



**AAAI-25 / IAAI-25 / EAAI-25**

**FEBRUARY 25 – MARCH 4, 2025 | PHILADELPHIA, USA**



# Qua<sup>2</sup>SeDiMo: Quantifiable Quantization Sensitivity of Diffusion Models

Keith G. Mills<sup>1,2</sup>, Mohammad Salameh<sup>2</sup>, Ruichen Chen<sup>1</sup>,  
Negar Hassanpour<sup>2</sup>, Wei Lu<sup>3</sup> and Di Niu<sup>1</sup>

<sup>1</sup>Dept. ECE, University of Alberta <sup>2</sup>Huawei Technologies Canada Co., Ltd <sup>3</sup>Huawei Kirin Solution, Shanghai, China



**UNIVERSITY  
OF ALBERTA**



**HUAWEI**

**ALBERTA**   
**INNOVATES**

# Motivation

- Diffusion Models generate great visual content!
  - Examples: SDXL, PixArt, Hunyuan, etc.



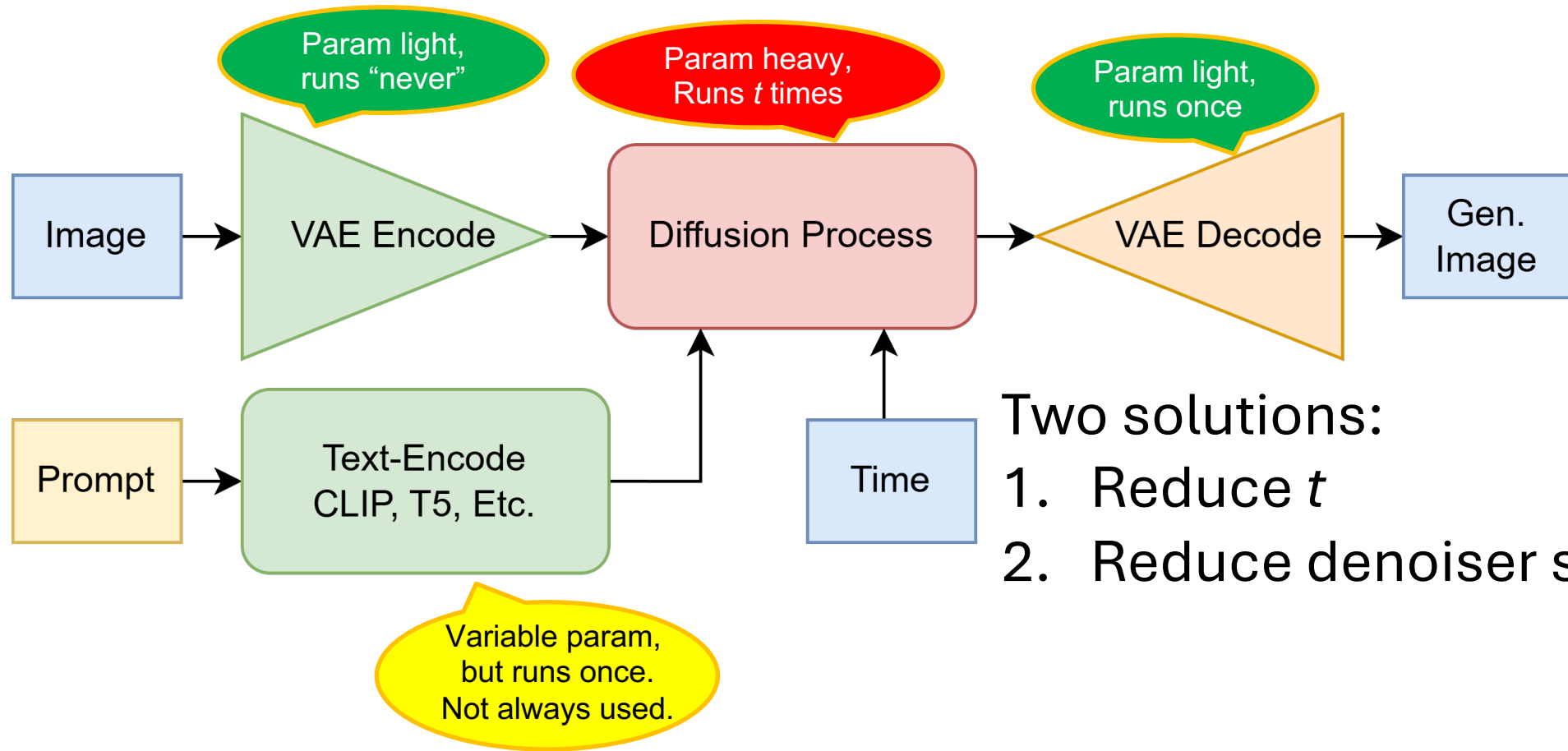
PixArt- $\Sigma$

Prompt:  
“A western-  
style medieval  
dragon with  
large white  
wings spread  
wide”



