**

---

### 5.1 STAGE 1: OUROMAMBA-GEN

Following the typical DFQ setup Ramachandran et al. (2024a); Li et al. (2023b), OuroMamba-`Gen` requires an input batch of $\mathcal{B}$ random Gaussian noise images $X_\mathcal{B}$, and corresponding randomly generated task-specific targets $T_{G_\mathcal{B}}$ (detailed in §6). $X_\mathcal{B}$ is fed to the VMM to minimize the generation loss $L^{gen}$ (described below), updating $X_\mathcal{B}$ via backpropagation for $G$ iterations. Upon convergence, the updated $X_\mathcal{B}^*$ is the generated synthetic data employed for calibration. The algorithm of OuroMamba-`Gen` pipeline is given in Alg. 1.

Building on the observations in §4.1, OuroMamba-`Gen` first constructs an implicit attention representation that possesses enhanced global, spatial interactions for improved synthetic data generation.

**Patched hidden state** $(h_p(t))$**.** Once the original latent space for a layer is determined, for each time-step of the hidden state $h(t)$ in the original latent space, we define a neighborhood patch $\mathcal{N}(t)$ of size $p \times p$. $\mathcal{N}(t)$ covers all the spatially adjacent tokens $h(k), k \in \mathcal{N}(t)$ in the corresponding 2D representation of the hidden states, where the cardinality of $\mathcal{N}(t)$ denotes the number of neighboring states. For instance, for the $\tau^{th}$ state with $p = 3$, yields total $|\mathcal{N}(\tau)| = 9$ states, with $h(\tau)$ at the center of the patch. For $\tau^{th}$ state, we compute its corresponding *patched state* $(h_p(\tau))$ via a weighted sum of the states in $\mathcal{N}(\tau)$ as,

$$h_p(\tau) = \sum_{k \in \mathcal{N}(\tau)} w_k h(k) \tag{7}$$

where, $w_k$ is a weighting factor that modulates the contribution of each neighboring state. Through our experiments and insights from prior work Han et al. (2025), we observe that $\Delta(t)$ exhibits higher

magnitude responses in informative regions, such as foreground, while suppressing less relevant regions (Figure 1(b)). This property makes $\Delta(t)$ a natural choice for weighting the linear aggregation in Equation 7, enabling adaptive feature aggregation. To leverage this, we determine the weighting factor for all states in $\mathcal{N}(\tau)$ by performing a mean reduction along the channel dimension ($E$) of each $\Delta(\tau)$, as given by, $w_k = mean_E(\Delta(k))$ for $k \in \mathcal{N}(\tau)$.

**Enhanced Implicit Attention.** By substituting $h(t)$ with $h_p(t)$ in §3.2, we obtain an enhanced representation of the N-dimensional implicit attention matrix ($\alpha$), denoted as $\alpha_p$. Interestingly, empirical validations show its ability to effectively separate the foreground from background (refer to Figure 1(d)), making it a potentially more informative representation compared to $\alpha$.

**Generation Loss ($\mathcal{L}^{gen}$).** Ramachandran et al. (2024a) applied patch-level contrastive loss to self-attention in generating synthetic data for ViT. Inspired by this, we apply the contrastive loss (Equation 6) on $\alpha_p$. We specifically choose the patch-level contrastive loss as it has been proven by prior work to generate high quality synthetic data through capturing semantic relations. For a VMM layer $l$, each $\alpha_p^l[i, j]$ serves as the anchor patch i.e, $\lambda_{l,(i,j)}^* = \alpha_p^l[i, j]$ with positive and negative patches selected from spatially adjacent $\alpha_p^l[x, y], (x, y) \in \mathcal{N}(i, j)$. Similar to Ramachandran et al. (2024a), we leverage the cosine similarity metric to identify positive and negative patches in $\mathcal{N}(i, j)$, where we select the top-$n$ patches in $\mathcal{N}(i, j)$ with highest similarity as positive and rest as negative, where $n = \lfloor|\mathcal{N}(i, j)|/2\rfloor$. We then compute $\mathcal{L}_{l,(i,j)}^C$ for all anchor patches in $\alpha_p^l$ across all layers to get the final contrastive learning objective $\mathcal{L}^C$.

Furthermore, to direct the synthetic data generation process towards task-specific goals Li et al. (2023b); Ramachandran et al. (2024a), we employ an additional output loss $\mathcal{L}^O$, which is the mean absolute error (MAE) between the predicted model output with the task specific targets. The final data generation loss is thus $L^{gen} = \mathcal{L}^C + \mathcal{L}^O$.

---

**Algorithm 2:** OuroMamba-`Quant`

---

**Input** : Activation $X(t) \in \mathbb{R}^{N \times E}$, Static scale $S^I(t)$, Threshold $\theta$, Refresh rate $n_{\texttt{refresh}}$,
Outlier list $O_{\texttt{list}}$, Inlier and outlier bit-precision $b_a^I, b_a^O$
**Output:** Quantized activation $X_q(t)$, Updated outlier list $O_{\texttt{list}}$
**if** $t \% n_{refresh} == 0$ **then**
| $\quad O_{\texttt{list}} = \{\phi\}$
**end**
$S^D(t) = \texttt{ComputeScale}(X(t)[:, c] \forall c \notin O_{\text{list}})$
**if** $S^D(t) > S^I(t)$ **then**
| $\quad$ **for** *each channel c in $X(t)$ **not in** $O_{list}$* **do**
| $\quad\quad$ **if** $max(|X(t)[:, c]|) \geq \theta$ **then**
| $\quad\quad\quad$ | $O_{\texttt{list}} = O_{\texttt{list}} \cup \{c\}$
| $\quad\quad$ **end**
| $\quad$ **end**
**end**
$I(t), O(t) = \texttt{Separate}(X(t), O_{\texttt{list}})$
$I_q(t) = \texttt{InlierQuant}(I(t), S^I(t), b_a^I)$
$O_q(t) = \texttt{OutlierQuant}(O(t), b_a^O)$
$X_q(t) = \texttt{Merge}(I_q(t), O_q(t))$
**return** $X_q(t), O_{\texttt{list}}$

---

## 5.2 STAGE 2: OUROMAMBA-QUANT

Building on the insights from §4.2, OuroMamba-`Quant` is designed to accommodate the dynamic nature of VMMs. OuroMamba-`Quant`'s algorithm for activation quantization per-time-step is given in Alg. 2.

**Offline Calibration.** We use synthetic data generated by OuroMamba-`Gen` for calibration. **During calibration**, for an input activation tensor $X(t)$, we identify the per-time-step inlier scale factor $S^I(t)$ shared over each tensor of size $N \times E$. We then compute the magnitude-based threshold $\theta$ that distinguishes outliers by identifying values exceeding this threshold. While outliers exhibit

Figure 5: Synthetic data samples generated by OuroMamba.

inter-channel variations across time steps, their magnitude and dynamic range remain consistent, allowing a static determination of $\theta$.

