



**ICLR**

# SSDI8: Accurate and Efficient 8-Bit Quantization for State Space Duality

Hyunwoo Kim<sup>1</sup>, Byoungchan Ko<sup>2</sup>, Minseok Kang<sup>2</sup>, Minwoo Kim<sup>2</sup>, Dongin Lee<sup>3</sup>, Jaehoon Lee<sup>4</sup>, Sungroh Yoon<sup>3,4,5†</sup>, Dahuin Jung<sup>1†</sup>

<sup>1</sup>Department of Artificial Intelligence, Chung–Ang University

<sup>2</sup>School of Computer Science and Engineering, Soongsil University

<sup>3</sup>Department of Electrical and Computer Engineering, Seoul National University

<sup>4</sup>Interdisciplinary Program in Artificial Intelligence, Seoul National University

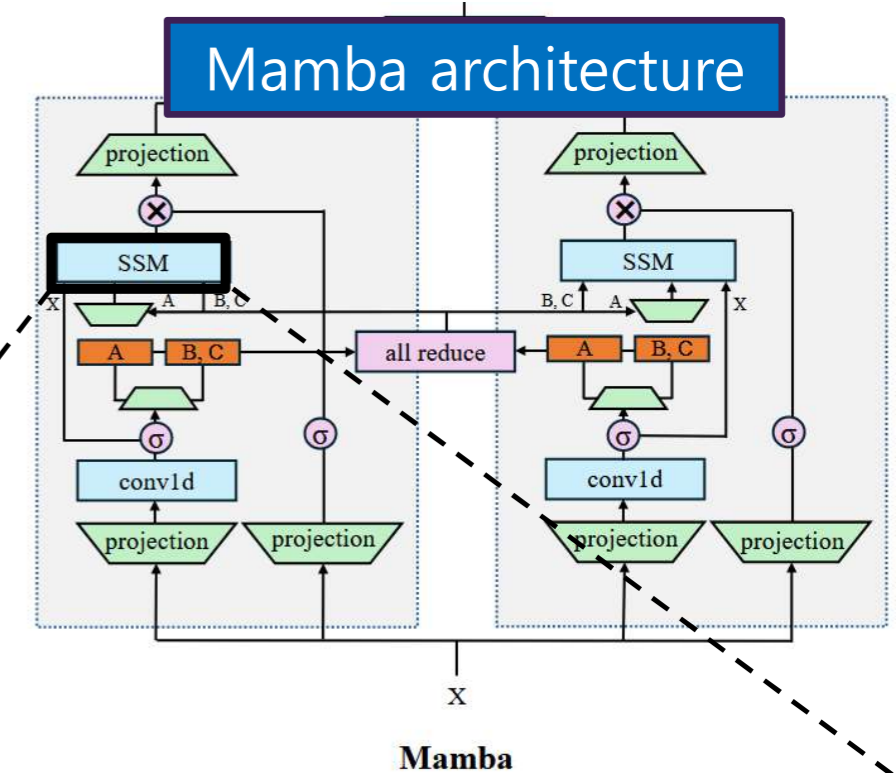
<sup>5</sup>AIIS, ASRI, INMC, and ISRC, Seoul National University

\*Equal contribution

†Corresponding Authors

# Background & Motivation

- Transformer-based LLMs incur computation and memory costs that scale quadratically with sequence length.
- Mamba, an SSM-based architecture, achieves performance comparable to or better than recent architectures.
- However, SSM is difficult to parallelize on modern accelerators and challenging to scale to larger model sizes.



