

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

ICCV 2025

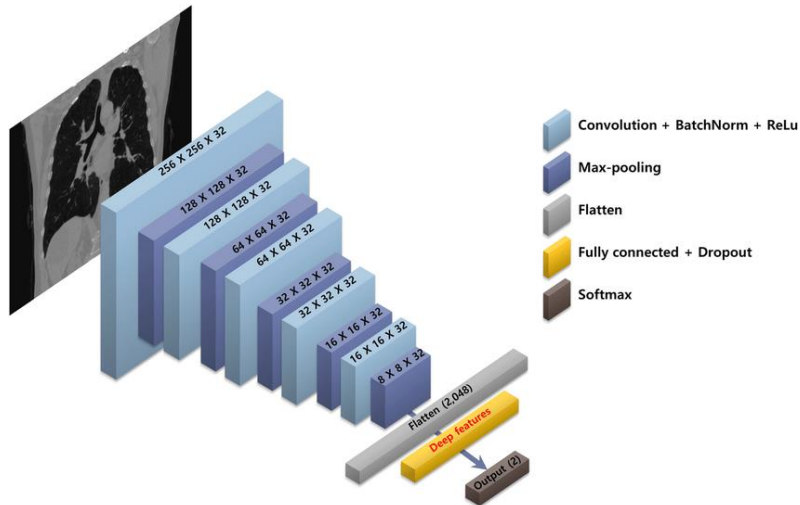
Akshat Ramachandran^{1*}, Mingyu Lee^{1*}, Huan Xu¹, Souvik Kundu², Tushar Krishna¹

^{*}Equal Contribution

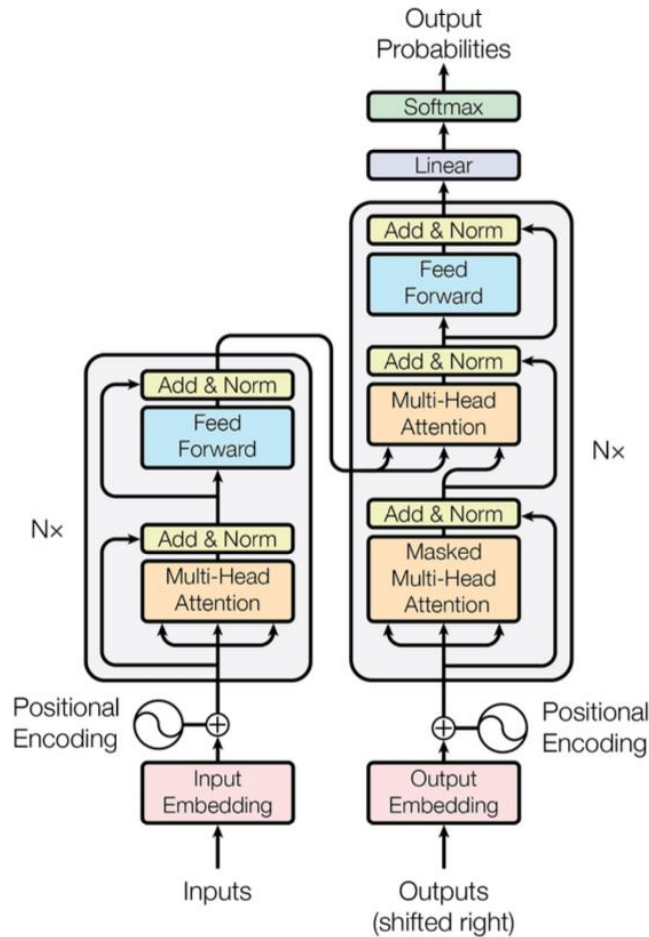
¹Georgia Institute of Technology ²Intel Labs

