

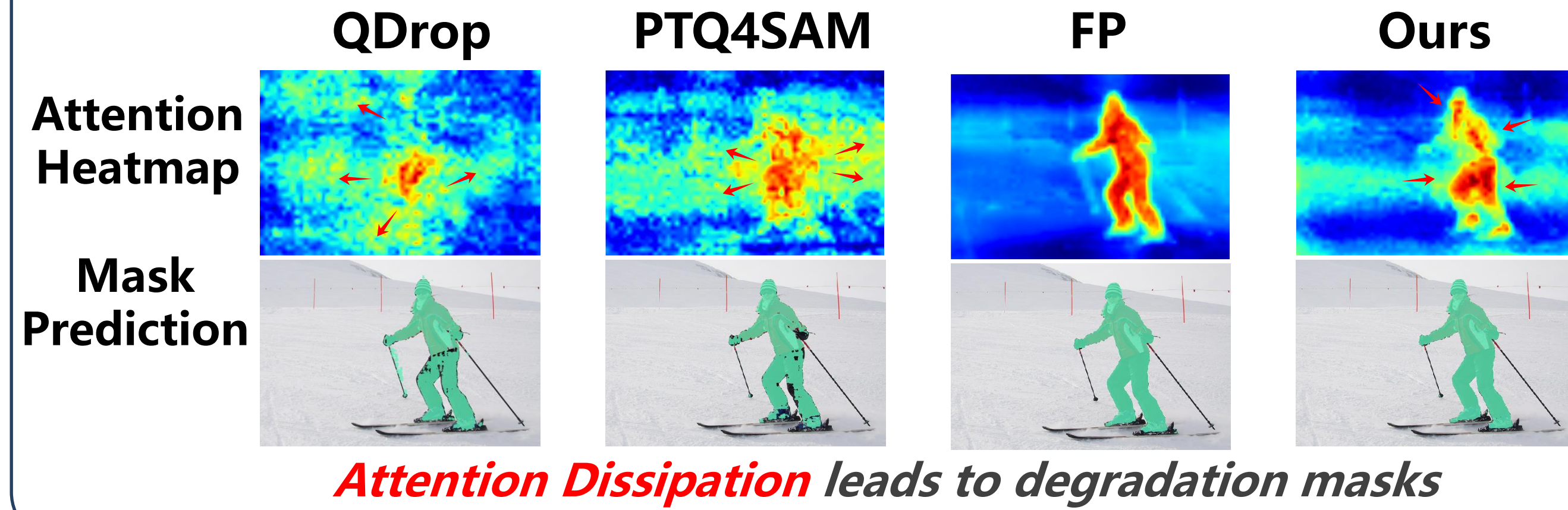


## Motivation

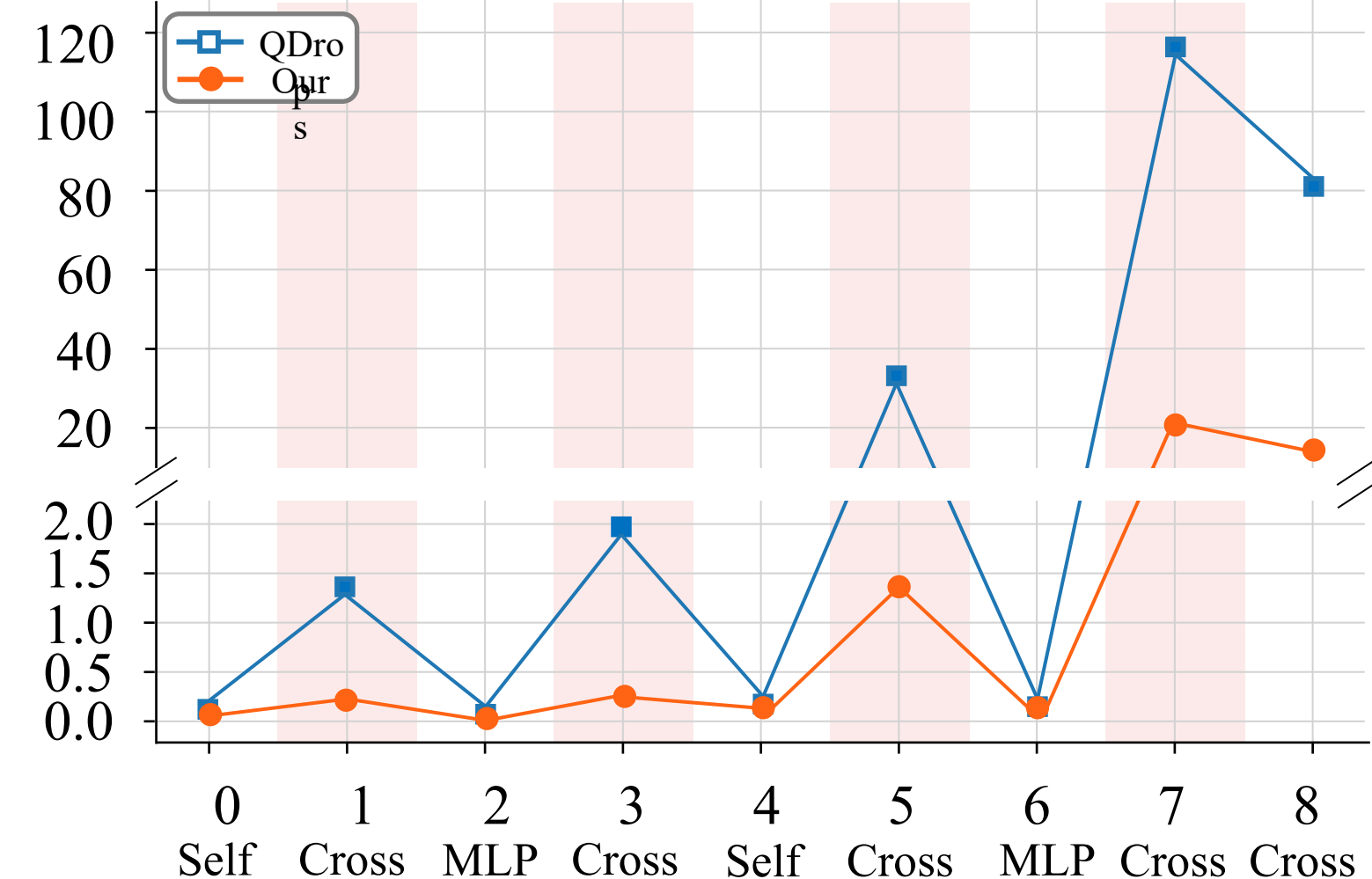
SAMs achieve strong prompt-guided segmentation, but their large memory and compute cost hinder deployment on edge devices. Existing PTQ methods mainly target encoder-only ViTs and do not explicitly handle the bidirectional cross-attention structure in the SAM decoder.