## State Space Model Equation

Continuous

$$h'(t) = Ah(t) + Bx(t)$$

$$y(t) = Ch(t)$$



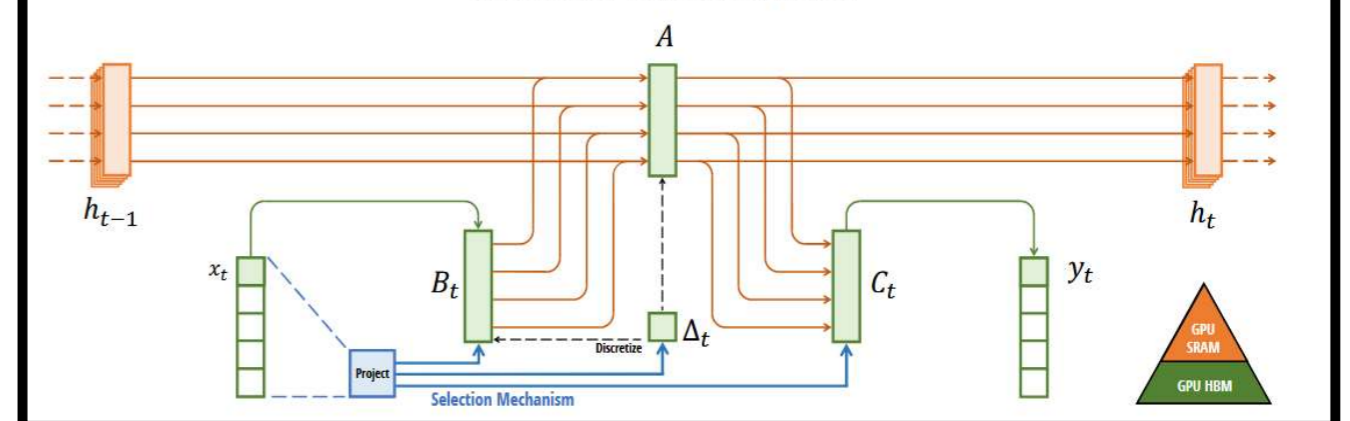
Discretized (ZOH)

$$h_t = \bar{A}h_{t-1} + \bar{B}x_t$$

$$y_t = Ch_t$$

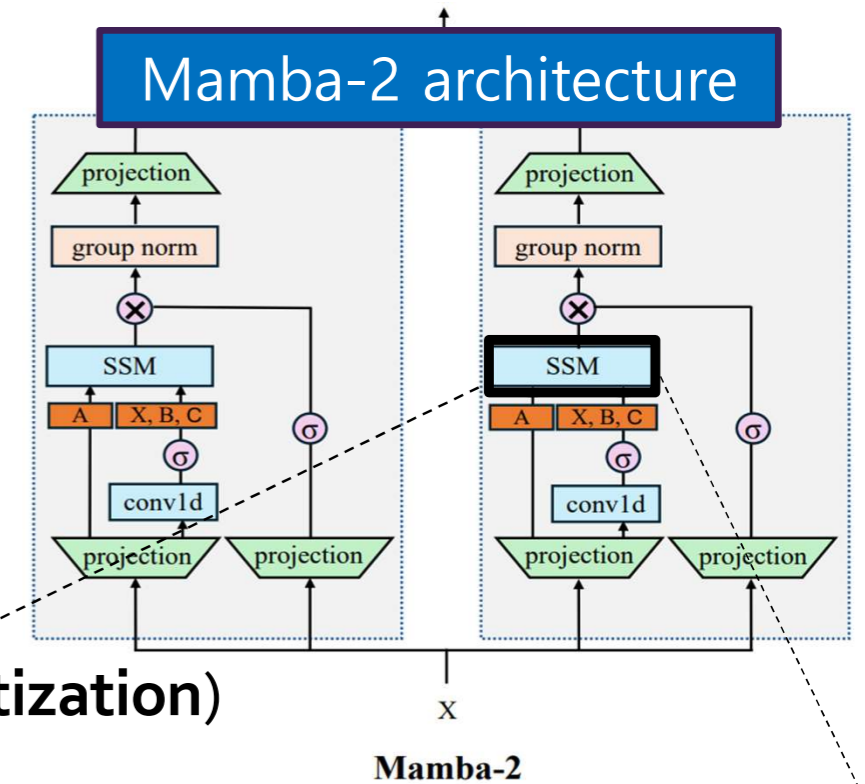
$$\bar{A} = \exp(\Delta A), \quad \bar{B} = (\Delta A)^{-1}(\exp(\Delta A) - I) \cdot \Delta B$$