Contact: akshat.r@gatech.edu

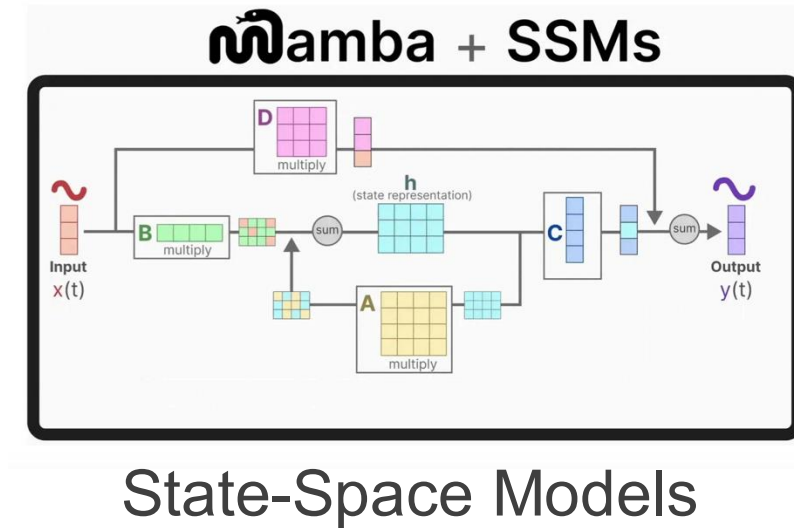
Diverse Space of AI Models



Convolutional Neural Networks

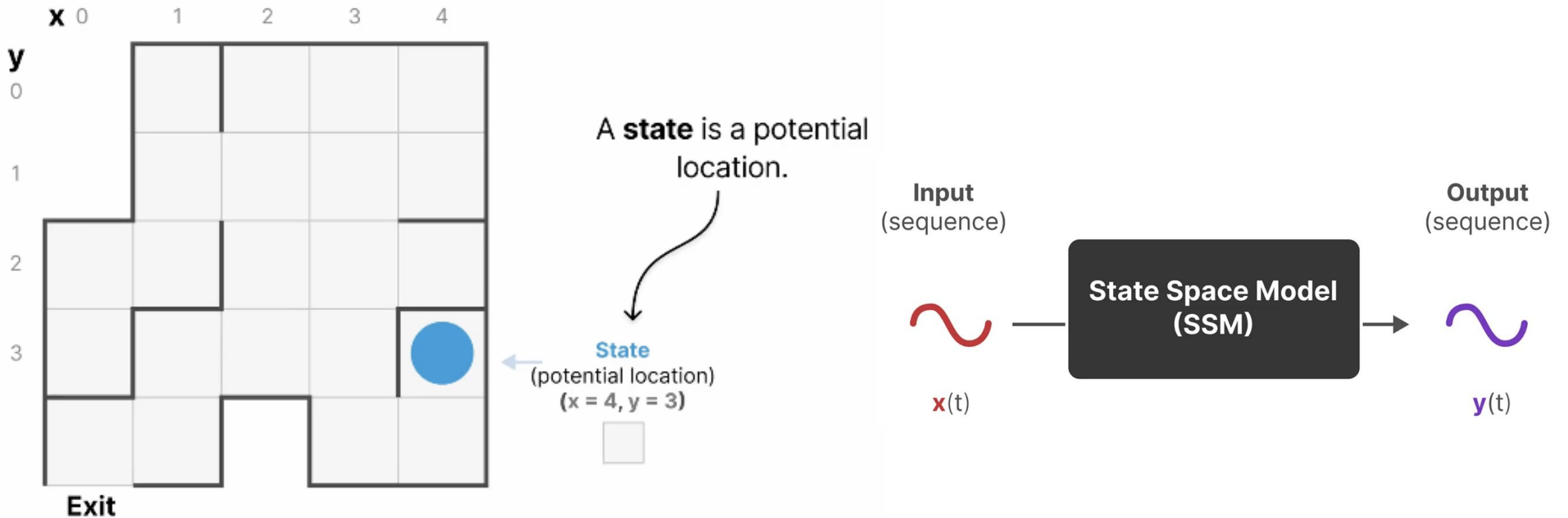


Transformer-based Models



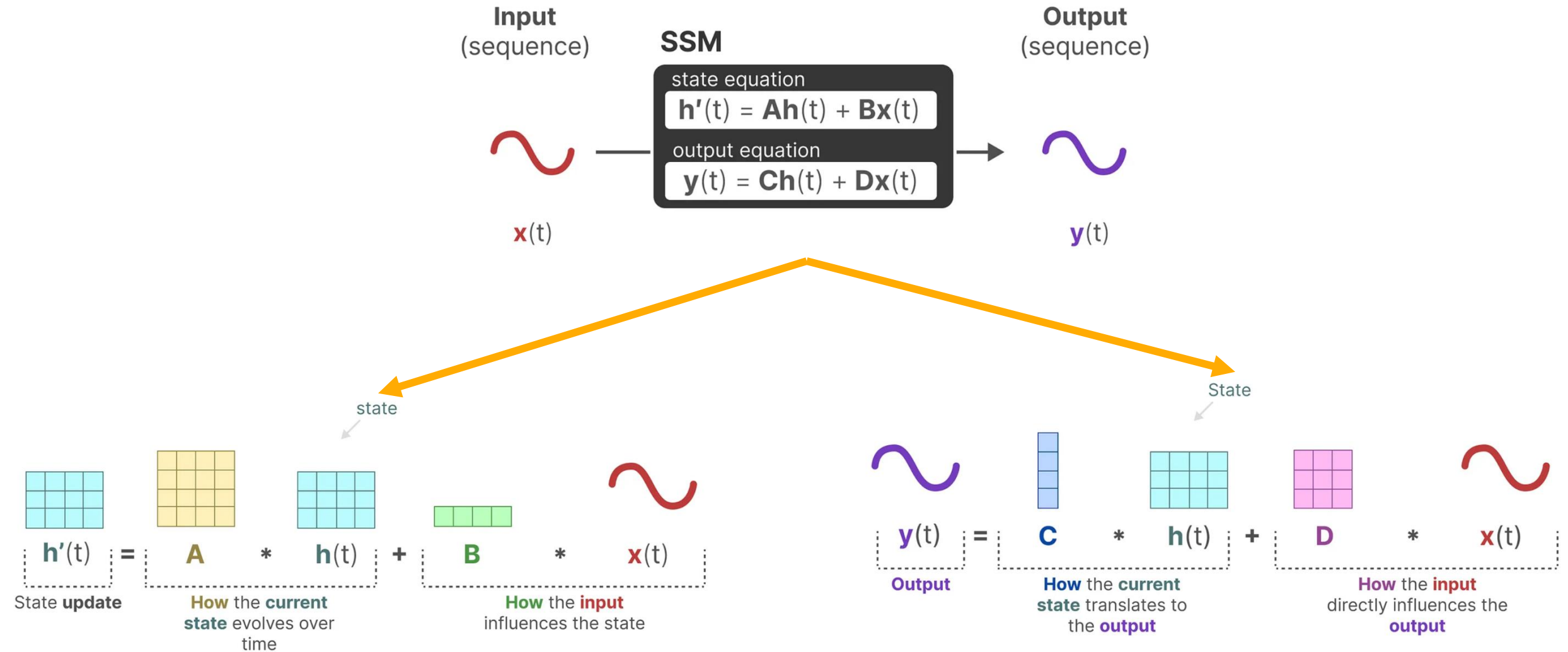
State-Space Models

State Space Model (SSM) Family

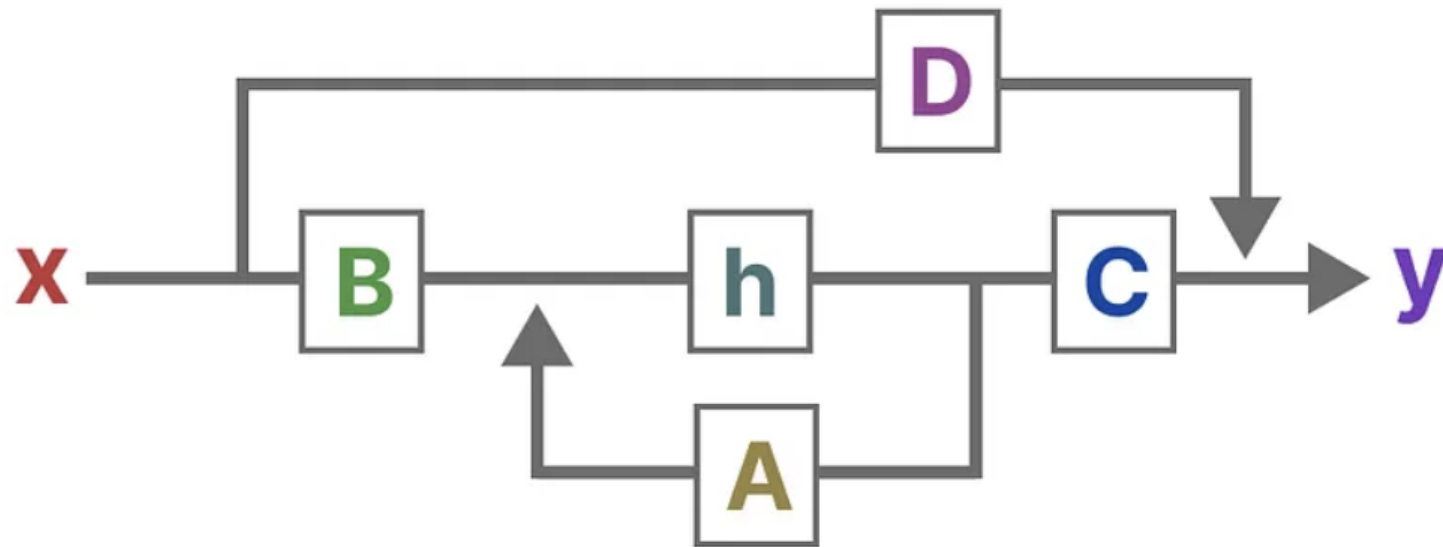


A State Space contains the minimum number of variables that fully describe a system. It is a way to mathematically represent a problem by defining a system's possible states.

The Two Core Equations of SSMs

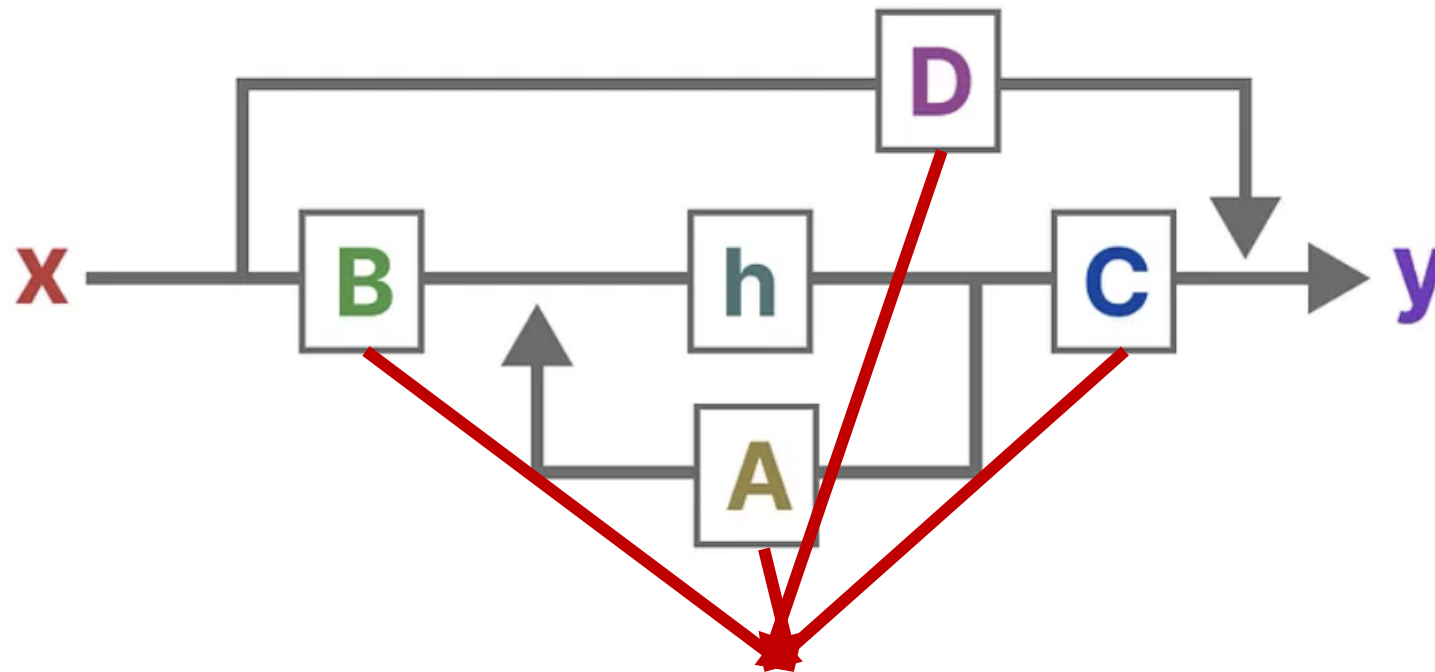


Simplified SSM Model Representation



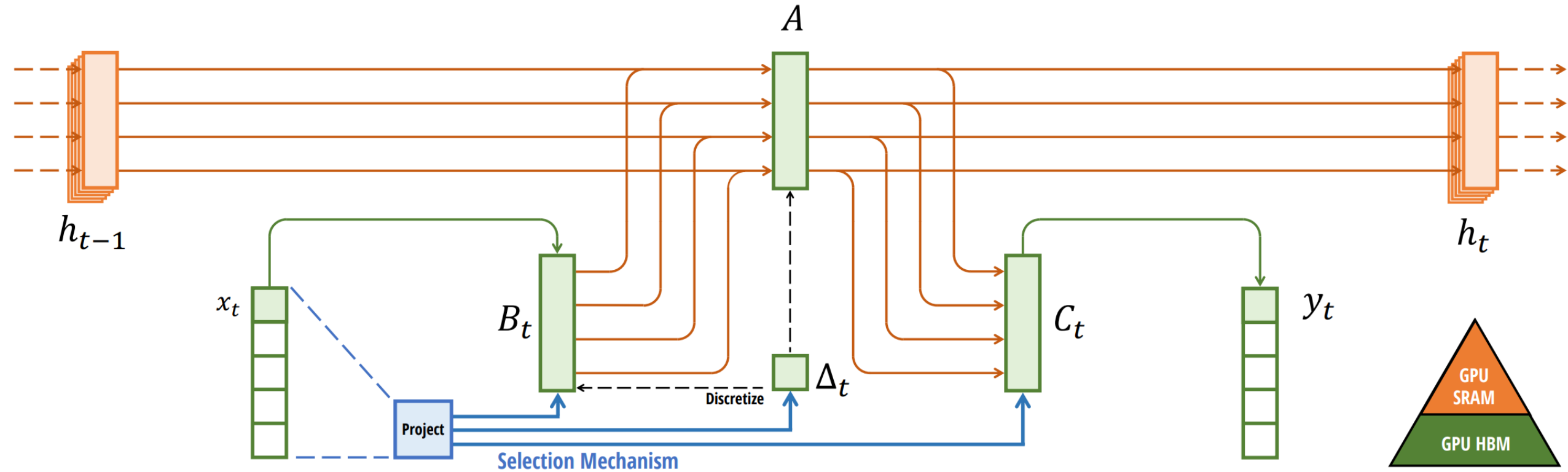
Continuous domain representation of SSMs

Simplified SSM Model Representation



Constant with different inputs x

Selective SSM or S6



Convert to discrete domain for LLMs!