## Attention dissipation



## Reconstruction Loss Oscillation

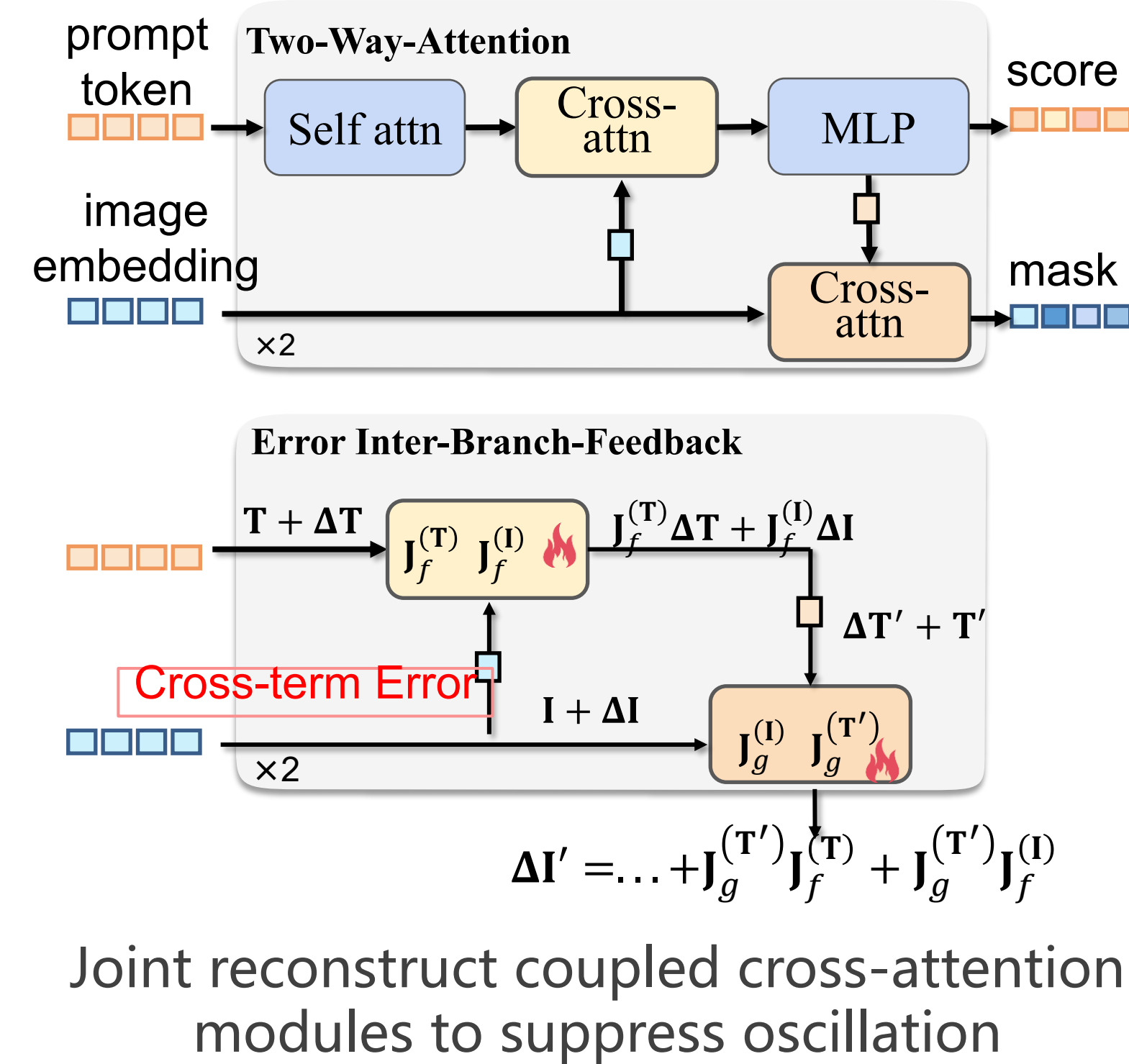


Cross-attention activations fluctuate across decoder blocks, causing unstable reconstruction under low-bit quantization

