

**DENSO**  
Crafting the Core



CVPR  
JUNE 3-7, 2026



DENVER  
COLORADO

# Rethinking Asymmetric Quantization: Hidden Symmetry in Vision Model Weights

Masafumi Mori<sup>1</sup>, Shinya Gongyo<sup>2</sup>, Mitsuru Ambai<sup>2</sup>

<sup>1</sup>DENSO CORPORATION, <sup>2</sup>DENSO IT Laboratory

CVPR Poster: Day3 (Morning)



# Background: PTQ for Efficient Vision Inference

- Many applications require both high accuracy and real-time inference on edge devices



- Efficient deployment of vision models has become essential
  - **Post-Training Quantization (PTQ)** enables efficient inference without retraining
- Recent PTQ methods commonly use **Asymmetric Quantization (AsymQ)** for weights, because it improves low-bit accuracy

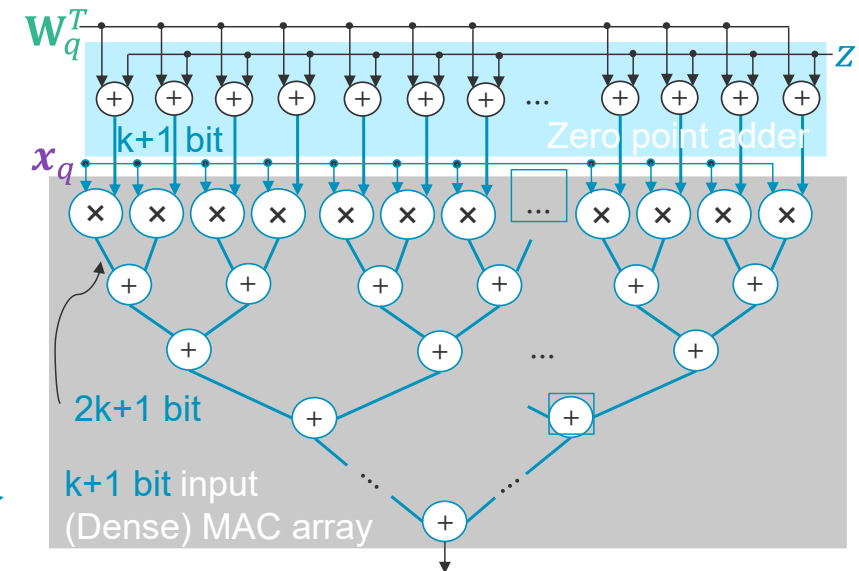
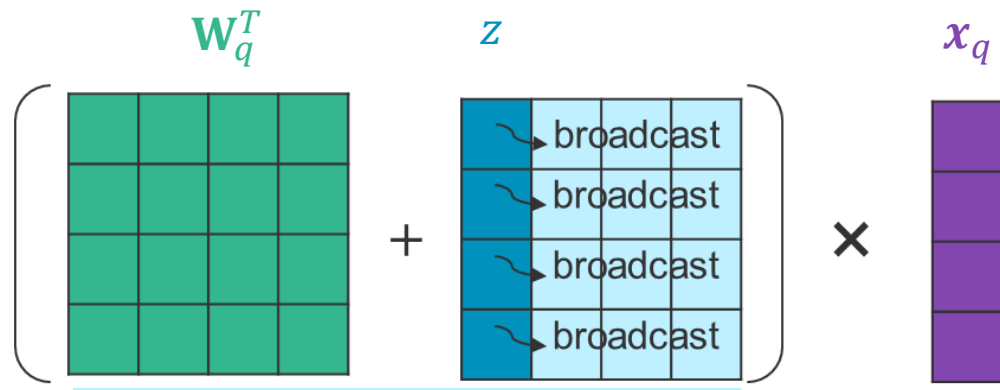
$$\hat{W} \equiv Q(W, \Delta, z, k) = \Delta(W_q - z)$$

$$W_q \equiv \text{clip}(\lceil W \rceil / \Delta + z, -2^{k-1}, 2^{k-1} - 1)$$

**PTQ is attracting increasing attention in both industry and academia**

# Motivation

- Asymmetric Quantization (AsymQ)
  - Low-bit symmetric  $W_q$  + zero-point  $z$  addition
  - Commonly adopted for better accuracy than SymQ ( $z = 0$ )