## Selective State Space Model



# Background & Motivation

- Mamba-2 introduces the **Structured State Space Duality(SSD)**, a hybrid design that integrates recurrent model with attention mode.
- It employs a dual representation that improves GEMM utilization, yielding higher throughput on GPUs and TPUs.
- Mamba-2 scales effectively to over 8B parameters, but the resulting memory and latency overhead calls for efficient compression. (**SSD-Quantization**)



## State Space Model Equation

Continuous

$$h'(t) = Ah(t) + Bx(t)$$

$$y(t) = Ch(t)$$



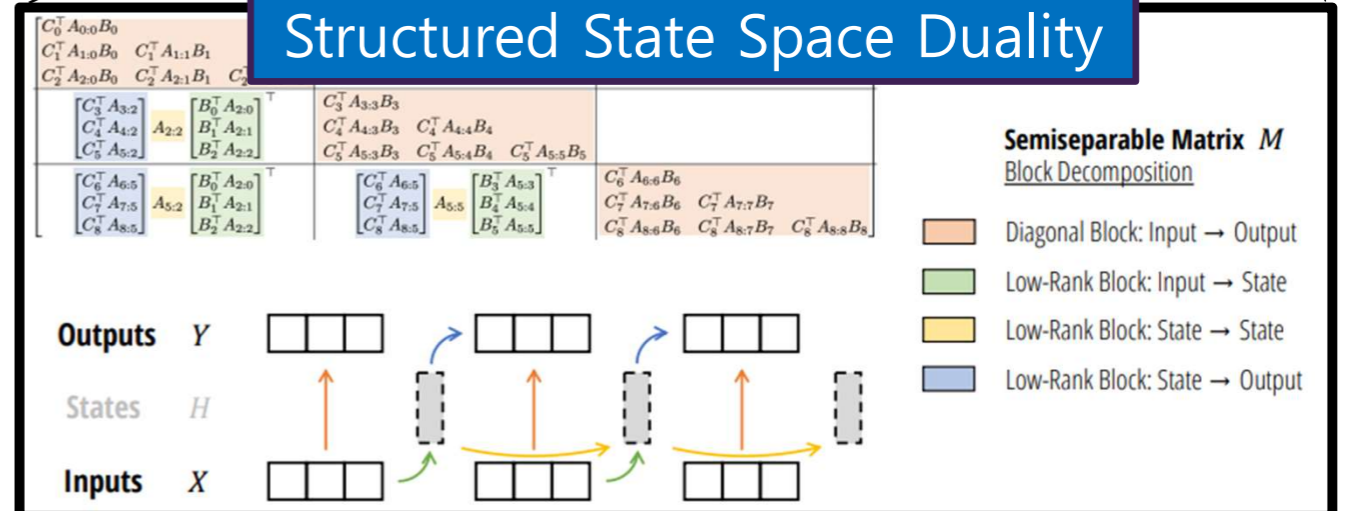
Discretized (ZOH)