HunYuan-DiT

# Problem



Two solutions:

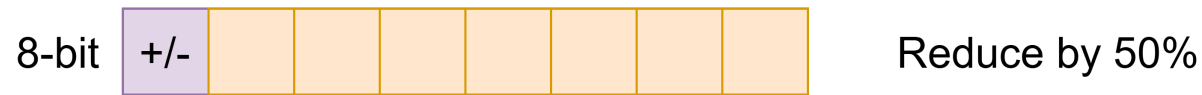
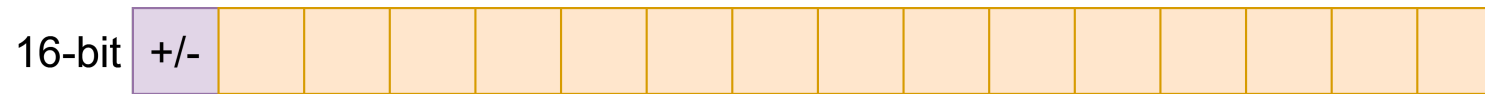
1. Reduce  $t$
2. Reduce denoiser size

# Quantization

Reduces bit precision of weights/activations.

Quantization-Aware Training (QAT) (costly)

Post-Training Quantization (PTQ) (feasible)

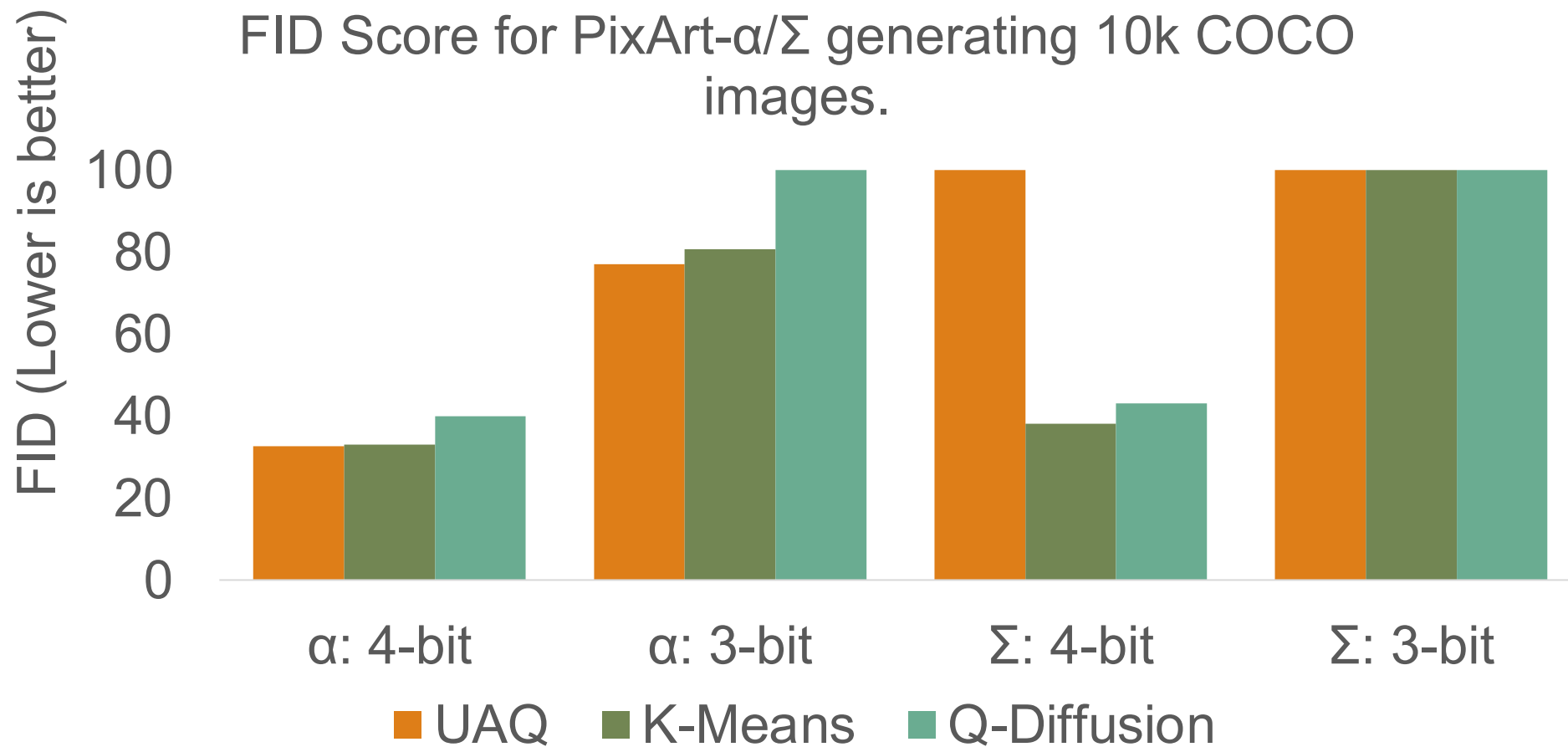


**PTQ struggles below this point**

---



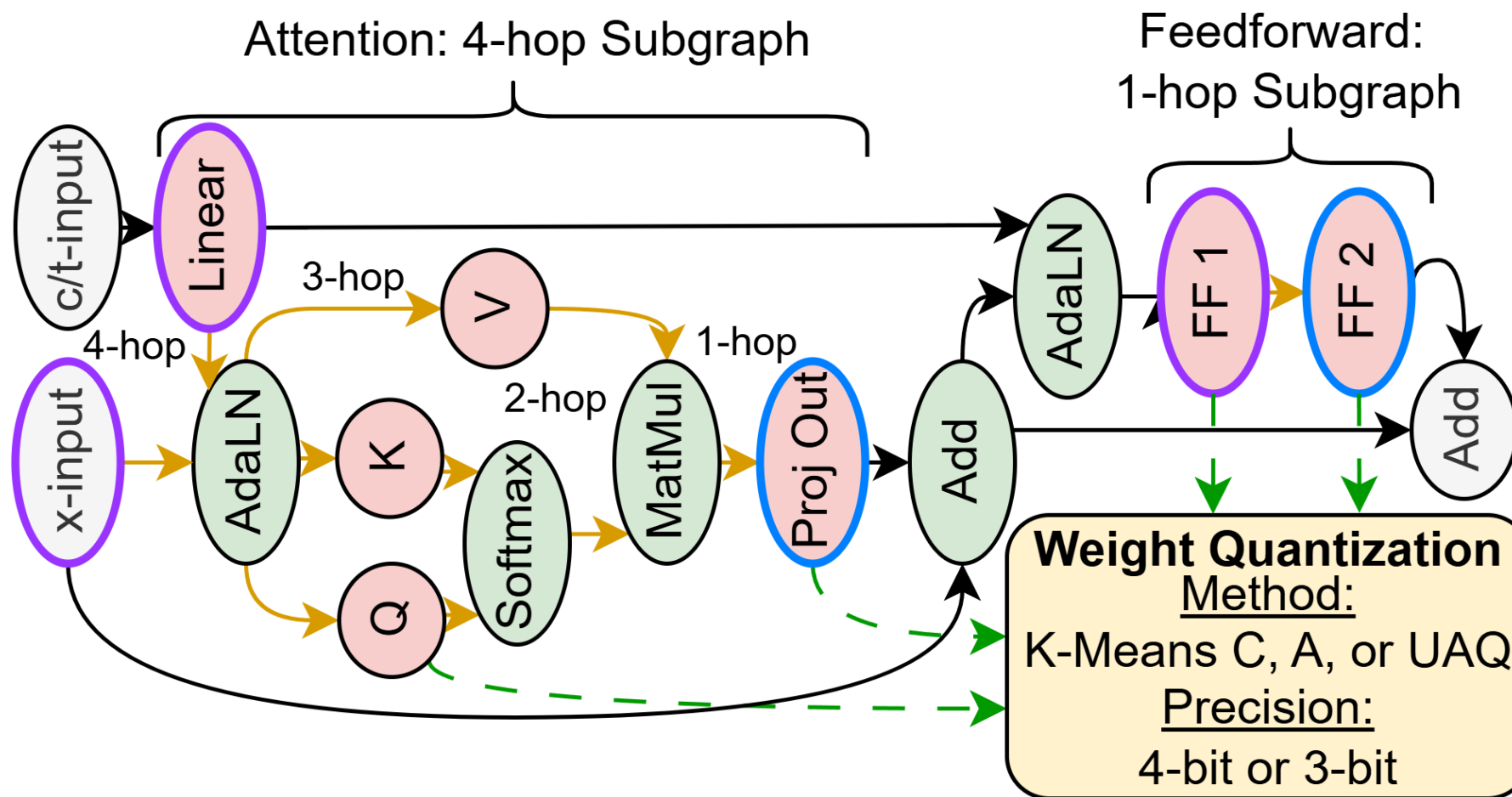
# 4-bit and 3-bit PTQ



# Why? Sensitivity Hypothesis

- Not all weights cause 3-bit performance loss.
- Assert that some *sensitive* weights are the culprit.
  - Motivates mixed-precision approach!
- “Sensitive”?