Selective SSM or S6

State space equation

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

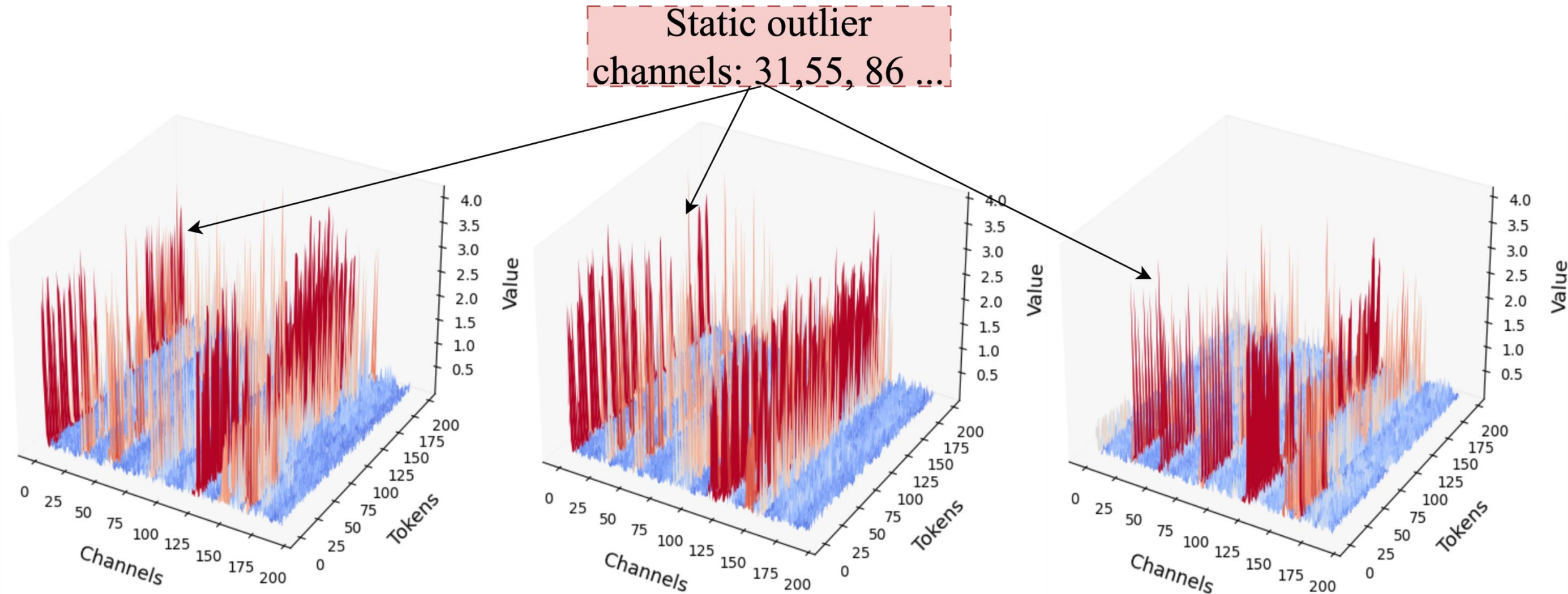
Discretization

$$\bar{A}(t) = e^{(A\Delta(t))}; B(t) = W_B(u(t)); \bar{B}(t) = B\Delta(t)$$

$$\Delta(t) = S^+(u(t)\Delta(t)_{\text{proj}}); C(t) = (W_C(z(t)))^T$$

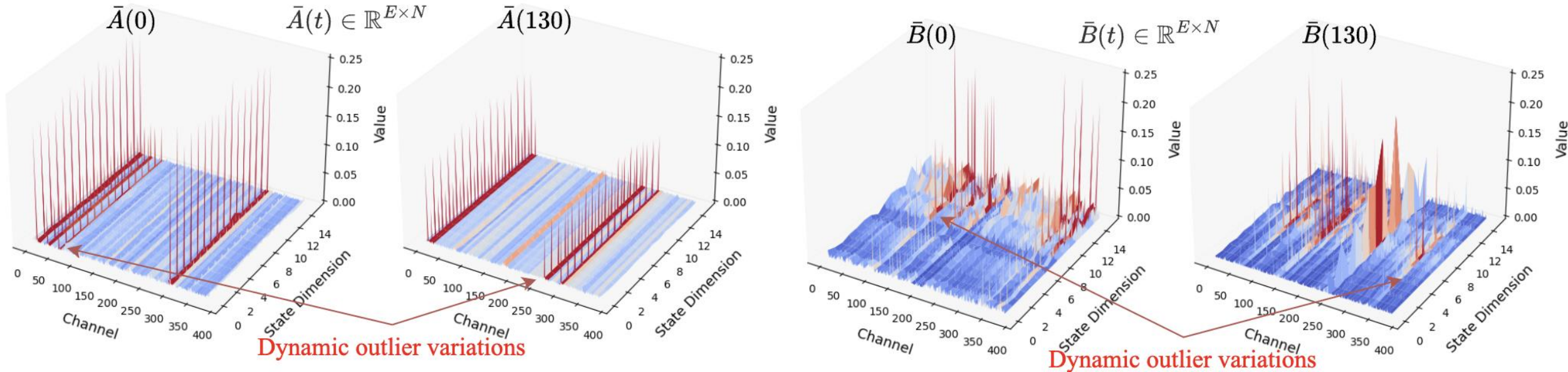
Convert to discrete domain for LLMs!

Transformers' Static Outlier Patterns



DeiT Layer 9 for different inputs

Dynamic Activation Variations



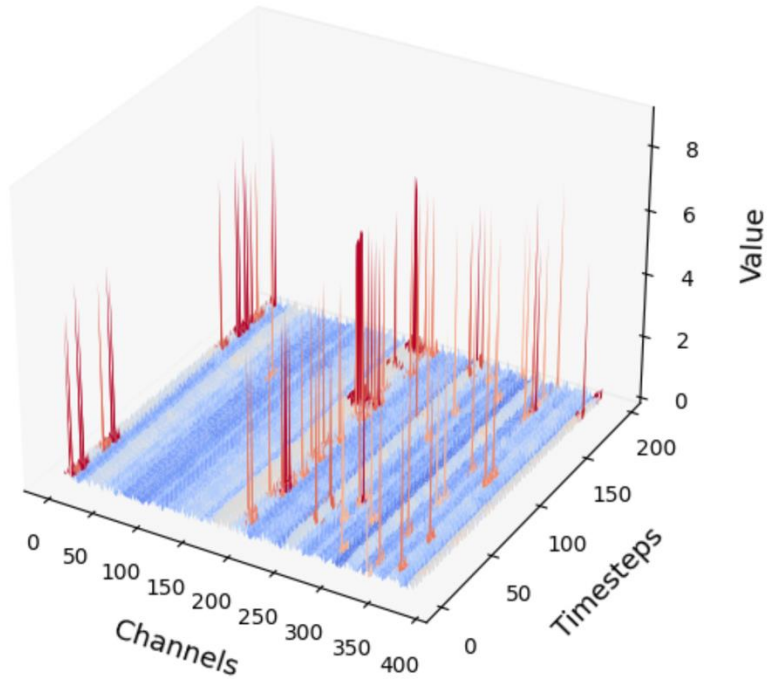
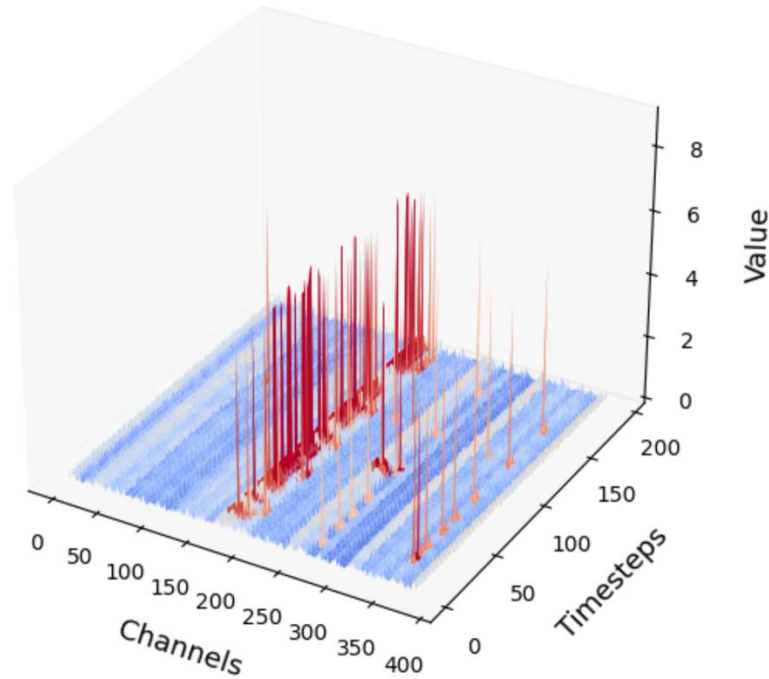
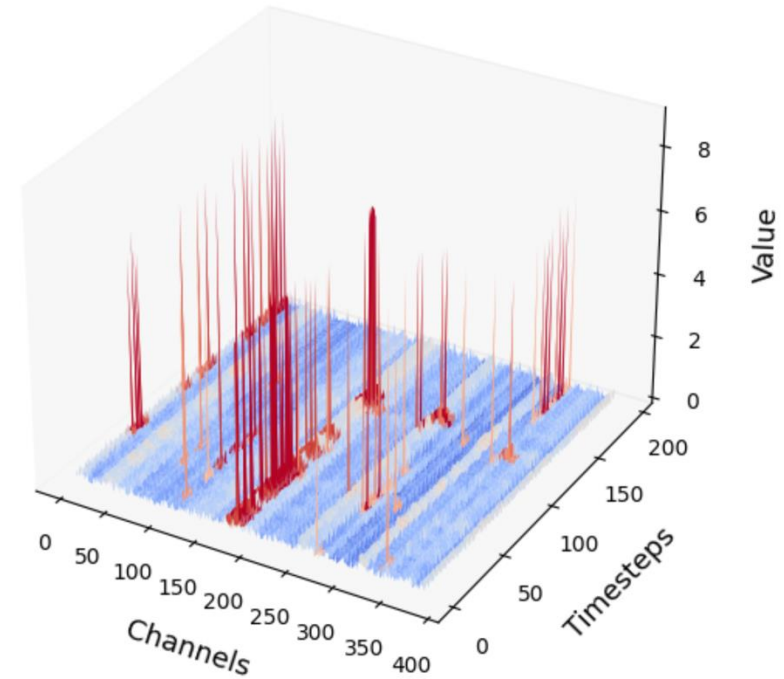
$$\bar{A}(t) = e^{(A\Delta(t))}; B(t) = W_B(u(t)); \bar{B}(t) = B\Delta(t)$$