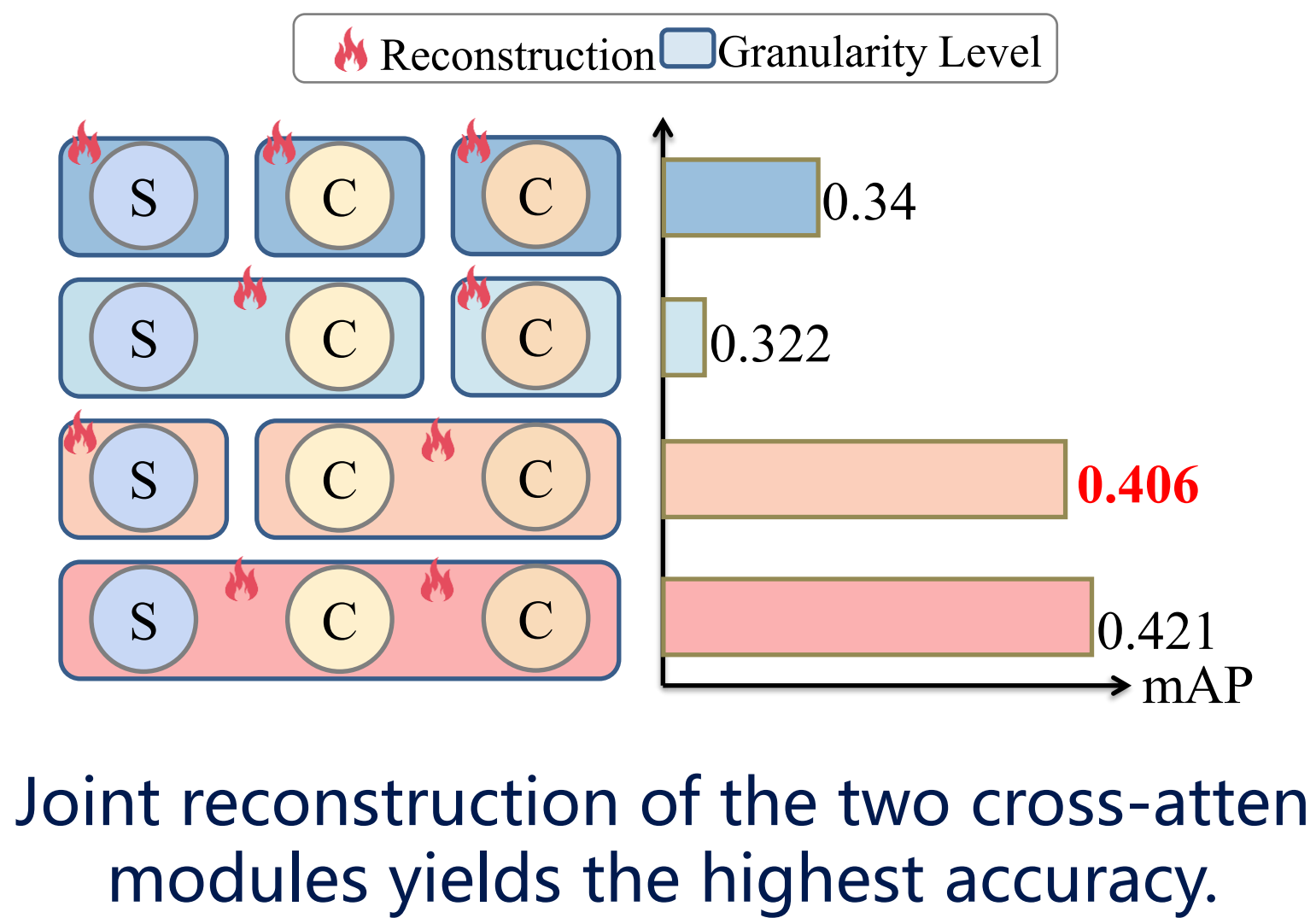


## Method Overview

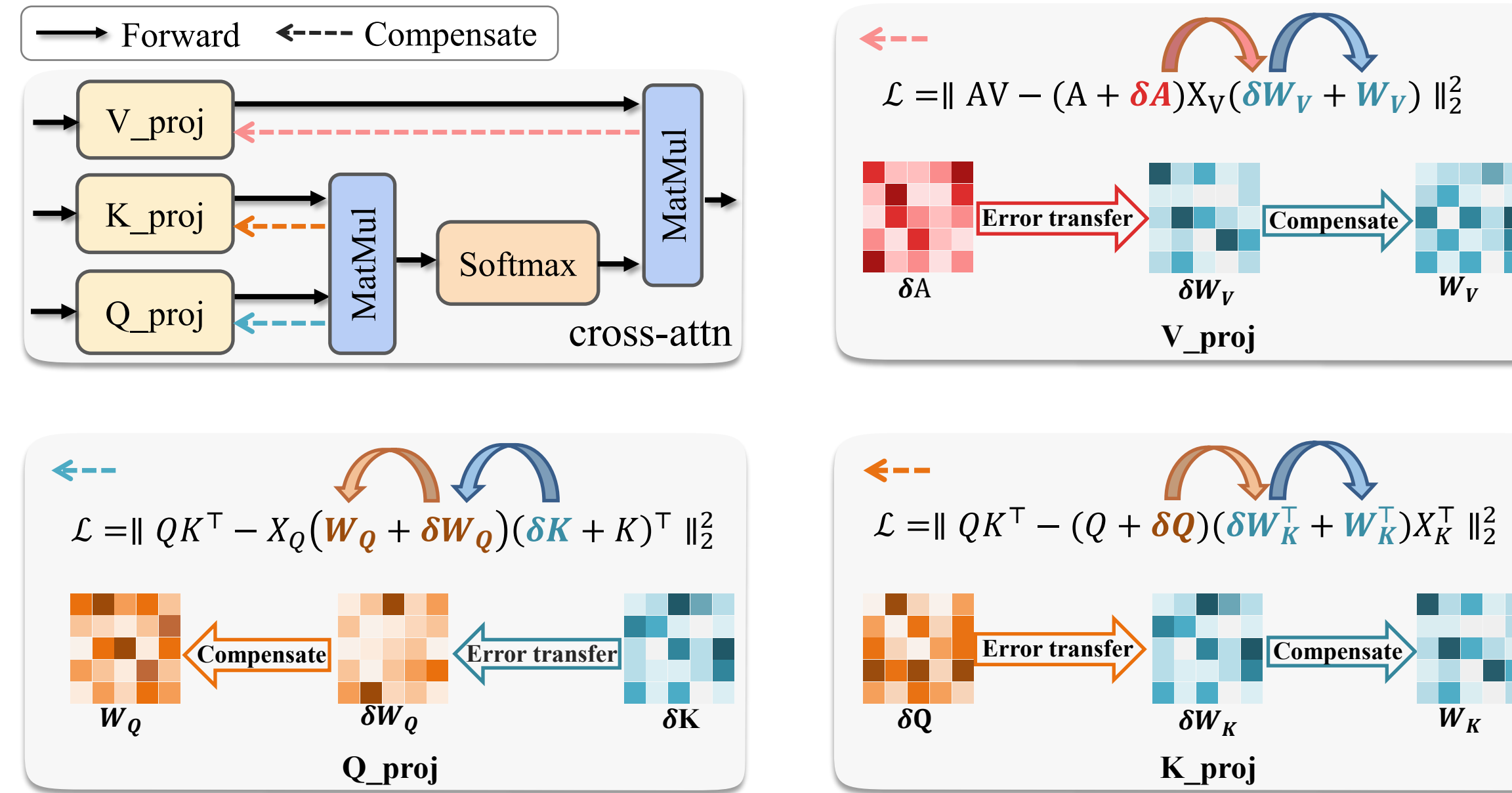
### A) Joint Reconstruction



### Why Joint Reconstruction?



### B) MatMul-Aware Compensation



Compensate Matmul-Induced activation error by updating Q/K/V projection weight