  - Individual weights? Too granular.
  - Weight categories? E.g., time-embed vs. caption-embed.
  - Weights in specific transformer blocks, like first/last?
- How to find sensitive weights?

# Our Solution



# Predictor with Hop-Level Ranking Loss

## Preliminary: Graphs and GNNs

- $(arch, perf) = (G_1, y_1)$
- Learn  $y'_1 = GNN(G_1)$

Building Optimal Neural Architectures using Interpretable Knowledge

CVPR'24

Keith G. Mills<sup>1,2</sup> Fred X. Han<sup>2</sup> Mohammad Salameh<sup>2</sup> Shengyao Lu<sup>1</sup>  
Chunhua Zhou<sup>3</sup> Jiao He<sup>3</sup> Fengyu Sun<sup>3</sup> Di Niu<sup>1</sup>

<sup>1</sup>Dept. ECE, University of Alberta <sup>2</sup>Huawei Technologies Canada <sup>3</sup>Huawei Kirin Solution, China

{kgmills, shengyao, dniu}@ualberta.ca sunfengyu@hisilicon.com

{fred.xuefei.han1, mohammad.salameh, zhouchunhua, hejiao4}@huawei.com

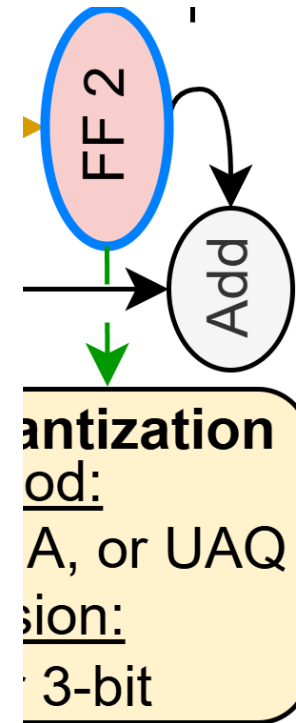
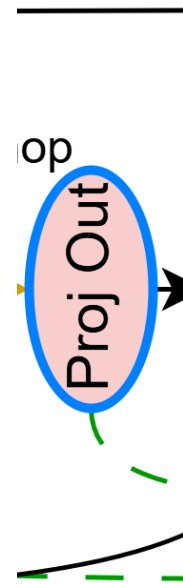
## Intermediate workings: Node and Graph Embeddings

- $GNN(G) = MLP(h_G^m)$ ;  $h_G^m = \frac{1}{|V_G|} \sum_{v \in V_G} h_v^m$
- $m$  is hop-level  $\Rightarrow h_v^m$  represents an *entire* subgraph/module!

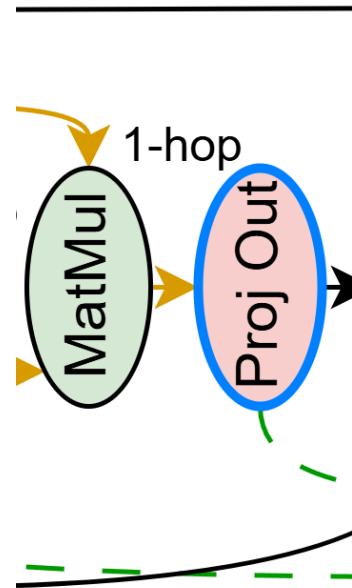
Key learning constraint: if  $y_1 > y_2$ , then  $\|h_{G_1}\|_1 > \|h_{G_2}\|_1$