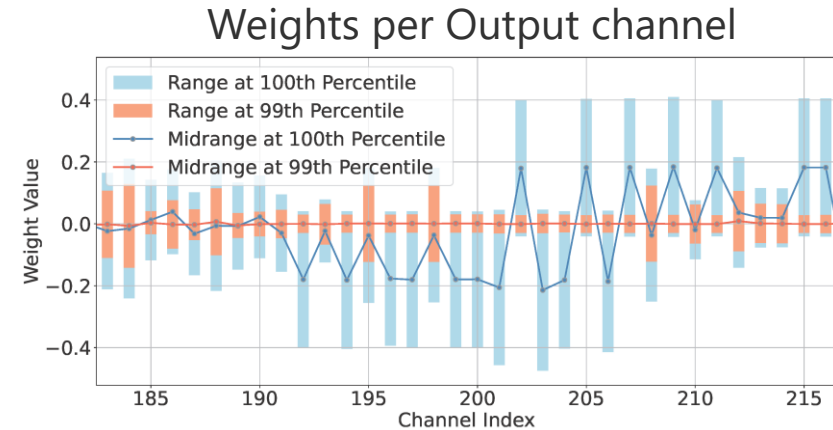
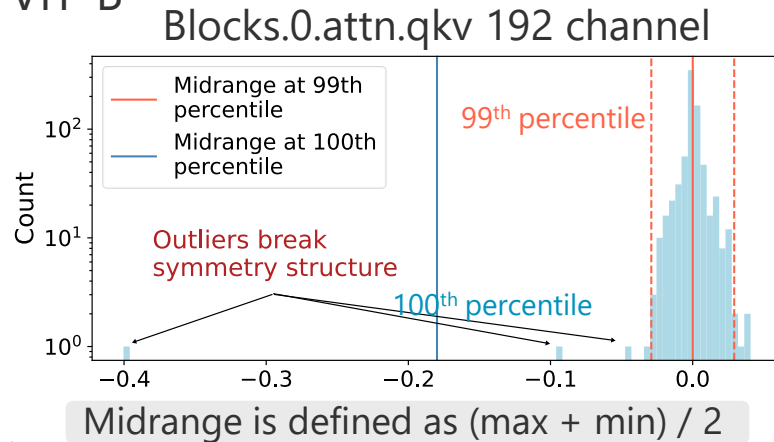
Zero-point addition increases multiplier bit-width → higher hardware cost