### 1) Objections: Minimize MatMul Error

$$\mathcal{L} = \| QK^T - X_Q(W_Q + \delta W_Q)(\delta K + K)^T \|_2^2$$

### 2) Sylvester Reformulation

$$AX + XB = C$$

where

$$A = (X_Q^T X_Q)^{-1} \lambda I$$

$$B = \hat{K}^T \hat{K}$$

$$C = W_Q \delta K^T \hat{K}$$

$$X = \delta W_Q$$

### 3) Weight update

$$W_Q = W_Q + \delta W_Q$$

$$W_K = W_K + \delta W_K$$

$$W_V = W_V + \delta W_V$$



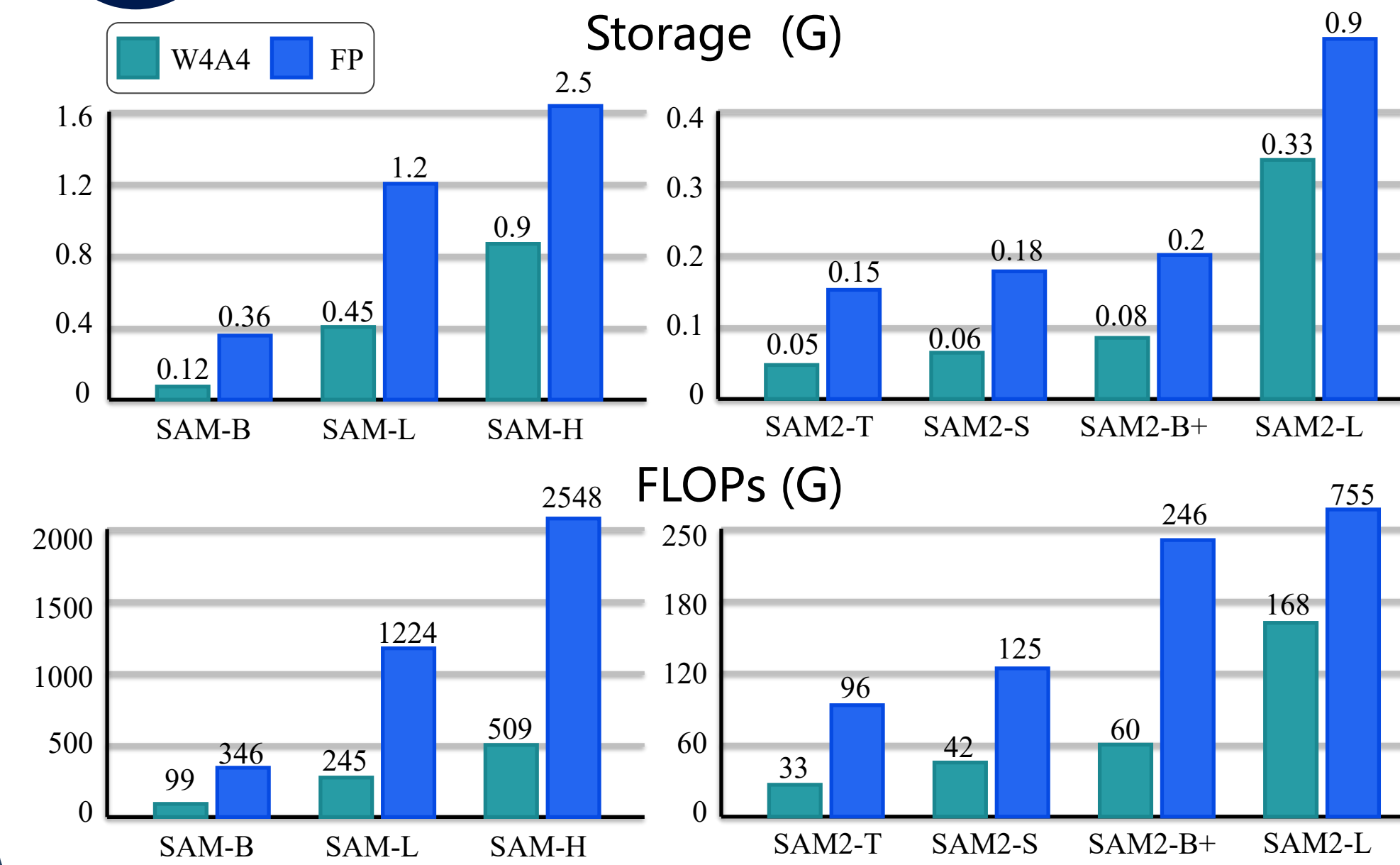
## Main result (COCO W4A4)

Model	QDrop	PTQ4SAM	CAR-SAM (Ours)
SAM-B	26.4	24.7	39.3
SAM-L	38.0	41.9	48.5
SAM-H	48.8	50.7	51.6