$$h_t = \bar{A}h_{t-1} + \bar{B}x_t$$

$$y_t = Ch_t$$

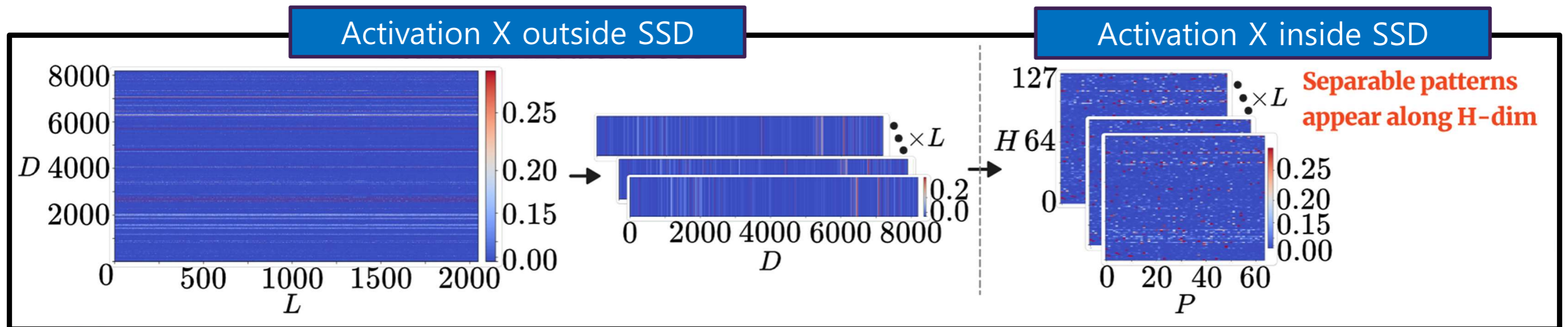
$$\bar{A} = \exp(\Delta A), \quad \bar{B} = (\Delta A)^{-1}(\exp(\Delta A) - I) \cdot \Delta B$$