## Do we really need AsymQ?

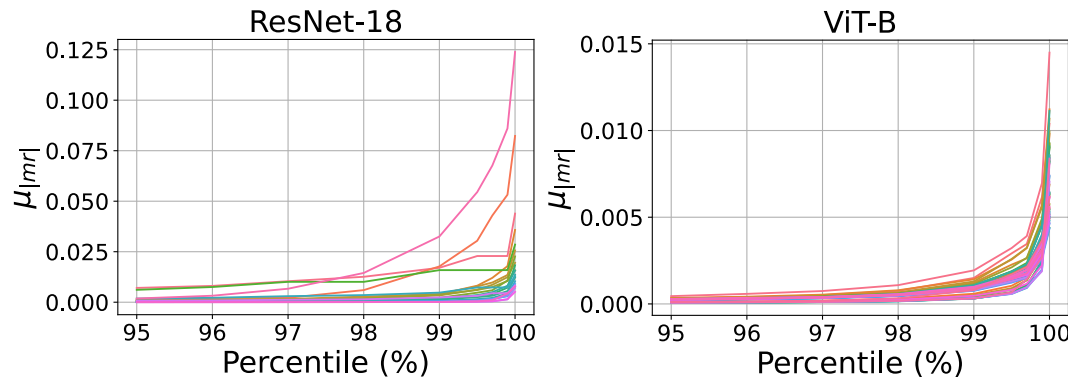


# Observation

- We analyze the distribution of network weights
  - Weights of ViT-B



- CNN and ViT



➔ Weight become **symmetric** across CNN and Vision Transformer

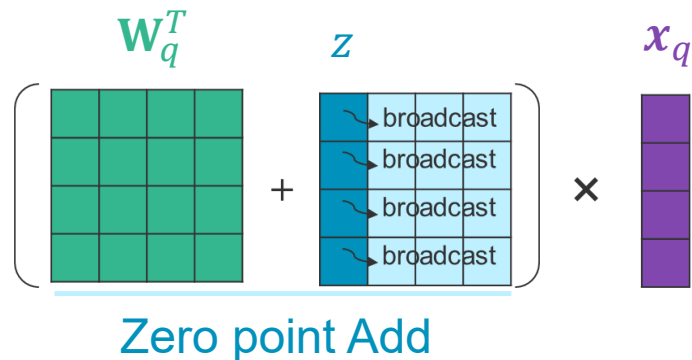
$\mu_{|mr|}$  represents the mean of the absolute midrange averaged over output channels in each layer.