**Dynamic Outlier Detection.** To dynamically detect outlier channels **during inference** at each time step, we employ an adaptive selection mechanism alongside an outlier list ($O_{\texttt{list}}$) initialized to NULL that tracks outlier channels. At each time-step $t$, we analyze the activation tensor *online* and compute a dynamic scale factor $S^D(t)$ of the whole tensor not including channels in $O_{\texttt{list}}$. For $S^D(t) \leq S^I(t)$, it indicates that no new outlier channels are present. However, if $S^D(t) > S^I(t)$, the increase in scale factor suggests the presence of high magnitude outliers. In this case, we iterate through all activation channels, and compute the per-channel scale factor over each channel of size $N \times 1$. We then compare the per-channel scale factor with the threshold $\theta$ to identify outliers channels. The identified outlier channels are added to $O_{\texttt{list}}$. $O_{\texttt{list}}$ is propagated across time-steps accumulating new outlier channels at every time-step while also retaining outlier channels of previous time-steps. However, such an approach may result in outdated or transient outlier channels being retained. To prevent this, we introduce a periodic refresh mechanism that updates the $O_{\texttt{list}}$ every $n_{\texttt{refresh}}$ time-steps to NULL.

**Mixed Precision Quantization.** To prevent significant accuracy degradation, we keep dynamically identified outliers at higher precision as compared to the inliers. Therefore, at every time-step each outlier channel having dedicated scale factor, are symmetric quantized to $b_a^O$-bits via per-channel quantization. On the other hand, inliers having a shared scale factor $S^I(t)$ over the tensor, are symmetric group quantized to $b_a^I$-bits ($b_a^I < b_a^O$) following Equation 5.

**Weight Quantization.** We apply per-channel symmetric group quantization for the weights. In specific, we follow Equation 5 to apply $b_w$-bit per-channel static quantization, where each channel $c$ has a shared scale factor $S_c^W$.

# 6 EXPERIMENTAL RESULTS

## 6.1 EXPERIMENTAL SETUP

**Models and Datasets.** We evaluate OuroMamba across a range of VMM and hybrid Transformer-Mamba models, covering fundamental CV tasks including classification, detection, segmentation, and image generation tasks, as detailed below. Please see appendix on details of extension of OuroMamba-Quant to transformer layers.

*Image Classification.* We employ ImageNet-1K Deng et al. (2009) with 50K tests for ViM-T/S/B Zhu et al. (2024), VMamba-T/B Liu et al. (2025), LVMamba-S Huang et al. (2024) and hybrid MambaVision-T/B Hatamizadeh & Kautz (2025).

*Object Detection.* We use the COCO 2017 dataset Lin et al. (2015) having 20K test data. Following Ramachandran et al. (2024a); Li et al. (2023b), we use the Mask R-CNN He et al. (2017) framework with VMamba-S Liu et al. (2025) and MambaVision-T Hatamizadeh & Kautz (2025) as the backbones.

*Semantic Segmentation.* We avail the ADE20K dataset Zhou et al. (2019) with 3K test data encompassing 150 categories with VMamba-S Liu et al. (2025) and MambaVision-T Hatamizadeh & Kautz (2025) as the backbones. We employ the UperNet framework Xiao et al. (2018).

Table 1: Quantization accuracy comparison for ImageNet classification. 'R', 'S' signifies real and synthetic calibration data.

| Model | Method | Data | #Images | W/A | Top-1 | W/A | Top-1 |
|---|---|---|---|---|---|---|---|
| Vim-S | FP Baseline | - | - | 32/32 | 81.60 | 32/32 | 81.60 |
| | BRECQ Li et al. (2021) | R | 1024 | 4/8 | 57.32 | 4/4 | 22.46 |
| | DopQ-ViT Yang et al. (2024) | R | 1024 | 4/8 | 68.96 | 4/4 | 61.58 |
| | PTQ4VM Cho et al. (2025) | R | 256 | 4/8 | 74.37 | 4/4 | 69.60 |
| | QMamba Li et al. (2025b) | R | 1024 | 4/8 | 74.12 | 4/4 | 33.64 |
| | **OuroMamba (Ours)** | S | 128 | 4/8 | **79.81** | 4/4 | **75.93** |
| Vim-B | Baseline | - | - | 32/32 | 81.90 | 32/32 | 81.90 |
| | BRECQ Li et al. (2021) | R | 1024 | 4/8 | 61.52 | 4/4 | 27.34 |
| | DopQ-ViT Yang et al. (2024) | R | 1024 | 4/8 | 68.47 | 4/4 | 62.33 |
| | PTQ4VM Cho et al. (2025) | R | 256 | 4/8 | 71.02 | 4/4 | 55.60 |
| | QMamba Li et al. (2025b) | R | 1024 | 4/8 | 75.46 | 4/4 | 66.73 |
| | **OuroMamba (Ours)** | S | 128 | 4/8 | **80.17** | 4/4 | **77.34** |
| VMamba-B | FP Baseline | - | - | 32/32 | 83.90 | 32/32 | 83.90 |
| | BRECQ Li et al. (2021) | R | 1024 | 4/8 | 62.34 | 4/4 | 25.65 |
| | DopQ-ViT Yang et al. (2024) | R | 1024 | 4/8 | 68.12 | 4/4 | 61.48 |
| | PTQ4VM Cho et al. (2025) | R | 256 | 4/8 | 78.95 | 4/4 | 75.67 |
| | QMamba Li et al. (2025b) | R | 1024 | 4/8 | 76.12 | 4/4 | 59.35 |
| | **OuroMamba (Ours)** | S | 128 | 4/8 | **82.03** | 4/4 | **78.91** |
| LVMamba-S | FP Baseline | - | - | 32/32 | 83.70 | 32/32 | 83.70 |
| | BRECQ Li et al. (2021) | R | 1024 | 4/8 | 55.34 | 4/4 | 11.62 |
| | DopQ-ViT Yang et al. (2024) | R | 1024 | 4/8 | 61.37 | 4/4 | 54.96 |
| | PTQ4VM Cho et al. (2025) | R | 256 | 4/8 | 81.46 | 4/4 | 78.21 |
| | QMamba Li et al. (2025b) | R | 1024 | 4/8 | 75.19 | 4/4 | 64.24 |
| | **OuroMamba (Ours)** | S | 128 | 4/8 | **82.94** | 4/4 | **80.11** |
| **Hybrid Model** MambaVision-B | FP Baseline | - | - | 32/32 | 84.20 | 32/32 | 84.20 |
| | BRECQ Li et al. (2021) | R | 1024 | 4/8 | 64.59 | 4/4 | 47.38 |
| | DopQ-ViT Yang et al. (2024) | R | 1024 | 4/8 | 70.68 | 4/4 | 61.43 |
| | PTQ4VM Cho et al. (2025) | R | 256 | 4/8 | 76.51 | 4/4 | 71.27 |
| | QMamba Li et al. (2025b) | R | 1024 | 4/8 | 76.08 | 4/4 | 68.59 |
| | **OuroMamba (Ours)** | S | 128 | 4/8 | **82.97** | 4/4 | **79.24** |

*Image Generation.* We evaluate generation performance using two high-resolution datasets FacesHQ Karras et al. (2019) and LandscapesHQ Skorokhodov et al. (2021) on Zigma, a zigzag Mamba backbone based diffusion model Hu et al. (2024).

**Baselines for Comparison.** We compare OuroMamba with state-of-the-art (SoTA) VMM PTQ techniques, PTQ4VM Cho et al. (2025) and QMamba Li et al. (2025b). Additionally, we compare with ViT PTQ methods, BREC-Q Li et al. (2021) and DopQ-ViT Yang et al. (2024), applying them to VMMs for a comprehensive evaluation.