## Structured State Space Duality



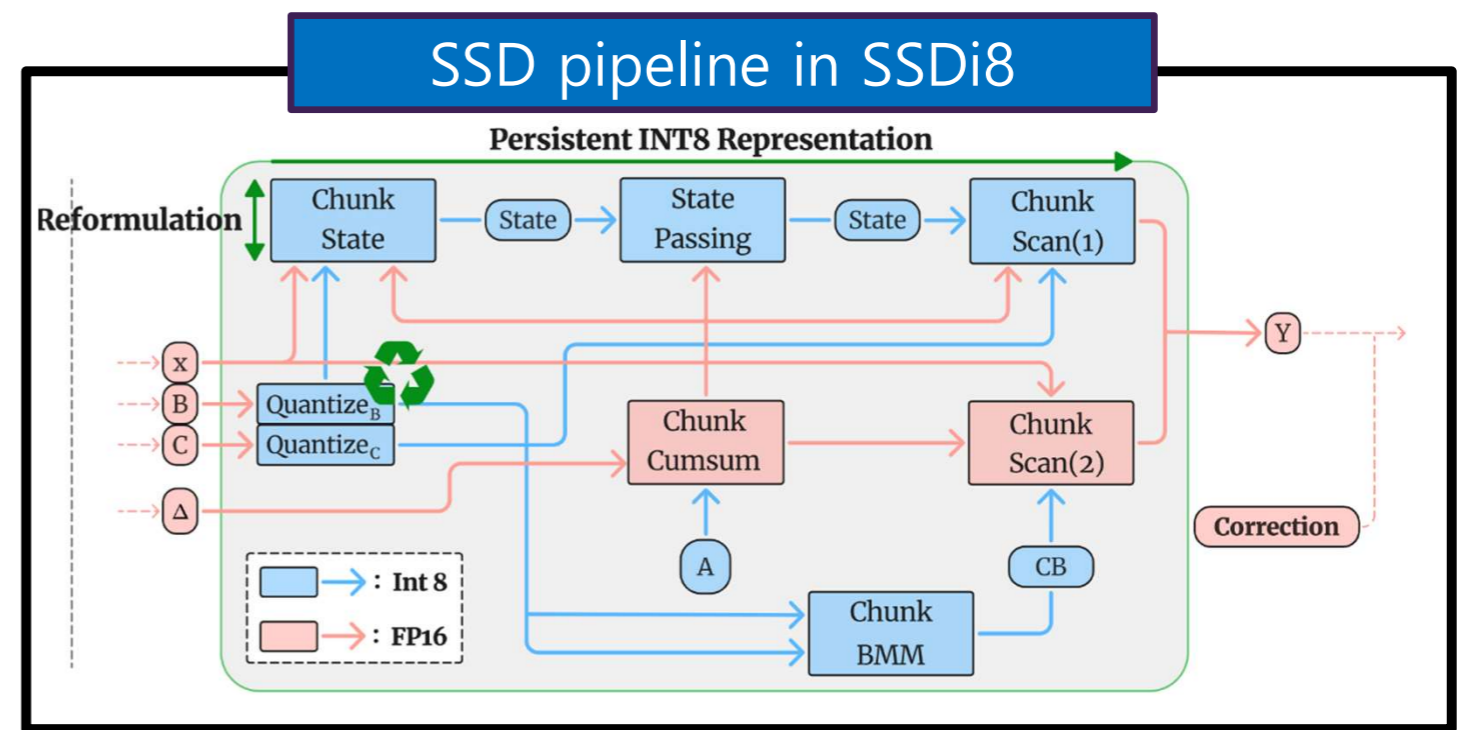
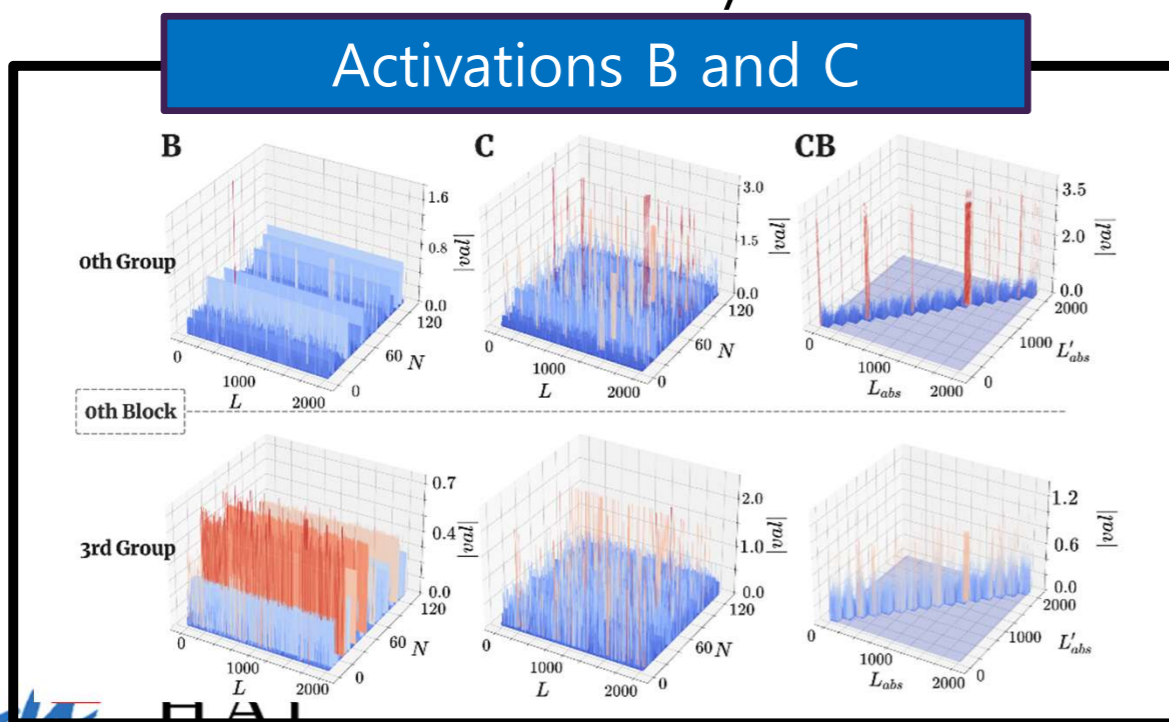
# Problems

- SSD involves external-to-internal dimension changes that must be considered to preserve quantization performance.
- Activations are repeatedly reused across multiple modules and exhibit distinct outlier characteristics.
- Element-wise multiplications are tightly intertwined with matrix multiplications, complicating quantization.