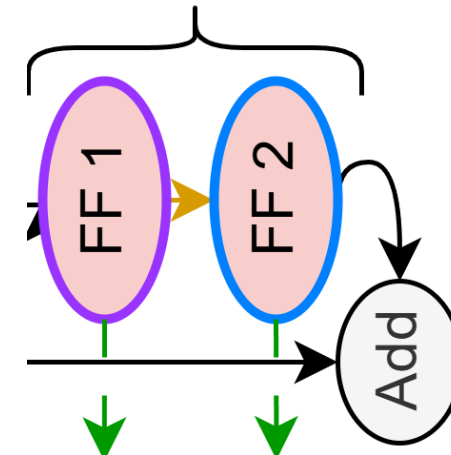
# Visual Example



# Visual Example

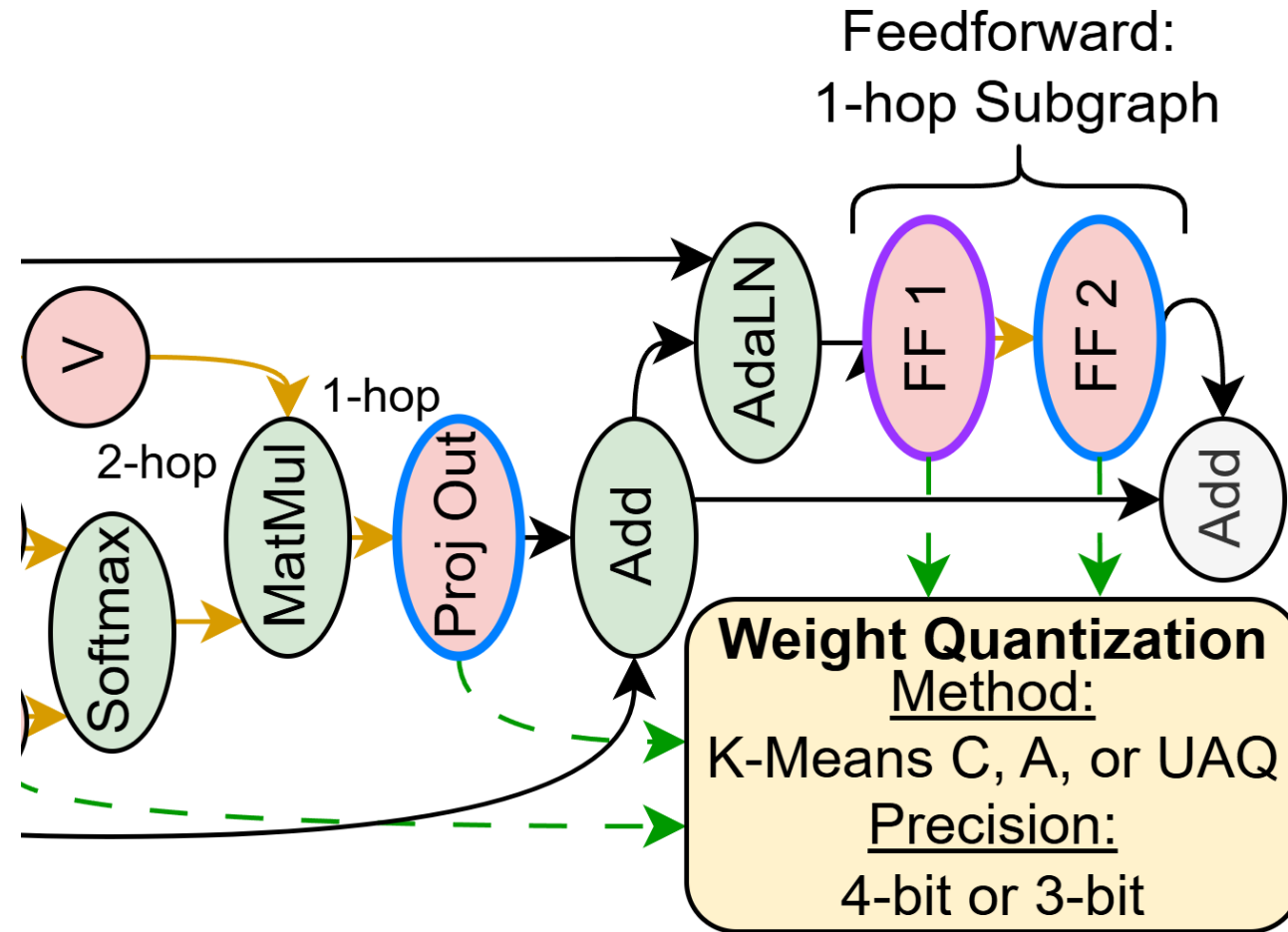


Feedforward:  