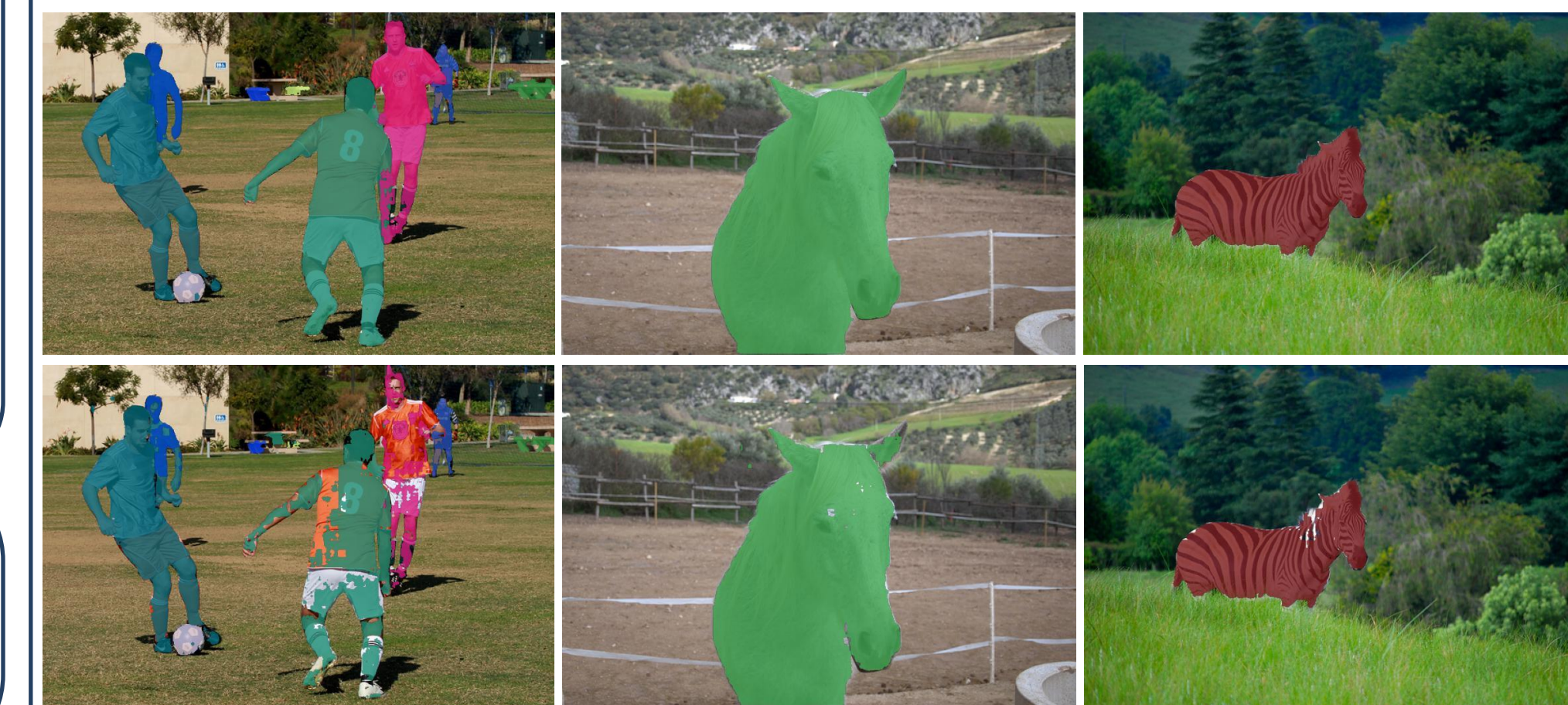
★ +14.6% mAP on SAM-B and +6.6% on SAM-L over prior methods at 4bit



## Efficiency (W4A4 vs FP16)



## Qualitative Result



## Conclusion

CAR-SAM identifies two decoder-specific quantization failures in SAM: attention dissipation and reconstruction oscillation. MAC compensates MatMul-induced errors, while JCAR stabilizes coupled cross-attention reconstruction, enabling robust 4-bit quantization for SAM and SAM2.