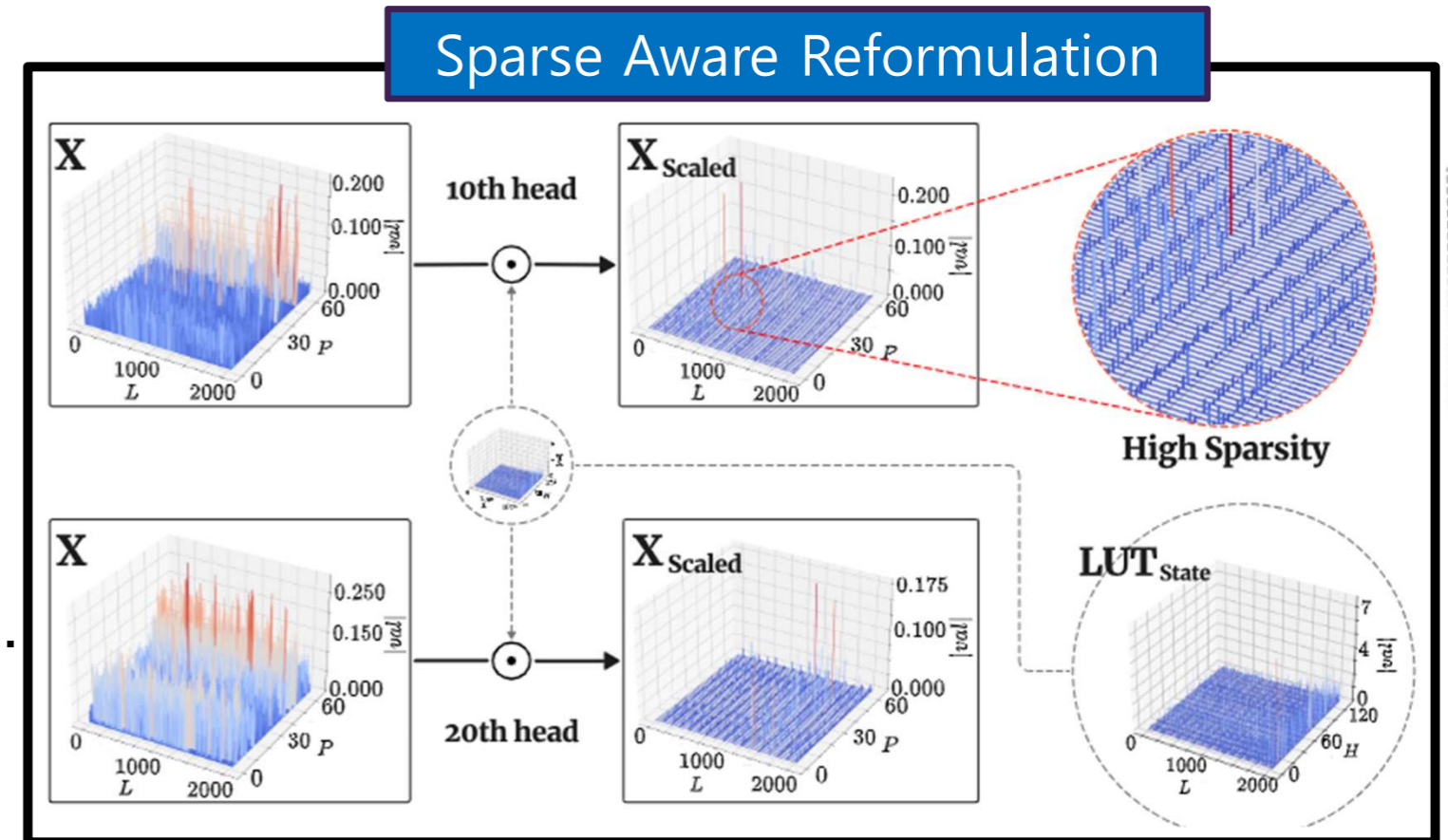
# SSDi8 : Persistent INT8 Representation