$$\Delta(t) = S^+(u(t)\Delta(t)_{\text{proj}}); C(t) = (W_C(z(t)))^T$$

VMM activations exhibit dynamic inter-time-step channel variations

Dynamic Activation Variations

ViM-S Layer 3

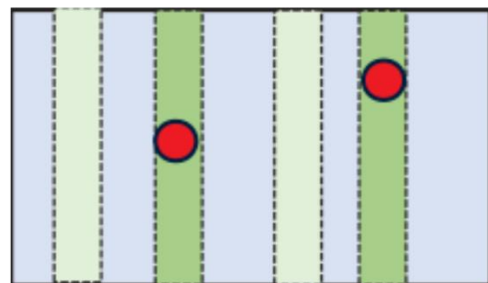
 Δ  A  B

Dynamic variations across different activations too!

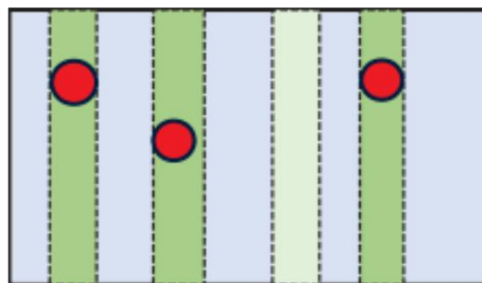
Shortcoming of Static Determination

Overprovisioning for outliers during calibration

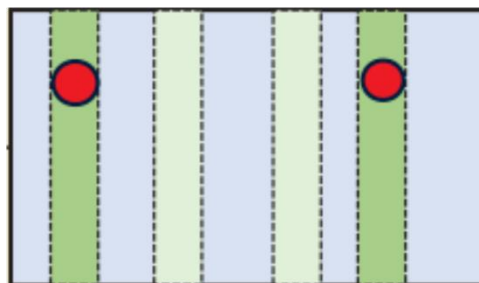
During Calibration



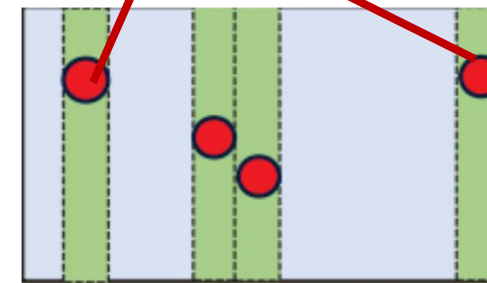
Time step 1



Time step 2

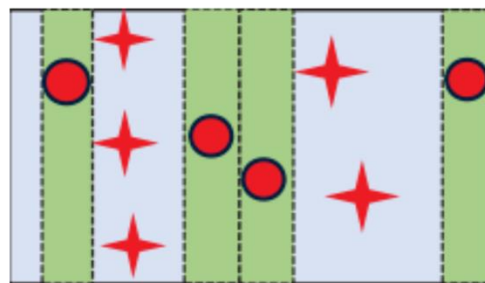


Time step 3

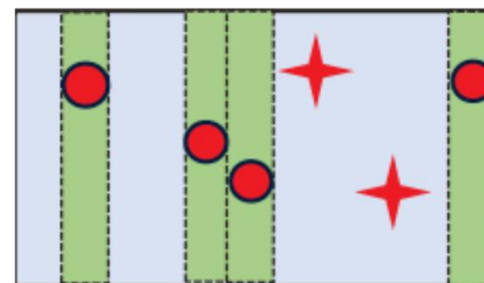


Time step 4

During inference



Time step 1



Time step 2

Misses new outlier channels

OuroMamba-Quant

Input : Activation $X(t) \in \mathbb{R}^{N \times E}$, Static scale $S^I(t)$, Threshold θ , Refresh rate n_{refresh} ,
Outlier list O_{list} , Inlier and outlier bit-precision b_a^I, b_a^O

Output: Quantized activation $X_q(t)$, Updated outlier list O_{list}

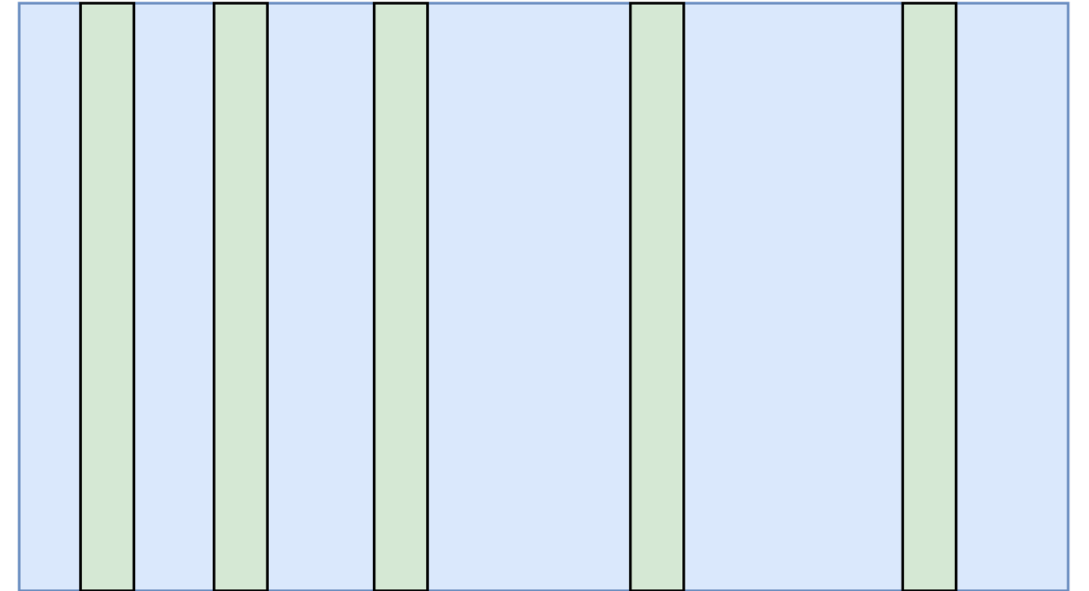
```

if  $t \% n_{\text{refresh}} == 0$  then
  |  $O_{\text{list}} = \{\phi\}$ 
end
 $S^D(t) = \text{ComputeScale}(X(t)[:, c] \forall c \notin O_{\text{list}})$ 
if  $S^D(t) > S^I(t)$  then
  | for each channel  $c$  in  $X(t)$  not in  $O_{\text{list}}$  do
  |   | if  $\max(|X(t)[:, c]|) \geq \theta$  then
  |   |   |  $O_{\text{list}} = O_{\text{list}} \cup \{c\}$ 
  |   | end
  | end
end
 $I(t), O(t) = \text{Separate}(X(t), O_{\text{list}})$ 
 $I_q(t) = \text{InlierQuant}(I(t), S^I(t), b_a^I)$ 
 $O_q(t) = \text{OutlierQuant}(O(t), b_a^O)$ 
 $X_q(t) = \text{Merge}(I_q(t), O_q(t))$ 
return  $X_q(t), O_{\text{list}}$ 

```

Outlier List: $\{\}$

Inliers



Offline: Determine threshold to detect outliers and inlier scale factor

OuroMamba-Quant

Input : Activation $X(t) \in \mathbb{R}^{N \times E}$, Static scale $S^I(t)$, Threshold θ , Refresh rate n_{refresh} ,
Outlier list O_{list} , Inlier and outlier bit-precision b_a^I, b_a^O