**Implementation Details.** OuroMamba is implemented in PyTorch. All experiments are conducted on a single NVIDIA A100 GPU. Please see Table 5 for different hyperparameter values used in the evaluations.

**Task-Specific Synthetic Data Generation.** For calibration, we use $\mathcal{B} = 128$ synthetic samples. Following Ramachandran et al. (2024a), for image classification on the ImageNet-1K, we create $T_{G_\mathcal{B}} \in \mathcal{R}^{\mathcal{B} \times 1000}$, where the class-wise probabilities are randomly determined and assigned. Similarly, the target for object detection is $T_{G_\mathcal{B}} \in \mathcal{R}^{\mathcal{B} \times bb \times 5}$ where $bb$ is the number of bounding boxes in the image that is randomly selected from the integer set $[1, 3]$. $T_{G_\mathcal{B}}[\mathcal{B}, :, 0]$ corresponds to the bounding box category and $T_{G_\mathcal{B}}[\mathcal{B}, :, 1 : 4]$ is the bounding box coordinates $x, y, w, h$ Huang et al. (2018). For segmentation, the target is a pixel-wise classification map of the same size as $X_\mathcal{B}$.

**Quantization Setting Notation.** Following prior work Lin et al. (2024); Chen et al. (2024c), the notation for a mixed-precision scheme such as OuroMamba is, $W(b_w)A(b_a^I)O(b_a^O)$. Since, in our evaluations we fix $b_a^O = 8$, we omit $O8$ from the results for brevity.

Table 2: Quantization performance comparison for detection on COCO 2017. 'R', 'S' signifies real and synthetic calibration data.

| Method | Data | VMamba-S | | | MambaVision-T | | |
|---|---|---|---|---|---|---|---|
| | | W/A | $AP^{box}$ | $AP^{mask}$ | W/A | $AP^{box}$ | $AP^{mask}$ |
| Baseline | - | 32/32 | 48.7 | 43.7 | 32/32 | 46.4 | 41.8 |
| BREC-Q Li et al. (2021) | R | 4/8 | 35.6 | 33.7 | 4/8 | 39.9 | 40.1 |
| DopQ-ViT Yang et al. (2024) | R | 4/8 | 39.2 | 39.1 | 4/8 | 40.8 | 40.5 |
| PTQ4VM Cho et al. (2025) | R | 4/8 | **48.4** | 42.0 | 4/8 | **46.2** | 40.9 |
| QMamba Li et al. (2025b) | R | 4/8 | 47.2 | 41.3 | 4/8 | 45.3 | 39.9 |
| **OuroMamba (Ours)** | S | 4/8 | 48.3 | **42.9** | 4/8 | 46.1 | **41.4** |
| BREC-Q Li et al. (2021) | R | 4/4 | 27.3 | 25.4 | 4/4 | 29.6 | 25.1 |
| DopQ-ViT Yang et al. (2024) | R | 4/4 | 31.5 | 31.1 | 4/4 | 30.9 | 29.7 |
| PTQ4VM Cho et al. (2025) | R | 4/4 | 45.6 | 41.3 | 4/4 | 43.8 | 39.4 |
| QMamba Li et al. (2025b) | R | 4/4 | 43.7 | 39.1 | 4/4 | 42.1 | 38.6 |
| **OuroMamba (Ours)** | S | 4/4 | **47.8** | **42.5** | 4/4 | **44.9** | **40.9** |

## 6.2 ANALYSIS OF GENERATED SYNTHETIC DATA

In Figure 5, we visualize the synthetic samples generated by OuroMamba-Gen for Vim-B Zhu et al. (2024) for image classification (refer appendix for additional visualizations). Evidently, OuroMamba-Gen generates, realistic and clear class-specific foreground objects in contextually suitable background, showcasing a sophisticated understanding of semantic relationships between patches. This can be attributed to proposed enhanced impliict attention representation that is able to capture complex spatial relationships. This is in stark contrast to the naive synthetic samples in fig. 4(b), which are noisy and lack semantic structure.

## 6.3 QUANTIZATION RESULTS FOR CLASSIFICATION

As shown in Table 1, existing ViT PTQ techniques Yang et al. (2024); Li et al. (2021) suffer from significant accuracy degradations–particularly at W4A4–of up to 72% on VMM models and up to 37% on hybrid models. Clearly demonstrating the limitations of ViT PTQ methods in handling the dynamic outlier characteristics of VMMs. On the other hand, baseline VMM PTQ methods PTQ4VM Cho

Table 3: Quantization performance comparison for segmentation on ADE20K. 'R', 'S' signifies real and synthetic calibration data.

| Method | Data | VMamba-S | | MambaVision-T | |
|---|---|---|---|---|---|
| | | W/A | mIoU | W/A | mIoU |
| Baseline | - | 32/32 | 50.6 | 32/32 | 46.6 |
| PTQ4VM Cho et al. (2025) | R | 4/4 | 42.6 | 4/4 | 40.6 |
| QMamba Li et al. (2025b) | R | 4/4 | 40.7 | 4/4 | 38.5 |
| **OuroMamba (Ours)** | S | 4/4 | **47.3** | 4/4 | **44.1** |

et al. (2025), QMamba Li et al. (2025b) achieve acceptable performance at W4A8. However, at W4A4, they face large accuracy drops of up to 47%, due to static temporal grouping, per-token static scale factors. Furthermore, PTQ4VM's migration of outliers from activations to weights, complicates weight quantization, contributing to its accuracy degradation. In contrast, OuroMamba consistently outperforms baselines across different models and W/A settings, achieving near-lossless performance, with an accuracy drop of $\leq 1.8\%$ at W4A8. Similarly, at W4A4 OuroMamba yields an **accuracy improvement of 7.84% over PTQ4VM and 19.40% over QMamba on average, at W4A4**. It is important to note that, OuroMamba's superior performance is achieved in a data-free scenario requiring only 128 synthetic data samples, compared to PTQ4VM and QMamba requiring 256 and 1024 real calibration data.

## 6.4 QUANTIZATION RESULTS FOR OBJECT DETECTION

In Table 2 we present the quantization performance of OuroMamba and the baselines for object detection. Across different quantization settings and models, OuroMamba consistently outperforms the alternatives yielding up to **21.1** and **18.1 higher box AP and mask AP**, respectively.

## 6.5 QUANTIZATION RESULTS FOR SEGMENTATION

We present our evaluation for segmentation in Table 3. OuroMamba at W4A4 achieves upto **6.6 higher mIoU as compared to PTQ4VM and QMamba**.

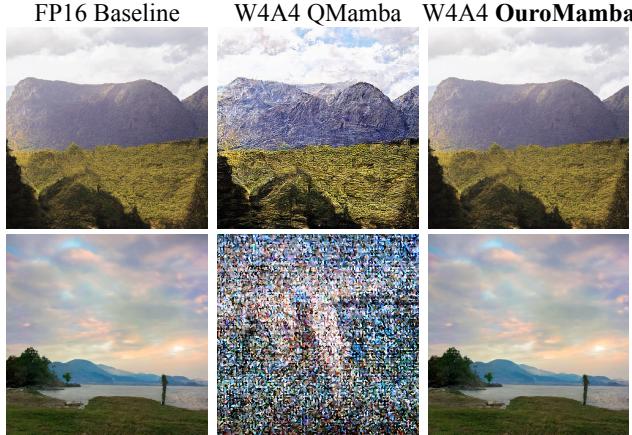FP16 Baseline  W4A4 QMamba  W4A4 **OuroMamba**



Figure 6: Visualization of landscape Skorokhodov et al. (2021) image generation by Zigma under different W4A4 quantization techniques.