- **Our Goal:** Establish a **Persistent INT8 path Representation (PIR)** within SSD to minimize latency while preserving accuracy.
- **Quantize once and reuse:** Quantize reused activation B and C once and reuse them across SSD submodules to reduce memory traffic and quantization overhead.
- **Mean correction:** Apply a per-channel mean correction after SSD to mitigate accumulated quantization error with minimal latency overhead.



# SSDi8 : Sparse Aware Reformulation

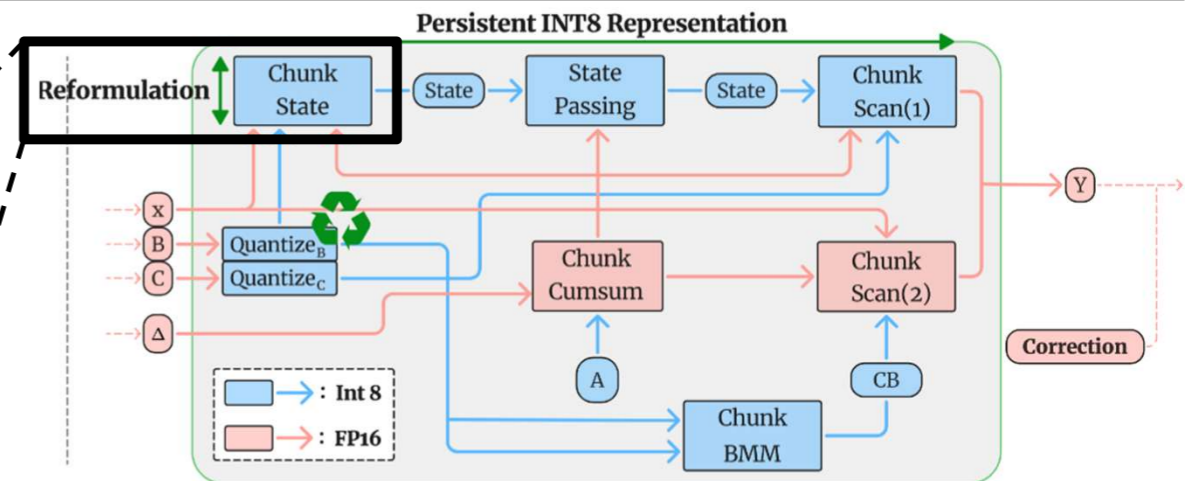
- Multiplying reused INT8  $B$  by FP16  $LUT_{state}$  incurs additional overhead.
- **Sparse Aware Reformulation(SAR)** preserves INT8 execution path and avoids the overhead.
- Although  $X_{scaled}$  is highly sparse, its quantization does not increase error and is theoretically justified.
- As a result, PIR is completed, enabling the SSD internal path to remain in INT8.



## SAR Equation

$$State = X \times (B \odot LUT_{state}) \longrightarrow Q(State) = Q(X_{scaled}) \times Q(B)$$