## Hidden Symmetry in Vision Model weights

# Method: DASQ (Dense and Additive Sparse Quantization)

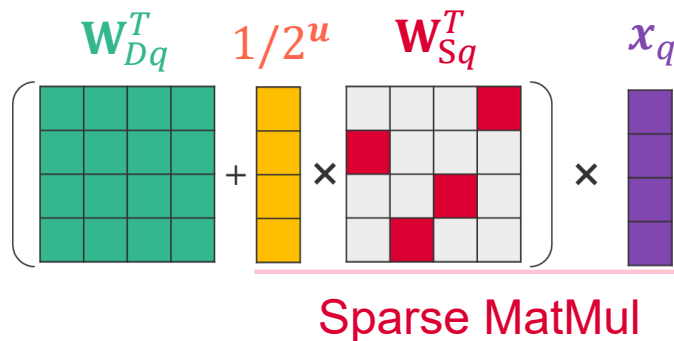
## Weight Representation

**AsymQ**  
(Baseline)



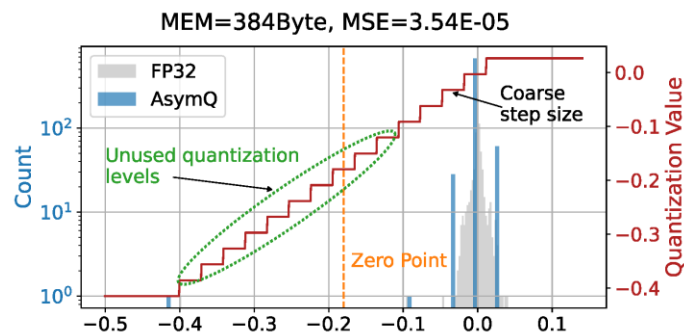
Asymmetry handled by **zero-point**

**DASQ**  
(Ours)

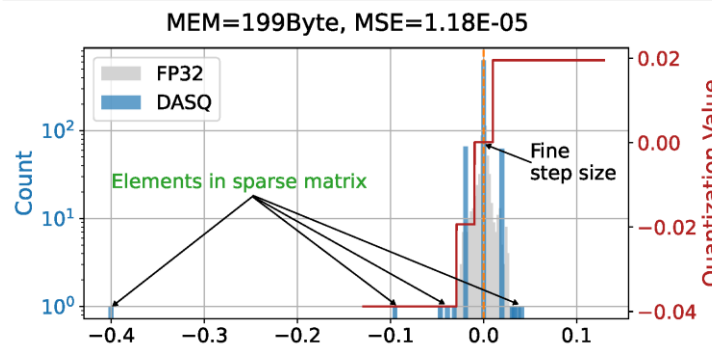


Asymmetry handled by **sparse outliers**

## Quantization Efficiency

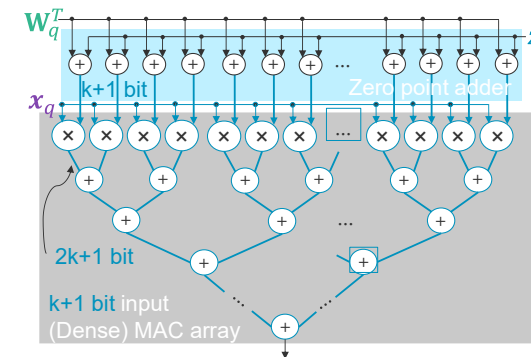


**Larger Quantization Error**

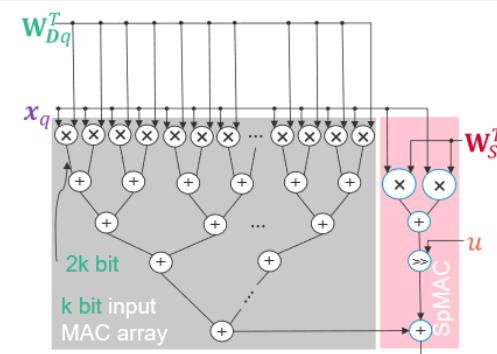


**Smaller Quantization Error**

## Hardware Efficiency



**Larger Multiplier tree**



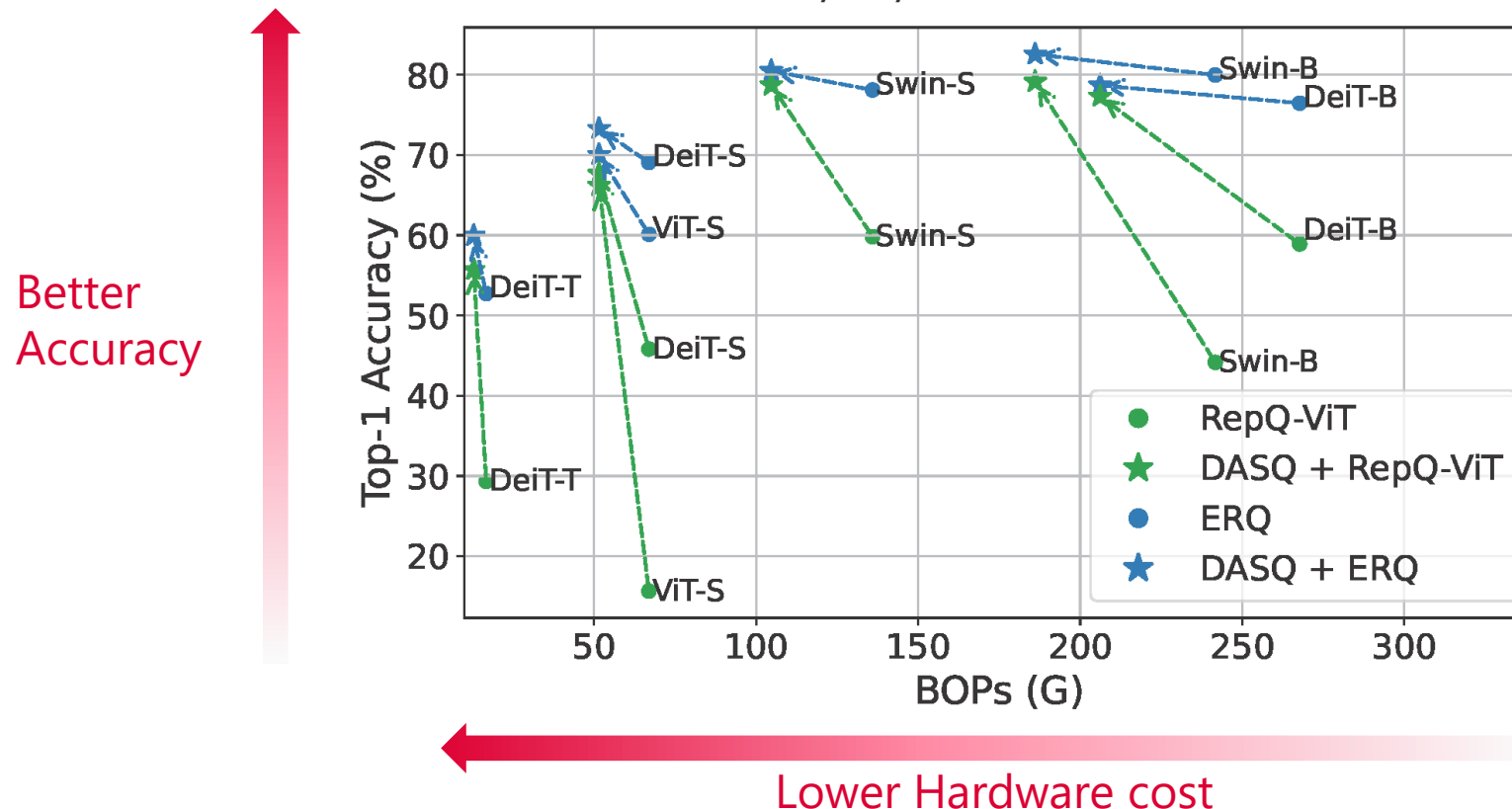
**Smaller Multiplier tree**

# Experimental Results

## ➤ ImageNet Evaluation

➤ Weight Quantization has been changed **AsymQ** to **DASQ**