## 6.6 QUANTIZATION RESULTS FOR DIFFUSION MODELS

**Quantitative Analysis.** Following Hu et al. (2024), we assess the image generation fidelity of Zigma, a VMM-based diffusion model, under W4A4 quantization. Fidelity is measured using the Fréchet Inception Distance (FID) score Heusel et al. (2017) across 5k samples, with results summarized in Table 4. Notably, across both evaluated high-resolution datasets, the **OuroMamba-quantized W4A4 Zigma model exhibits the smallest FID increase, maintaining high-quality image generation.** In contrast, PTQ4VM, QMamba have significant degradation, highlighting their limitations in preserving fidelity under ultra low-precisions.

**Qualitative Analysis.** In Figure 6, we qualitatively compare the generated image quality of OuroMamba with the second best technique from Table 4, QMamba, for the landscape image generation quality Skorokhodov et al. (2021). As observed in Figure 6, **OuroMamba generates images indistinguishable from the FP16 baseline**, whereas QMamba produces significant noise artifacts and/or unreadable content.

Table 4: Zigma image generation fidelity comparison under different quantization techniques.

| Method | Faces HQ | | Landscape HQ | |
|---|---|---|---|---|
| | W/A | FID ($\downarrow$) | W/A | FID ($\downarrow$) |
| Baseline | 16/16 | 37.8 | 16/16 | 10.7 |
| PTQ4VM Cho et al. (2025) | 4/4 | 89.6 | 4/4 | 46.2 |
| QMamba Li et al. (2025b) | 4/4 | 90.3 | 4/4 | 41.3 |
| **OuroMamba (Ours)** | 4/4 | **39.2** | 4/4 | **11.1** |

## 6.7 GEMM IMPLEMENTATION AND EVALUATION

**Implementation.** Efficient inference of Ouro-Mamba in VMMs requires addressing dynamic outlier extraction for real-time identification and mixed-precision GEMM operation for optimized computation. To this end, we implement the OuroMamba-`Quant` kernel using CUTLASS Thakkar et al. (2023). For efficient outlier channel extraction, activations are partitioned across channels and mapped to the thread blocks, where each block independently compares its assigned channels against



Figure 7: Overview of W4A4 hybrid GEMM kernel implementation.

the threshold $\theta$ to identify outliers and updates the $O_{\text{list}}$. To efficiently map the mixed-precision computation, we introduce a hybrid GEMM kernel. For a specific scenario of 4-bit weights with inlier and outlier activations of 4-bit and 8-bit, respectively, the hybrid GEMM executes one $4b \times 4b$ GEMM for inliers and one mixed-precision $4b \times 8b$ GEMM for outliers, as shown in Figure 7. We pack two consecutive 4-bit inlier activations into one byte, with outlier positions set to zero in the inlier buffer, and leverage the INT4 tensor cores for inlier GEMM. The 8-bit outlier channels are extracted and stored *compactly* in a separate contiguous outlier buffer. The outliers execute GEMM

on INT8 tensor cores performing computations only on the combined tensor of extracted outlier columns. Finally, the outputs from the inlier and outlier GEMMs are dequantized to FP16 and summed together, with all steps fused into a single pipeline for efficiency.

**Evaluation.** As shown in Figure 8, OuroMamba-`Quant` achieves speedup of **2.36× on Vim-B**, **1.61× on VMamba-S and 1.63× on Zigma** over FP16 baselines with a $\mathcal{B} = 128$. The speedup primarily stems from optimizations performed for efficient outlier extraction, mixed-precision GEMM through the INT4 and INT8 tensor cores. Unlike PTQ4VM, that uses a standalone kernel for dequantization, *we fuse dequantization directly into the GEMM pipeline, eliminating associated overhead and improving latency*. Please refer appendix for additional details.

## 6.8 Discussions and Ablations

**Effect of Batch Size $\mathcal{B}$.** We present the accuracy comparison across batch sizes from 16 to 1024 in Figure 9(a) for W4A8 Vim-S on image classification. OuroMamba exhibits minimal accuracy gains beyond 128, validating this as the optimal batch size. In contrast, VMM PTQ techniques PTQ4VM and QMamba rely on significantly larger batch sizes of 256 and 1024 real-data samples, respectively. Notably, even at batch sizes below 128, OuroMamba consistently outperforms PTQ4VM and QMamba, highlighting the effectiveness of its synthetic



Figure 8: End-to-end speedup over FP16 baseline.

data generation which is tailored for quantization of the respective model.

**Effect of Neighborhood Size $\mathcal{N}$ and Top-$n$ Patches.** Figure 9(b) demonstrates the effect of different neighborhood sizes $\mathcal{N}$ and top-$n$ positive patch selection on the quality of synthetic data generated by OuroMamba-`Gen`. To assess this, we measure the top-1 accuracy of W4A8 Vim-S using the generated samples are as for calibration. Notably, a $5 \times 5$ neighborhood size with the top-12 positive patches yields the highest accuracy. This choice of $\mathcal{N}$ effectively captures complex spatial dependencies while remaining computationally efficient compared to using all patches in $\alpha_p$.

**Impact of Scanning Direction.** To demonstrate that OuroMamba-`Gen` is scanning direction agnostic, we generate synthetic data by applying it independently to each scanning direction of VMamba-B. Using the same evaluation setup, we measure the top-1 accuracy for W4A4 quantization. As shown in Fig. fig. 9(c), under delta weighting, we observe that regardless of the scanning direction used for synthetic data generation, the resulting quantized model achieves similar accuracy.

**Impact of Weighting Factor.** Figure 9(c) shows the effect of uniform weighting compared to the $\Delta(t)$ weighting on the synthetic data generation by OuroMamba-`Gen`. For W4A4 VMamba-B, following the same evaluation setup as above, we see independent of different scanning directions delta weighting consistently yields accuracy improvements.

**Choice of Objective Function.** The study in Table 6 evaluates the impact of different loss function components in $\mathcal{L}^{gen}$ on synthetic data generation effectiveness and their influence on the top-1 accuracy of W4A8 quantization for Vim-S. We empirically validate our choice of *patch-level contrastive learning* ($\mathcal{L}^C$) over the patch-similarity metric ($\mathcal{L}^{PSE}$) used in prior ViT DFQ techniques Li et al. (2022; 2023b), demonstrating its superior effectiveness in generating high-quality calibration data. As shown in Table 6, a linear combination of $\mathcal{L}^C$ and $\mathcal{L}^O$ yields the highest accuracy, while using only $\mathcal{L}^O$ results in the lowest, highlighting its limited role in capturing semantic relationships for synthetic data generation. Although $\mathcal{L}^{PSE}$ combined with $\mathcal{L}^O$ achieves acceptable accuracy, it still underperforms compared to $\mathcal{L}^{gen}$.