$$X_{scaled} = LUT_{state} \odot X$$



# Experiments

## Zero-shot task results

Model	Size	Methods	Bitwidth	LA	HS	PIQA	Arc-E	Arc-C	WG	Avg.
Mamba-2	1.3B	-	FP16	65.6%	59.9%	73.3%	64.1%	33.3%	60.8%	59.5%
		Quamba	W8A8	49.8%	58.5%	71.2%	61.9%	32.1%	58.1%	55.2%
		Quamba2	W8A8	62.0%	59.2%	72.5%	63.4%	32.7%	60.0%	58.3%
			W4A8	61.0%	58.8%	72.4%	62.7%	32.6%	59.1%	57.7%
		SSDi8 (Ours)	W8A8	<b>64.7%</b>	<b>59.7%</b>	<b>72.7%</b>	<b>64.0%</b>	<b>32.8%</b>	<b>60.9%</b>	<b>59.1%</b>
			W4A8	<b>63.6%</b>	<b>59.2%</b>	<b>72.7%</b>	<b>63.5%</b>	<b>33.5%</b>	<b>60.4%</b>	<b>58.8%</b>
	2.7B	-	FP16	69.5%	66.6%	76.4%	69.5%	36.4%	64.2%	63.8%
		Quamba	W8A8	52.4%	60.4%	71.6%	62.9%	33.7%	58.0%	56.5%
		Quamba2	W8A8	66.1%	65.5%	74.4%	68.4%	<b>37.1%</b>	<b>63.7%</b>	62.5%
			W4A8	65.6%	65.1%	74.7%	68.1%	<b>36.1%</b>	62.8%	62.1%
		SSDi8 (Ours)	W8A8	<b>68.3%</b>	<b>66.2%</b>	<b>75.6%</b>	<b>69.0%</b>	36.8%	63.4%	<b>63.2%</b>
			W4A8	<b>67.4%</b>	<b>65.3%</b>	<b>75.6%</b>	<b>68.9%</b>	35.2%	<b>63.5%</b>	<b>62.6%</b>
	8B	-	FP16	70.9%	77.7%	79.7%	76.0%	48.0%	72.0%	70.7%
		Quamba	W8A8	54.0%	74.6%	77.1%	73.5%	44.2%	65.5%	64.8%
Quamba2		W8A8	69.8%	<b>77.8%</b>	79.1%	<b>75.9%</b>	46.9%	69.0%	69.8%	
		W4A8	68.8%	<b>77.1%</b>	79.1%	75.0%	46.0%	68.7%	69.1%	
SSDi8 (Ours)		W8A8	<b>70.4%</b>	77.2%	<b>79.6%</b>	75.5%	<b>47.2%</b>	<b>71.2%</b>	<b>70.2%</b>	
		W4A8	<b>69.9%</b>	76.5%	<b>79.1%</b>	<b>75.4%</b>	<b>46.2%</b>	<b>70.6%</b>	<b>69.6%</b>	

## Perplexity results

Methods	Bitwidth	Wikitext2 Perplexity (↓)		
		1.3B	2.7B	8B
-	FP16	10.42	9.06	7.25
HAD	W8A8	11.31	11.42	8.57
	W4A8	11.63	11.85	8.79
Quamba2	W8A8	10.80	9.32	7.79
	W4A8	11.08	9.54	7.94
SSDi8 (Ours)	W8A8	<b>10.63</b>	<b>9.22</b>	<b>7.49</b>
	W4A8	<b>10.92</b>	<b>9.43</b>	<b>7.62</b>

- SSD8 preserves strong zero-shot accuracy under quantization.
- SSDi8 achieve perplexity **close to FP16**, preserving linguistic fluency and generalization.

# Experiments

- SSDi8 achieves substantial latency reduction, reaching up to **1.47x speedup over FP16**.
- With Persistent INT8 Representation, *ChunkScan* and *StatePassing* achieve up to **1.77x** and **2.25x speedup over FP16**, respectively.
- SSDi8 consistently reduces SSD latency on Orin NX 16G across sequence length, demonstrating robustness in resource-constrained environments.

## SSD latency on Orin NX 16G

GPU	Orin NX 16G			
	W4A8		W8A8	
Method	Quamba2	SSDi8	Quamba2	SSDi8
$L = 256$	55.30	44.71	51.03	41.30
$L = 512$	76.10	68.00	70.95	60.49
$L = 1024$	134.40	127.51	139.10	114.36
$L = 2048$	262.90	240.54	249.29	217.69