Output: Quantized activation $X_q(t)$, Updated outlier list O_{list}

if $t \% n_{\text{refresh}} == 0$ **then**

$O_{\text{list}} = \{\phi\}$

end

$S^D(t) = \text{ComputeScale}(X(t)[:, c] \forall c \notin O_{\text{list}})$

if $S^D(t) > S^I(t)$ **then**

for each channel c **in** $X(t)$ **not in** O_{list} **do**

if $\max(|X(t)[:, c]|) \geq \theta$ **then**

$O_{\text{list}} = O_{\text{list}} \cup \{c\}$

end

end

end

$I(t), O(t) = \text{Separate}(X(t), O_{\text{list}})$

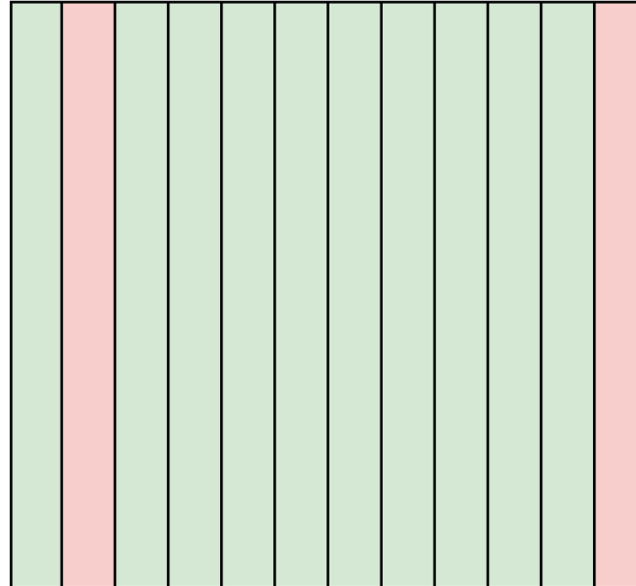
$I_q(t) = \text{InlierQuant}(I(t), S^I(t), b_a^I)$

$O_q(t) = \text{OutlierQuant}(O(t), b_a^O)$

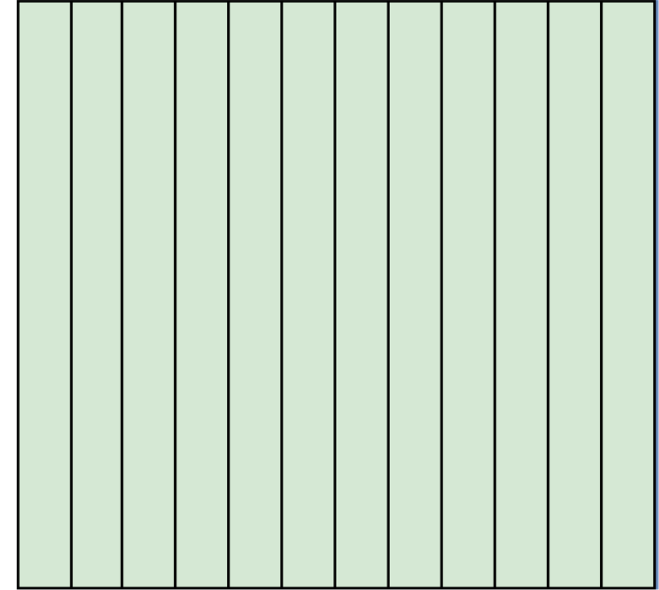
$X_q(t) = \text{Merge}(I_q(t), O_q(t))$

return $X_q(t), O_{\text{list}}$

Outlier List: {1,11}



Outlier List: {}



OuroMamba-Quant

We must now identify which channel is the outlier channel

Input : Activation $X(t) \in \mathbb{R}^{N \times E}$, Static scale $S^I(t)$, Threshold θ , Refresh rate n_{refresh} ,
Outlier list O_{list} , Inlier and outlier bit-precision b_a^I, b_a^O

Output: Quantized activation $X_q(t)$, Updated outlier list O_{list}

if $t \% n_{\text{refresh}} == 0$ **then**

$O_{\text{list}} = \{\phi\}$

end

$S^D(t) = \text{ComputeScale}(X(t)[:, c] \forall c \notin O_{\text{list}})$

if $S^D(t) > S^I(t)$ **then**

for each channel c **in** $X(t)$ **not in** O_{list} **do**

if $\max(|X(t)[:, c]|) \geq \theta$ **then**

$O_{\text{list}} = O_{\text{list}} \cup \{c\}$

end

end

end

$I(t), O(t) = \text{Separate}(X(t), O_{\text{list}})$

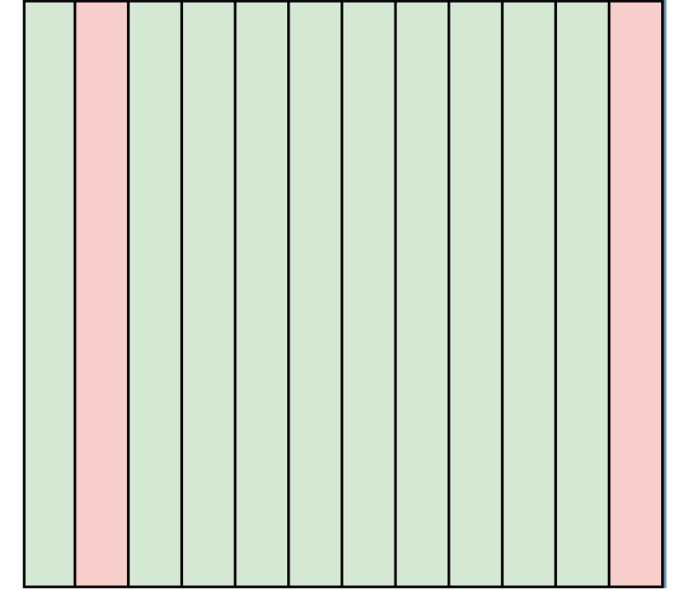
$I_q(t) = \text{InlierQuant}(I(t), S^I(t), b_a^I)$

$O_q(t) = \text{OutlierQuant}(O(t), b_a^O)$

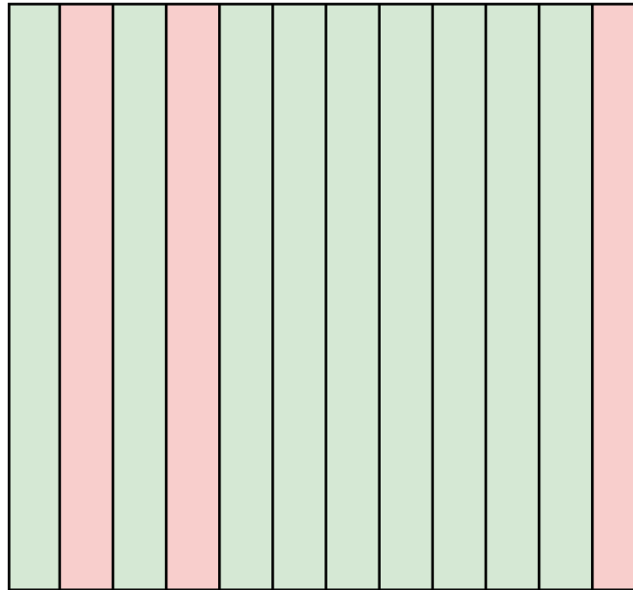
$X_q(t) = \text{Merge}(I_q(t), O_q(t))$

return $X_q(t), O_{\text{list}}$

Outlier List: $\{1, 11\}$



Outlier List: $\{1, 3, 11\}$



OuroMamba-Quant

16

Input : Activation $X(t) \in \mathbb{R}^{N \times E}$, Static scale $S^I(t)$, Threshold θ , Refresh rate n_{refresh} ,
Outlier list O_{list} , Inlier and outlier bit-precision b_a^I, b_a^O

Output: Quantized activation $X_q(t)$, Updated outlier list O_{list}

if $t \% n_{\text{refresh}} == 0$ **then**

$O_{\text{list}} = \{\phi\}$

end

$S^D(t) = \text{ComputeScale}(X(t)[:, c] \forall c \notin O_{\text{list}})$

if $S^D(t) > S^I(t)$ **then**

for each channel c **in** $X(t)$ **not in** O_{list} **do**

if $\max(|X(t)[:, c]|) \geq \theta$ **then**

$O_{\text{list}} = O_{\text{list}} \cup \{c\}$

end

end

end

$I(t), O(t) = \text{Separate}(X(t), O_{\text{list}})$

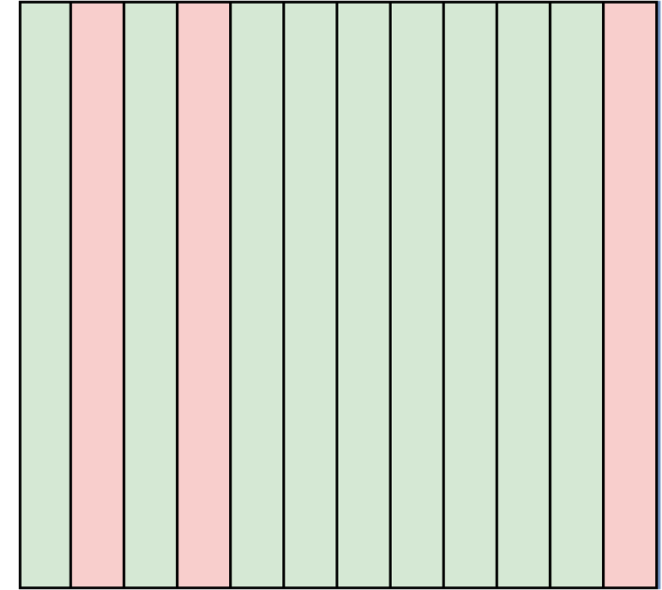
$I_q(t) = \text{InlierQuant}(I(t), S^I(t), b_a^I)$