**Effect of $n_{\texttt{refresh}}$.** In Figure 9(d), we present the normalized speedup for W4A4 Vim-B and VMamba-B across different values of the periodic refresh parameter $n_{\texttt{refresh}}$. We observe that $n_{\texttt{refresh}} = 10$ time-steps yields the highest speedup. Reducing $n_{\texttt{refresh}}$ leads to a sharp decline due to frequent reinitialization of $O_{\texttt{list}}$ and repeated outlier extraction over the entire tensor. Conversely, increasing $n_{\texttt{refresh}}$ beyond 10 results in excessive outlier accumulation, increasing computational complexity and reducing speedup. Notably, when $n_{\texttt{refresh}} = Full$ (i.e., no refresh),
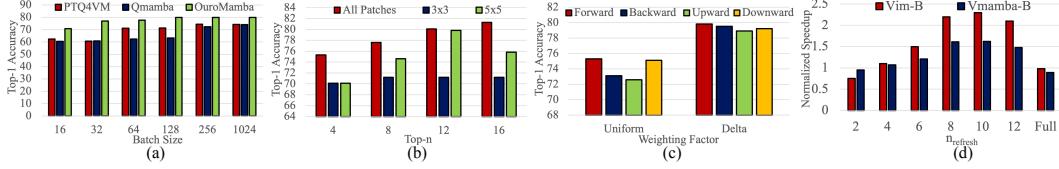
Figure 9: OuroMamba ablations for (a) Effect of batch size $\mathcal{B}$, (b) Effect of neighborhood size $\mathcal{N}$ and Top-$n$ positive patches, (c) Impact of weighting factor and scanning direction, (d) $n_{\texttt{refresh}}$ period.

Table 5: Hyperparameters used in evaluation.

| Parameter | Description | Value |
|---|---|---|
| $G$ | Gen. iterations | 1000 |
| $\mathcal{N}$ | Neighborhood size | $5 \times 5$ |
| $n$ | # Positive patches | 12 |
| $\mathcal{B}$ | Batch size | 128 |
| $n_{\texttt{ref.}}$ | Refresh period | 10 |
| $b_w$ | Weight precision | 4 |
| $b_a^I$ | Inlier act. precision | 4,8 |
| $b_a^O$ | Outlier act. precision | 8 |

Table 6: Impact of different loss components.

| $\mathcal{L}^{PSE}$ | $\mathcal{L}^C$ | $\mathcal{L}^O$ | W/A | Top-1 Acc. (%) |
|---|---|---|---|---|
| - | - | - | 32/32 | 81.60 |
| ✓ | ✗ | ✗ | 4/8 | 71.68 |
| ✗ | ✗ | ✓ | 4/8 | 21.65 |
| ✓ | ✗ | ✓ | 4/8 | 73.45 |
| ✗ | ✓ | ✗ | 4/8 | 75.52 |
| ✗ | ✓ | ✓ | 4/8 | **79.81** |
| ✓ | ✓ | ✓ | 4/8 | 74.32 |

OuroMamba-`Quant` becomes slower than the FP16 baseline due to the overhead of processing accumulated outliers.

**Limitation Discussion.** The dynamic outlier selection assumes that inlier distributions remain stable relative to their values determined during calibration. In our experiments across all evaluated models, we observed no significant variations in the dynamic range of inliers. However, if future model architectures exhibit substantial fluctuations in inlier values, additional investigation will be necessary to adapt the selection process accordingly.

## REFERENCES

Ameen Ali, Itamar Zimerman, and Lior Wolf. The hidden attention of mamba models. *arXiv preprint arXiv:2403.01590*, 2024.

Seyedarmin Azizi, Souvik Kundu, Mohammad Sadeghi, and Massoud Pedram. MambaExtend: A training-free approach to improve long context extension of mamba. *International Conference on Learning Representation*, 2025.

Junsong Chen, Chongjian Ge, Enze Xie, Yue Wu, Lewei Yao, Xiaozhe Ren, Zhongdao Wang, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart-sigma: Weak-to-strong training of diffusion transformer for 4k text-to-image generation, 2024a. URL https://arxiv.org/abs/2403.04692.

Lei Chen, Yuan Meng, Chen Tang, Xinzhu Ma, Jingyan Jiang, Xin Wang, Zhi Wang, and Wenwu Zhu. Q-dit: Accurate post-training quantization for diffusion transformers, 2024b. URL https://arxiv.org/abs/2406.17343.

Yidong Chen, Chen Zhang, Rongchao Dong, Haoyuan Zhang, Yonghua Zhang, Zhonghua Lu, and Jidong Zhai. Mixq: Taming dynamic outliers in mixed-precision quantization by online prediction. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage, and Analysis*, SC '24. IEEE Press, 2024c. ISBN 9798350352917. doi: 10.1109/SC41406.2024.00080. URL https://doi.org/10.1109/SC41406.2024.00080.

Hung-Yueh Chiang, Chi-Chih Chang, Natalia Frumkin, Kai-Chiang Wu, and Diana Marculescu. Quamba: A post-training quantization recipe for selective state space models. *arXiv preprint arXiv:2410.13229*, 2024.

Younghyun Cho, Changhun Lee, Seonggon Kim, and Eunhyeok Park. PTQ4VM: Post-training quantization for visual mamba. *Winter Conference on Application of Computer Vision*, 2025.

Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.

Peiyan Dong, Lei Lu, Chao Wu, Cheng Lyu, Geng Yuan, Hao Tang, and Yanzhi Wang. Packqvit: Faster sub-8-bit vision transformers via full and packed quantization on the mobile. *Advances in Neural Information Processing Systems*, 36:9015–9028, 2023.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.

Dongchen Han, Ziyi Wang, Zhuofan Xia, Yizeng Han, Yifan Pu, Chunjiang Ge, Jun Song, Shiji Song, Bo Zheng, and Gao Huang. Demystify mamba in vision: A linear attention perspective. *Advances in neural information processing systems*, 37:127181–127203, 2025.

Ali Hatamizadeh and Jan Kautz. Mambavision: A hybrid mamba-transformer vision backbone. *In Proceedings of the Computer Vision and Pattern Recognition*, 2025.

Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pp. 2961–2969, 2017.

Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.

Vincent Tao Hu, Stefan Andreas Baumann, Ming Gui, Olga Grebenkova, Pingchuan Ma, Johannes Fischer, and Björn Ommer. Zigma: A dit-style zigzag mamba diffusion model. In *European Conference on Computer Vision*, pp. 148–166. Springer, 2024.

He Huang, Philip S Yu, and Changhu Wang. An introduction to image synthesis with generative adversarial nets. *arXiv preprint arXiv:1803.04469*, 2018.

Tao Huang, Xiaohuan Pei, Shan You, Fei Wang, Chen Qian, and Chang Xu. Localmamba: Visual state space model with windowed selective scan. *arXiv preprint arXiv:2403.09338*, 2024.

Yanfeng Jiang, Ning Sun, Xueshuo Xie, Fei Yang, and Tao Li. Adfq-vit: Activation-distribution-friendly post-training quantization for vision transformers. *Neural Networks*, pp. 107289, 2025.

Tero Karras, Samuli Laine, and Timo Aila. Flickr-faces-hq dataset (ffhq). GitHub repository, 2019. URL https://github.com/NVlabs/ffhq-dataset.

Minjun Kim, Jongjin Kim, and U Kang. Synq: Accurate zero-shot quantization by synthesis-aware fine-tuning. In *The Thirteenth International Conference on Learning Representations*, 2025.

Souvik Kundu, Qirui Sun, Yao Fu, Massoud Pedram, and Peter Beerel. Analyzing the confidentiality of undistillable teachers in knowledge distillation. *Advances in Neural Information Processing Systems*, 34:9181–9192, 2021.