W3/A3/SW4-98.0%

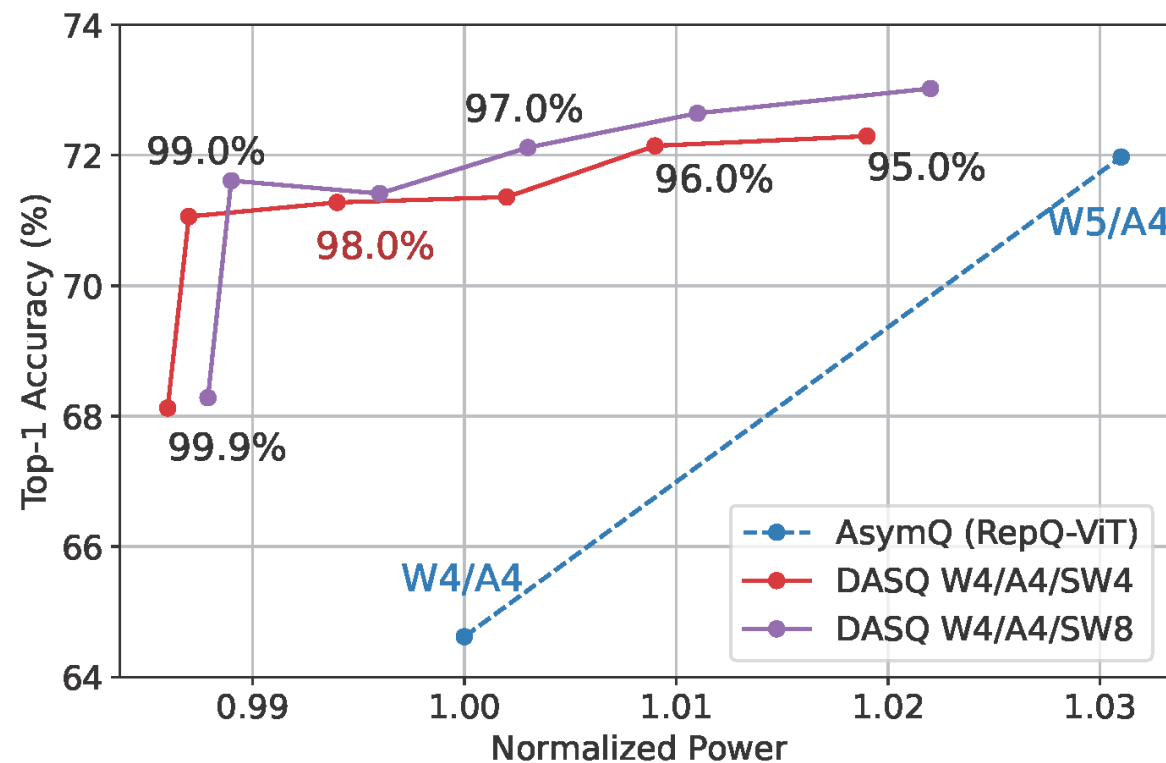


**DASQ achieves higher accuracy with lower hardware cost**

## FPGA Deployment



## Accuracy-Power Tradeoff: ViT-S



**DASQ outperforms AsymQ on FPGA**



# Conclusion

- Rethinking Asymmetric Quantization (AsymQ) for Vision Model Weights
  - ✓ Observation: Hidden symmetry in vision model weights
    - Vision model weights become symmetric after removing sparse outliers
  - ✓ Method: DASQ
    - Asymmetry is represented by sparse outliers, not zero-point additions
  - ✓ Result
    - Higher accuracy and better hardware efficiency than AsymQ

**DASQ enables accurate and hardware-friendly inference**



Thank you

CVPR Poster: Day3 (Morning)

***DENSO***  
Crafting the Core



**DENSO  
IT LAB**