$O_q(t) = \text{OutlierQuant}(O(t), b_a^O)$

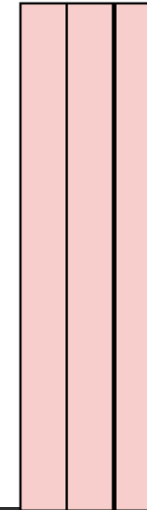
$X_q(t) = \text{Merge}(I_q(t), O_q(t))$

return $X_q(t), O_{\text{list}}$

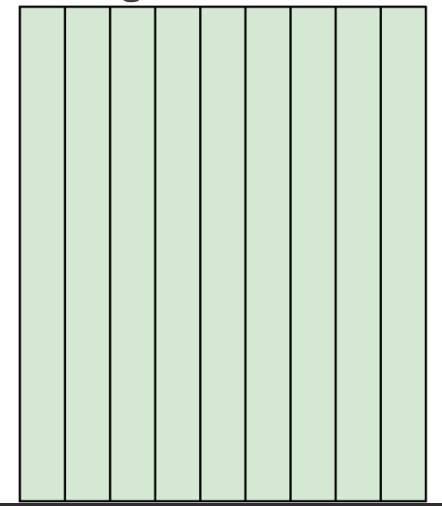
Outlier List: $\{1, 3, 11\}$



Outlier Quantization



Inlier Quantization



OuroMamba-Quant

Input : Activation $X(t) \in \mathbb{R}^{N \times E}$, Static scale $S^I(t)$, Threshold θ , Refresh rate n_{refresh} ,
Outlier list O_{list} , Inlier and outlier bit-precision b_a^I, b_a^O

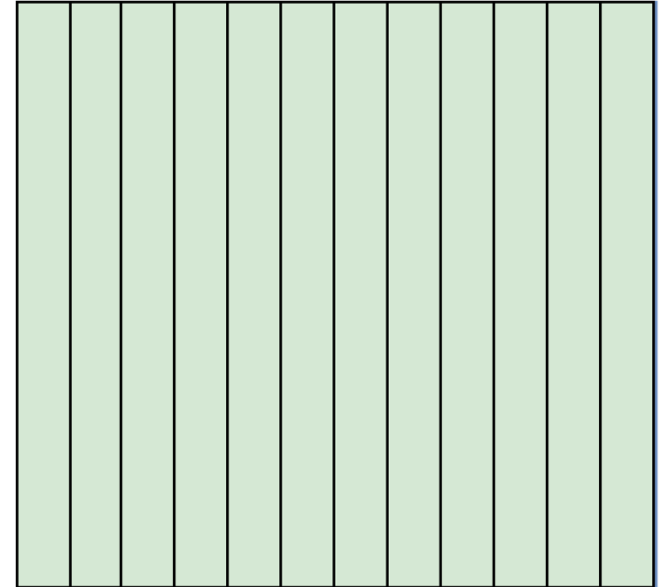
Output: Quantized activation $X_q(t)$, Updated outlier list O_{list}

```

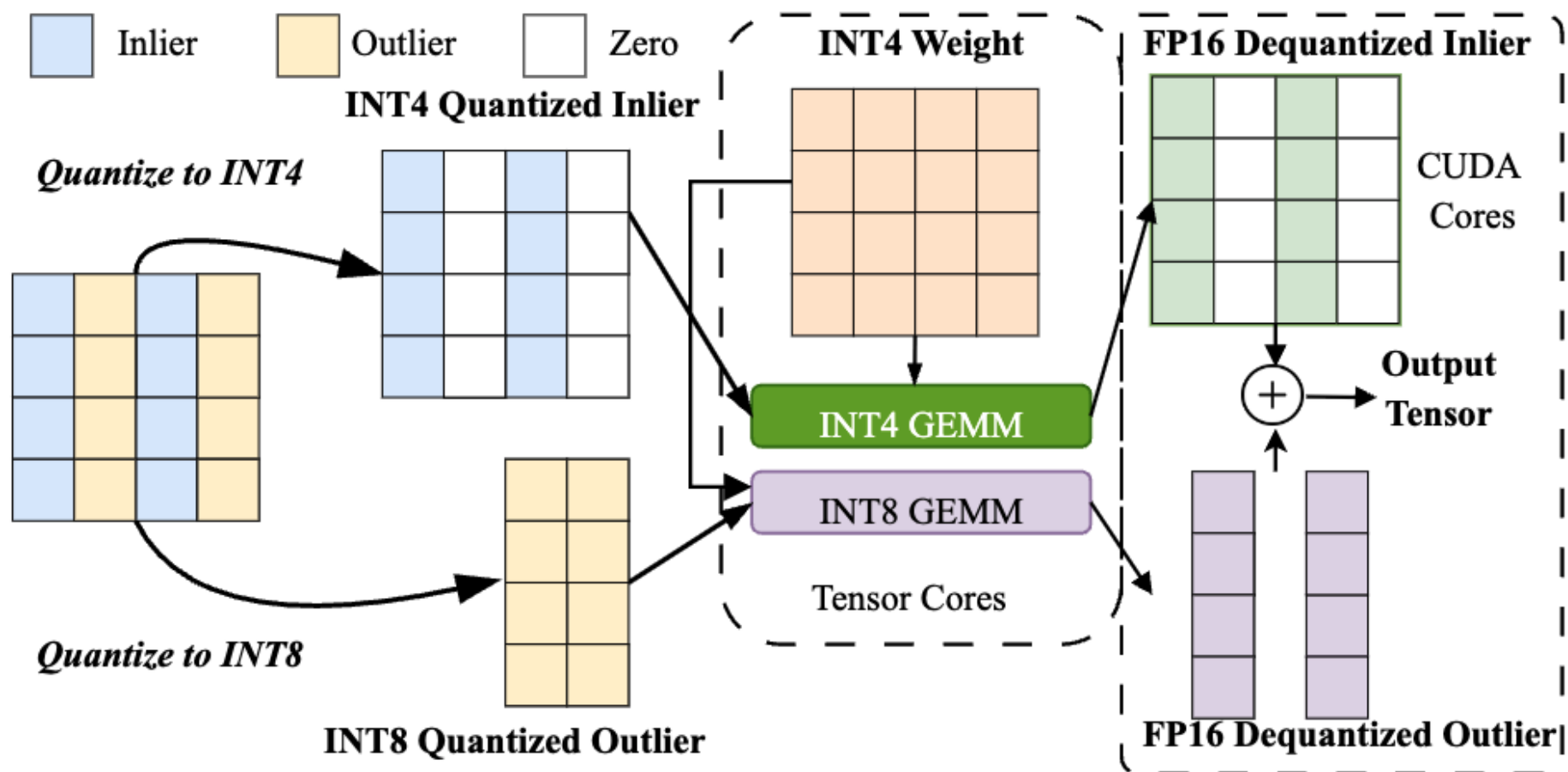
if  $t \% n_{\text{refresh}} == 0$  then
  |  $O_{\text{list}} = \{\phi\}$ 
end
 $S^D(t) = \text{ComputeScale}(X(t)[:, c] \forall c \notin O_{\text{list}})$ 
if  $S^D(t) > S^I(t)$  then
  | for each channel  $c$  in  $X(t)$  not in  $O_{\text{list}}$  do
    | if  $\max(|X(t)[:, c]|) \geq \theta$  then
      | |  $O_{\text{list}} = O_{\text{list}} \cup \{c\}$ 
    | end
  | end
end
 $I(t), O(t) = \text{Separate}(X(t), O_{\text{list}})$ 
 $I_q(t) = \text{InlierQuant}(I(t), S^I(t), b_a^I)$ 
 $O_q(t) = \text{OutlierQuant}(O(t), b_a^O)$ 
 $X_q(t) = \text{Merge}(I_q(t), O_q(t))$ 
return  $X_q(t), O_{\text{list}}$ 