Souvik Kundu, Shikai Wang, Qirui Sun, Peter A Beerel, and Massoud Pedram. Bmpq: bit-gradient sensitivity-driven mixed-precision quantization of dnns from scratch. In *2022 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, pp. 588–591. IEEE, 2022.

Muyang Li, Yujun Lin, Zhekai Zhang, Tianle Cai, Xiuyu Li, Junxian Guo, Enze Xie, Chenlin Meng, Jun-Yan Zhu, and Song Han. Svdquant: Absorbing outliers by low-rank components for 4-bit diffusion models, 2025a. URL https://arxiv.org/abs/2411.05007.

Xiuyu Li, Yijiang Liu, Long Lian, Huanrui Yang, Zhen Dong, Daniel Kang, Shanghang Zhang, and Kurt Keutzer. Q-diffusion: Quantizing diffusion models, 2023a. URL https://arxiv.org/abs/2302.04304.

Yinglong Li, Xiaoyu Liu, Jiacheng Li, Ruikang Xu, Yinda Chen, and Zhiwei Xiong. Qmamba: Post-training quantization for vision state space models. *arXiv preprint arXiv:2501.13624*, 2025b.

Yuhang Li, Ruihao Gong, Xu Tan, Yang Yang, Peng Hu, Qi Zhang, Fengwei Yu, Wei Wang, and Shi Gu. Brecq: Pushing the limit of post-training quantization by block reconstruction. *arXiv preprint arXiv:2102.05426*, 2021.

Zhikai Li, Liping Ma, Mengjuan Chen, Junrui Xiao, and Qingyi Gu. Patch similarity aware data-free quantization for vision transformers. In *European conference on computer vision*, pp. 154–170. Springer, 2022.

Zhikai Li, Mengjuan Chen, Junrui Xiao, and Qingyi Gu. Psaq-vit v2: Toward accurate and general data-free quantization for vision transformers. *IEEE Transactions on Neural Networks and Learning Systems*, 2023b.

Zhikai Li, Junrui Xiao, Lianwei Yang, and Qingyi Gu. Repq-vit: Scale reparameterization for post-training quantization of vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 17227–17236, 2023c.

Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common objects in context, 2015. URL https://arxiv.org/abs/1405.0312.

Yang Lin, Tianyu Zhang, Peiqin Sun, Zheng Li, and Shuchang Zhou. Fq-vit: Post-training quantization for fully quantized vision transformer. *arXiv preprint arXiv:2111.13824*, 2021.

Yujun Lin, Haotian Tang, Shang Yang, Zhekai Zhang, Guangxuan Xiao, Chuang Gan, and Song Han. Qserve: W4a8kv4 quantization and system co-design for efficient llm serving. *arXiv preprint arXiv:2405.04532*, 2024.

Yue Liu, Yunjie Tian, Yuzhong Zhao, Hongtian Yu, Lingxi Xie, Yaowei Wang, Qixiang Ye, Jianbin Jiao, and Yunfan Liu. Vmamba: Visual state space model. *Advances in neural information processing systems*, 37:103031–103063, 2025.

Yuexiao Ma, Huixia Li, Xiawu Zheng, Feng Ling, Xuefeng Xiao, Rui Wang, Shilei Wen, Fei Chao, and Rongrong Ji. Outlier-aware slicing for post-training quantization in vision transformer. In *Forty-first International Conference on Machine Learning*, 2024.

Alessandro Pierro and Steven Abreu. Mamba-ptq: Outlier channels in recurrent large language models. *arXiv preprint arXiv:2407.12397*, 2024.

Akshat Ramachandran, Souvik Kundu, and Tushar Krishna. Clamp-ViT: Contrastive data-free learning for adaptive post-training quantization of vits. In *European Conference on Computer Vision*, pp. 307–325. Springer, 2024a.

Akshat Ramachandran, Souvik Kundu, and Tushar Krishna. Microscopiq: Accelerating foundational models through outlier-aware microscaling quantization. *arXiv preprint arXiv:2411.05282*, 2024b.

Akshat Ramachandran, Zishen Wan, Geonhwa Jeong, John Gustafson, and Tushar Krishna. Algorithm-hardware co-design of distribution-aware logarithmic-posit encodings for efficient dnn inference. *arXiv preprint arXiv:2403.05465*, 2024c.

Bo-Yun Shi, Yi-Cheng Lo, et al. Post-training quantization for vision mamba with k-scaled quantization and reparameterization. *arXiv preprint arXiv:2501.16738*, 2025.

Ivan Skorokhodov, Grigorii Sotnikov, and Mohamed Elhoseiny. Aligning latent and image spaces to connect the unconnectable, 2021. URL https://arxiv.org/abs/2104.06954.

Vijay Thakkar, Pradeep Ramani, Cris Cecka, Aniket Shivam, Honghao Lu, Ethan Yan, Jack Kosaian, Mark Hoemmen, Haicheng Wu, Andrew Kerr, Matt Nicely, Duane Merrill, Dustyn Blasig, Fengqi Qiao, Piotr Majcher, Paul Springer, Markus Hohnerbach, Jin Wang, and Manish Gupta. CUTLASS, January 2023. URL https://github.com/NVIDIA/cutlass.

Jing Wang, Jiangyun Li, Wei Li, Lingfei Xuan, Tianxiang Zhang, and Wenxuan Wang. Positive–negative equal contrastive loss for semantic segmentation. *Neurocomputing*, 535:13–24, 2023.

Junyi Wu, Haoxuan Wang, Yuzhang Shang, Mubarak Shah, and Yan Yan. Ptq4dit: Post-training quantization for diffusion transformers, 2024. URL `https://arxiv.org/abs/2405.16005`.

Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. Smoothquant: Accurate and efficient post-training quantization for large language models. In *International Conference on Machine Learning*, pp. 38087–38099. PMLR, 2023.

Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 418–434, 2018.

Yicheng Xiao, Lin Song, Shaoli Huang, Jiangshan Wang, Siyu Song, Yixiao Ge, Xiu Li, and Ying Shan. Grootvl: Tree topology is all you need in state space model. *arXiv preprint arXiv:2406.02395*, 2024.

Lianwei Yang, Haisong Gong, and Qingyi Gu. Dopq-vit: Towards distribution-friendly and outlier-aware post-training quantization for vision transformers. *arXiv preprint arXiv:2408.03291*, 2024.

Zhihang Yuan, Chenhao Xue, Yiqi Chen, Qiang Wu, and Guangyu Sun. Ptq4vit: Post-training quantization for vision transformers with twin uniform quantization. In *European conference on computer vision*, pp. 191–207. Springer, 2022.

Yuke Zhang, Dake Chen, Souvik Kundu, Chenghao Li, and Peter A Beerel. SAL-ViT: Towards latency efficient private inference on vit using selective attention search with a learnable softmax approximation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 5116–5125, 2023.

Tianchen Zhao, Tongcheng Fang, Haofeng Huang, Enshu Liu, Rui Wan, Widyadewi Soedarmadji, Shiyao Li, Zinan Lin, Guohao Dai, Shengen Yan, Huazhong Yang, Xuefei Ning, and Yu Wang. Vidit-q: Efficient and accurate quantization of diffusion transformers for image and video generation, 2025. URL `https://arxiv.org/abs/2406.02540`.

Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ade20k dataset. *International Journal of Computer Vision*, 127:302–321, 2019.

Lianghui Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang, Wenyu Liu, and Xinggang Wang. Vision mamba: Efficient visual representation learning with bidirectional state space model, 2024.

## A  APPENDIX

## B  EXTENDED RELATED WORKS

Existing PTQ techniques for ViTs Lin et al. (2021); Jiang et al. (2025); Li et al. (2023c); Yang et al. (2024); Ma et al. (2024); Yuan et al. (2022) address long-tailed distributions and static activation outliers to enhance quantization accuracy. For instance, DopQ-ViT Yang et al. (2024) and ADFQ-ViT Jiang et al. (2025) mitigate outliers by optimizing per-channel and per-patch scale factors, respectively. FQ-ViT Lin et al. (2021) introduces Power-of-Two Factor (PTF) for inter-channel LayerNorm variation and Log-Int-Softmax (LIS) for 4-bit attention map quantization. Among the early PTQ methods for Mamba models, Mamba-PTQ Pierro & Abreu (2024) and Quamba Chiang et al. (2024) identified activation outliers as a key challenge but were tailored for language tasks. More recently, VMM PTQ techniques Cho et al. (2025); Li et al. (2025b) highlighted the highly dynamic activation distributions and inter-channel variations across time-steps. PTQ4VM Cho et al. (2025) adapts SmoothQuant Xiao et al. (2023) to migrate activation outliers into weights using a migration factor. However, as noted in Li et al. (2025a), this increases weight complexity, making both weights and activations more sensitive to dynamic variations, rendering it ineffective for ultra-low precision ($< 4$ bits) quantization. Additionally, PTQ4VM does not quantize SSM operators,

Table 7: Quantization accuracy comparison of SoTA techniques on ImageNet classification. 'R', 'S' signifies real and synthetic calibration data.

| | Method | Data | #Images | W/A | Top-1 | W/A | Top-1 | W/A | Top-1 |
|---|---|---|---|---|---|---|---|---|---|
| | Baseline | - | - | 32/32 | 76.10 | 32/32 | 76.10 | 32/32 | 76.10 |
| Vim-T | PTQ4VM Cho et al. (2025) | R | 256 | 4/8 | 74.15 | 6/6 | 73.94 | 4/4 | 56.29 |
| | QMamba Li et al. (2025b) | R | 1024 | 4/8 | 70.13 | 6/6 | 57.95 | 4/4 | 53.41 |
| | **OuroMamba (Ours)** | S | 128 | 4/8 | **74.98** | 6/6 | **74.84** | 4/4 | **63.49** |
| | Baseline | - | - | 32/32 | 82.60 | 32/32 | 82.60 | 32/32 | 82.60 |
| VMamba-T | PTQ4VM Cho et al. (2025) | R | 256 | 4/8 | 77.02 | 6/6 | 75.67 | 4/4 | 72.67 |
| | QMamba Li et al. (2025b) | R | 1024 | 4/8 | 76.51 | 6/6 | **80.49** | 4/4 | 51.48 |
| | **OuroMamba (Ours)** | S | 128 | 4/8 | **81.73** | 6/6 | 80.15 | 4/4 | **77.56** |
| | Baseline | - | - | 32/32 | 82.30 | 32/32 | 82.30 | 32/32 | 82.30 |
| **Hybrid Model** | PTQ4VM Cho et al. (2025) | R | 256 | 4/8 | 72.13 | 6/6 | 69.39 | 4/4 | 67.67 |
| MambaVision-T | QMamba Li et al. (2025b) | R | 1024 | 4/8 | 71.93 | 6/6 | 68.17 | 4/4 | 65.33 |
| | **OuroMamba (Ours)** | S | 128 | 4/8 | **80.57** | 6/6 | **79.05** | 4/4 | **74.92** |

limiting its scope to linear layer weights and output activations. QMamba Li et al. (2025b) addresses the dynamic inter-time-step variations in VMMs' hidden states by introducing fine-grained temporal grouped quantization, quantizing both weights and activations. Similarly, kSQ-VMM Shi et al. (2025) applies similarity-based k-scaled channel-wise and token-wise quantization to handle dynamic activation distributions. However, existing VMM PTQ methods rely on static scale factors Cho et al. (2025); Li et al. (2025b) or fixed temporal groupings Li et al. (2025b), leading to accuracy degradation at ultra-low bit precisions due to their inability to dynamically manage outlier channels at runtime.

## C   ADDITIONAL QUANTIZATION RESULTS

In Table 1 we provide additional quantization results of Vim-T Zhu et al. (2024), VMamba-T Liu et al. (2025) and the hybrid model MambaVision-T Hatamizadeh & Kautz (2025) for image classification.

## D   GEMM IMPLEMENTATION DETAILS

We first describe how the GEMM operation can be decomposed into separate computations for outliers and inliers. Consider an output element $Y[i, j]$ computed as

$$Y[i, j] = \sum_{k \in \mathcal{I}} A^I[i, k] \, W[k, j] + \sum_{k \in \mathcal{O}} A^O[i, k] \, W[k, j] \tag{8}$$

where the inlier activations $A^I$ have the outlier positions zeroed out, and the outlier activations $A^O$ have the inlier positions zeroed. This decomposition guarantees that the sum of the two partial GEMM results yields the same $Y[i, j]$ as the original full GEMM.

We now introduce our GEMM pipeline, which consists of the following five steps:

**Outlier Extraction.** Outlier values in the input activation are identified, and their corresponding positions are zeroed out. The outlier columns are then compacted into a small INT8 outlier buffer.

**Inlier Extraction.** With the outlier positions already zeroed out, the inlier values are extracted and packed into INT4 buffers, storing two values per byte.

**INT4 GEMM.** An INT4 GEMM is performed on the inlier data. During the CUTLASS epilogue, the results are immediately dequantized by multiplying with the activation and weight scales. This fusion is enabled by the use of per-tensor quantization for inliers, which offers greater efficiency compared to the per-token inlier quantization employed by PTQ4VM.

**INT8 GEMM.** A mixed-input INT4-INT8 GEMM is executed between the inlier and the compacted outlier matrices, utilizing INT8 tensor cores.

FP16 Baseline | W4A8 Q-DiT | W4A8 PTQ4DiT | W4A8 **OuroMamba** | W4A4 **OuroMamba**

Prompt: An astronaut relaxing on a beach chair, sipping coffee on Mars, with Earth visible in the sky
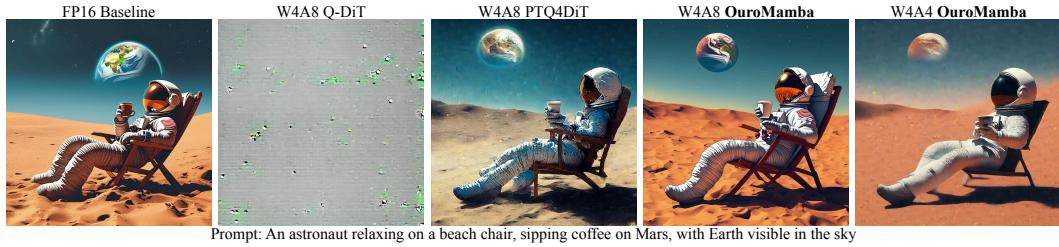
Figure 10: Quantization performance comparison for text-to-image generation task.