```

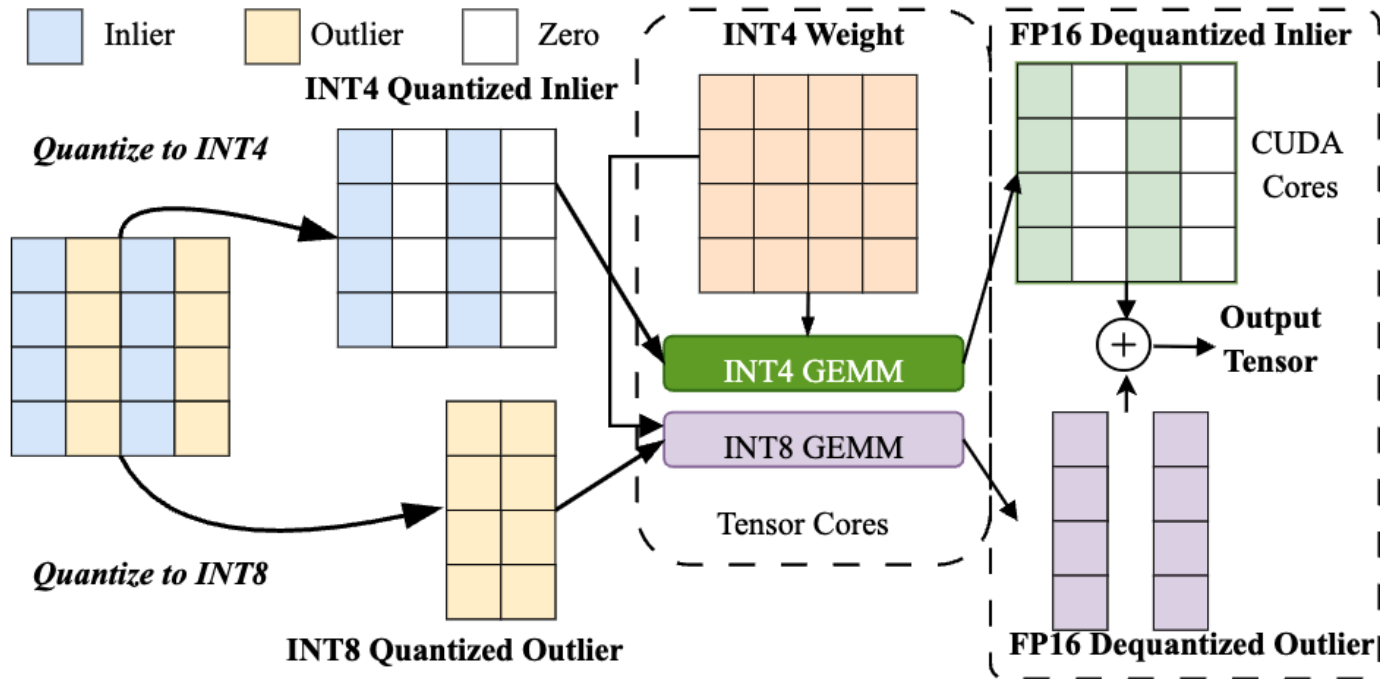
Outlier List: $\{\}$



OuroMamba-Quant: W₄A₄ Hybrid GEMM

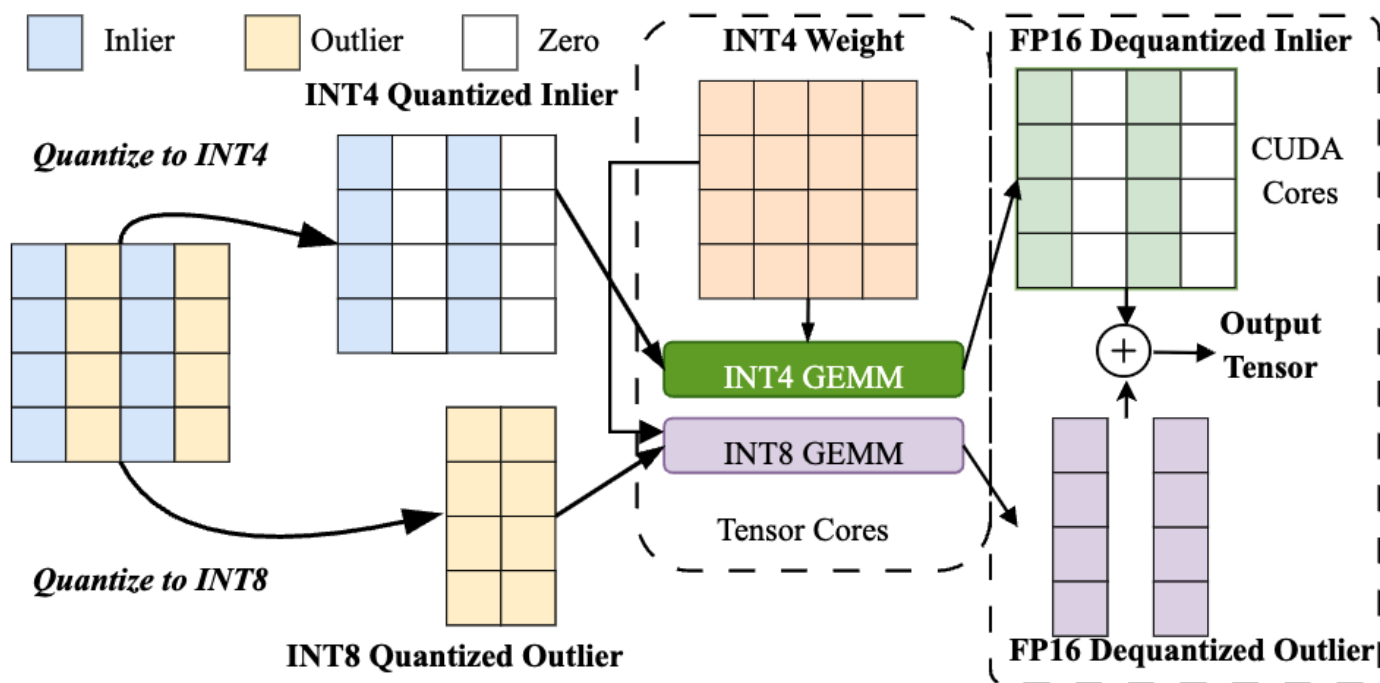


OuroMamba-Quant: W₄A₄ Hybrid GEMM



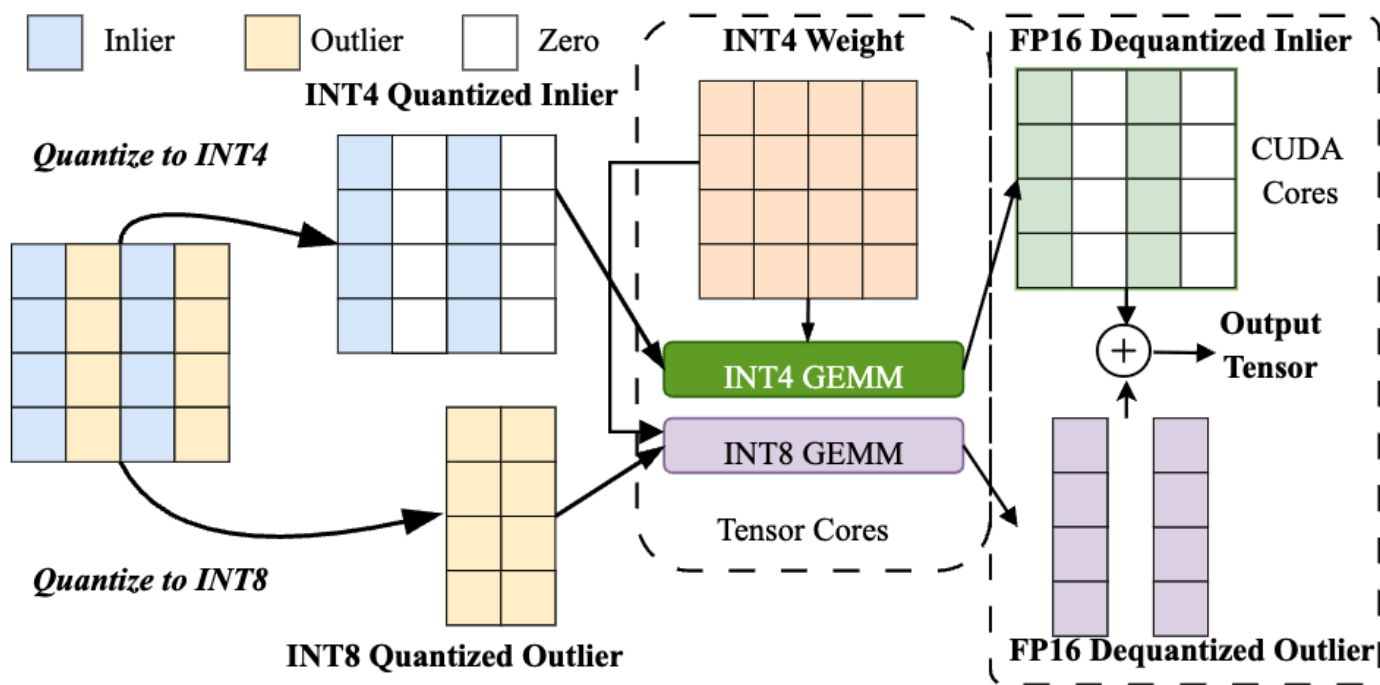
Activations are partitioned across channels and mapped to thread blocks, where each block independently compares its assigned channels against the threshold θ to identify outliers and updates the O_{list}

OuroMamba-Quant: W₄A₄ Hybrid GEMM



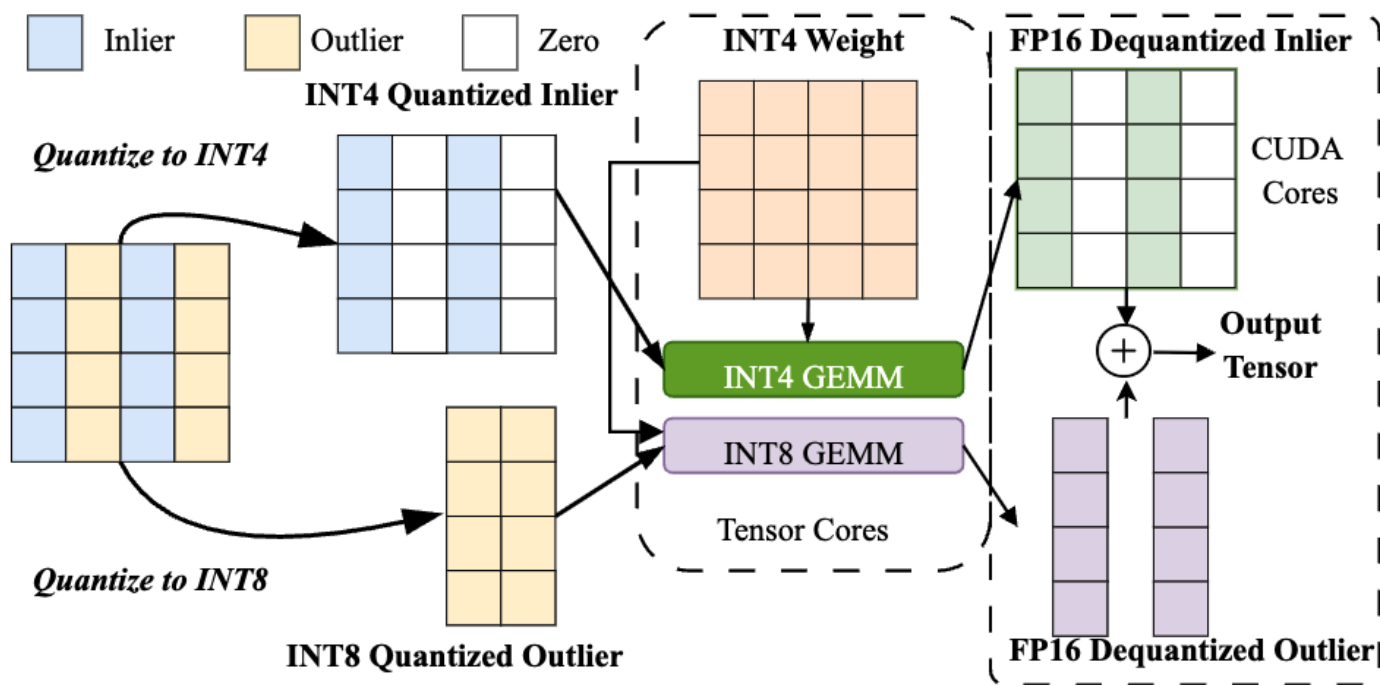
We pack two consecutive 4-bit inlier activations into one byte, with outlier positions set to zero, and leverage the **INT4 tensor cores for inlier GEMM**.