**Outlier Dequantization and Combination.** Finally, the dequantization of outliers and the combination of the two GEMM results are fused into a single kernel. This kernel is memory-bound because it writes the final result matrix.

## E    SPEED BREAKDOWN RESULTS

As shown in Figure 11, outlier extraction incurs minimal overhead. Specifically, we partition activation by channels, so outlier channel indices and scaling factors are calculated and recorded in parallel. Additionally, our compact extraction approach restricts the INT8 GEMM operation to a small subset of outlier channels, limiting its runtime contribution to less than 5%. As expected, the INT4 GEMM remains the dominant component. The primary performance bottleneck is the dequantization and combination step. This step writes to the entire output matrix, making it inherently memory-bound and therefore more expensive. Notably, the dequantization overhead is higher in Vmamba-B because it has a higher outlier rate (4.3%) compared to Vim-B (1.3%). Nonetheless, even in scenarios with higher outlier densities, our overall pipeline remains efficient due to the minimal costs associated with both outlier extraction and the outlier GEMM computations.
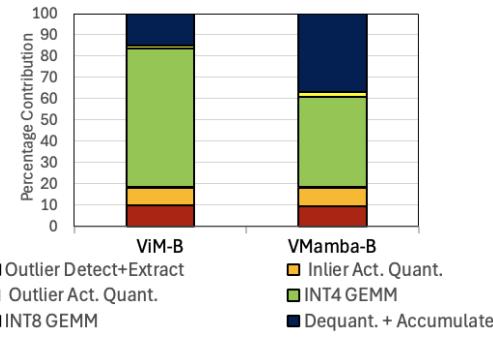


Figure 11: Kernel breakdown of OuroMamba-`Quant`.

## F    MEMORY COMPRESSION RESULTS

As shown in Table 8, we compare the memory compression factor of OuroMamba with PTQ4VM at W4A4, using the FP16 model as the baseline, on the Vim-S and Vim-B models Zhu et al. (2024). The results show that Ouro-Mamba consistently achieves a high memory compression factor of up to $3.80\times$ compared to the baseline FP16 model, while PTQ4VM achieves a memory compression factor of only $2.02\times$, as it quantizes only the Linear layers of VMMs.

Table 8: Memory compression comparison over PTQ4VM Cho et al. (2025).

| | Vim-S | | Vim-B | |
|---|---|---|---|---|
| Method | W/A | Mem. Comp. | W/A | Mem. Comp. |
| Baseline | 16/16 | 1.00 | 16/16 | 1.00 |
| PTQ4VM Cho et al. (2025) | 4/4 | 1.81 | 4/4 | 2.02 |
| **OuroMamba (Ours)** | 4/4 | **3.63** | 4/4 | **3.80** |

## G    EXTENSION OF OUROMAMBA-QUANT TO TRANSFORMER BASED MODELS

We extend OuroMamba-`Gen` to Transformer-based models and layers by mapping the time-step dimension to the token dimension. Outlier channels are identified per token, with $O_{\texttt{list}}$ propagated across tokens.
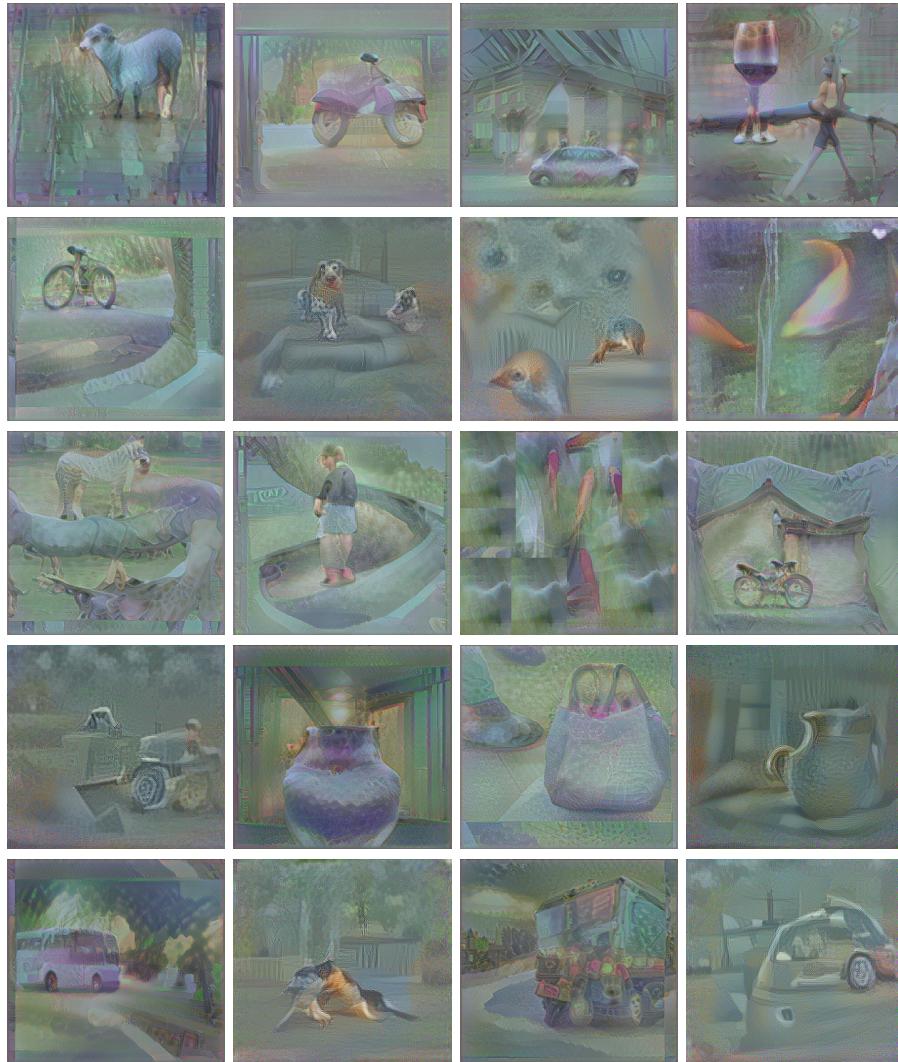
Figure 12: Generated synthetic data samples.

## H    TEXT-TO-IMAGE GENERATION RESULTS

**Implementation.**    We applied OuroMamba-`Quant` to PixArt-$\Sigma$ Chen et al. (2024a) with 20-iteration setting. Following ViDiT-Q Zhao et al. (2025), we quantize linear layers for query, key, and value projections and the second projection layer of feed-forward network to W4A4. Meanwhile, the first projection layer of feed-forward network and the output projection in self-attention are quantized in 8-bit for better numerical stability, while outliers bits are fixed at 8-bit and $n_{\text{refresh}}$ is set to 10. For calibration, we follow Q-Diffusion Li et al. (2023a) and randomly sample text prompts from MS-COCO dataset Lin et al. (2015) to obtain outlier threshold and inlier scale factors.

**Results.**    In Figure 10, we visualize the generated images of W4A8, W4A4 Ouromamba-`Quant` quantized PixArt-$\Sigma$ compared to W4A8 Q-DiT Chen et al. (2024b) and W4A8 PTQ4DiT Wu et al. (2024).

## I    ADDITIONAL SYNTHETIC DATA SAMPLES

In Figure 12, we additionally visualize synthetic samples generated by OuroMamba-`Gen` for image classification, object detection and segmentation tasks for Vim-B model Zhu et al. (2024).