OuroMamba-Quant: W₄A₄ Hybrid GEMM



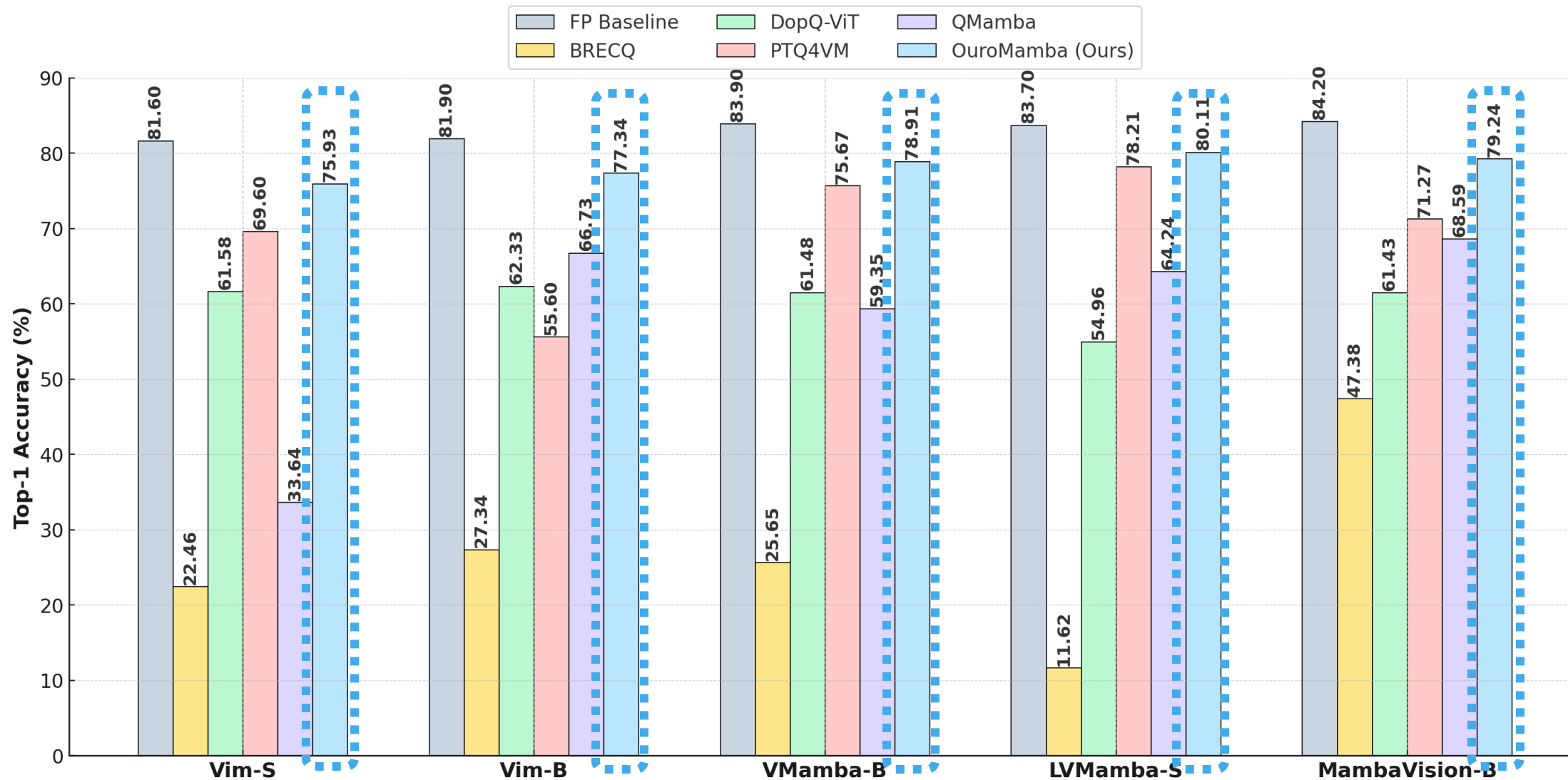
The outliers execute GEMM on **INT8 tensor cores**.

OuroMamba-Quant: W₄A₄ Hybrid GEMM



The outputs from the inlier and outlier GEMMs are dequantized to FP16 and summed together.

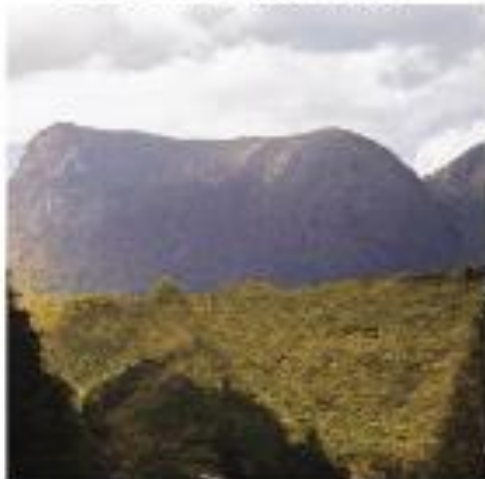
Quantization Results: W₄A₄



OuroMamba: Outliers in 8-bit

Quantization Results for Diffusion Models

FP16 Baseline



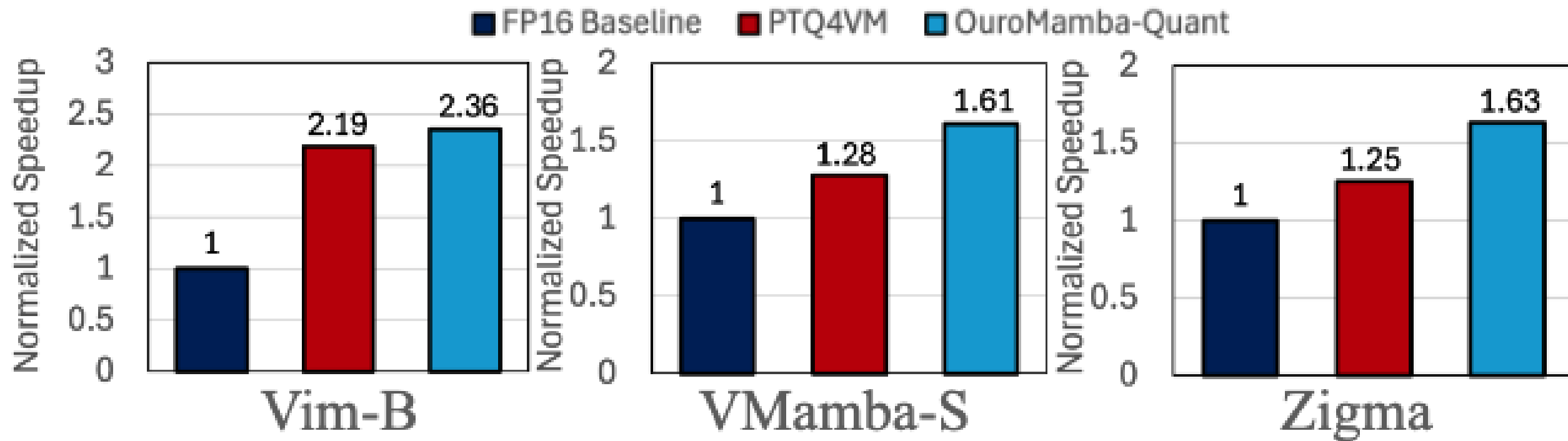
W4A4 QMamba



W4A4 OuroMamba



End-to-End Latency Comparison

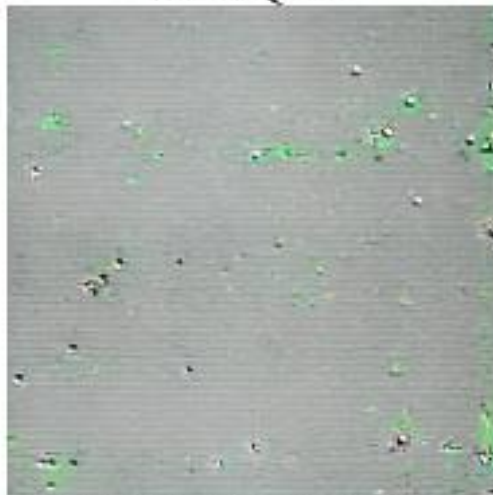


Generalization to ViTs

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