1-hop Subgraph

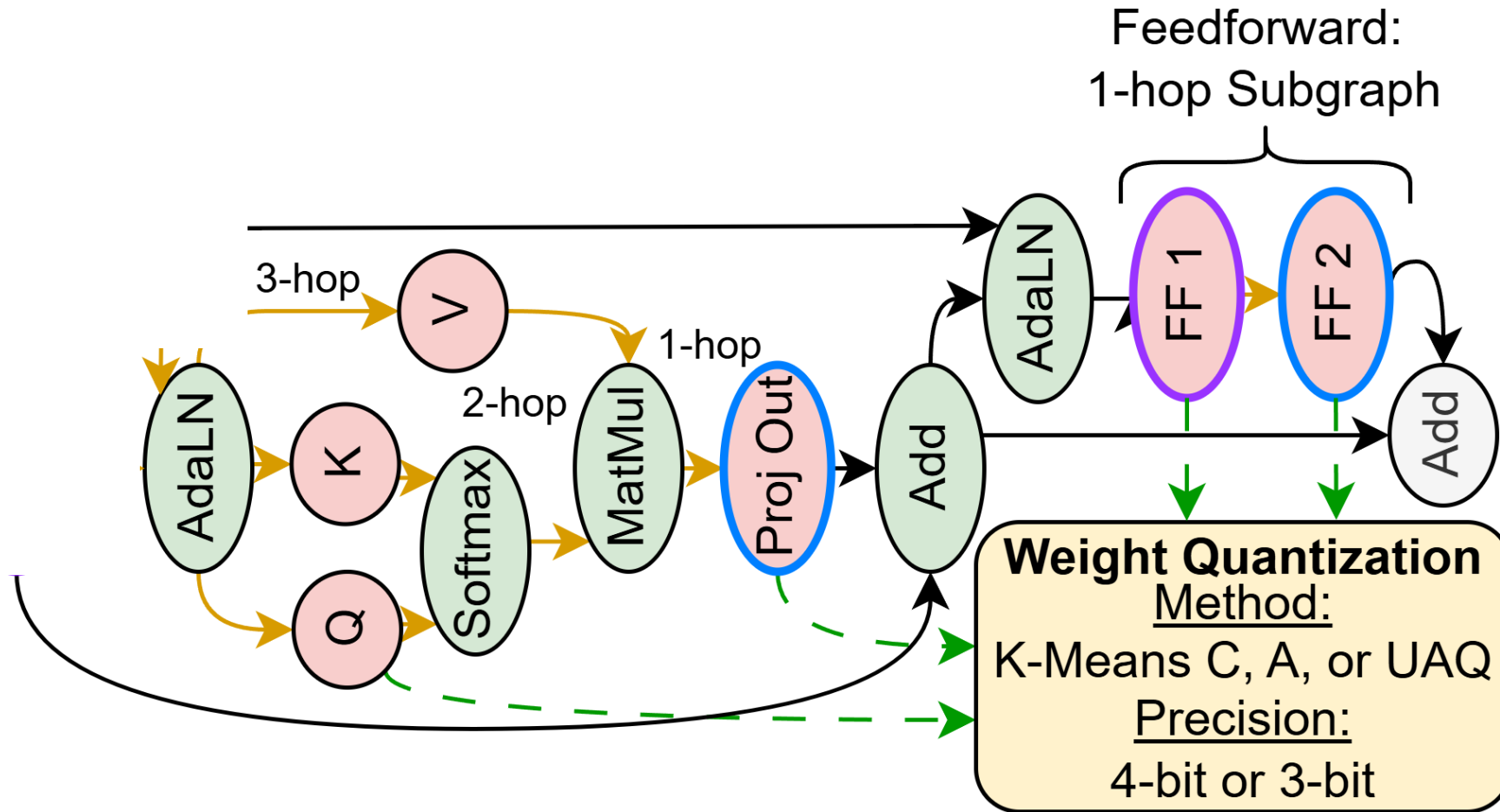


**ht Quantization**  
Method:  
ans C, A, or UAQ  
Precision:  
4-bit or 3-bit

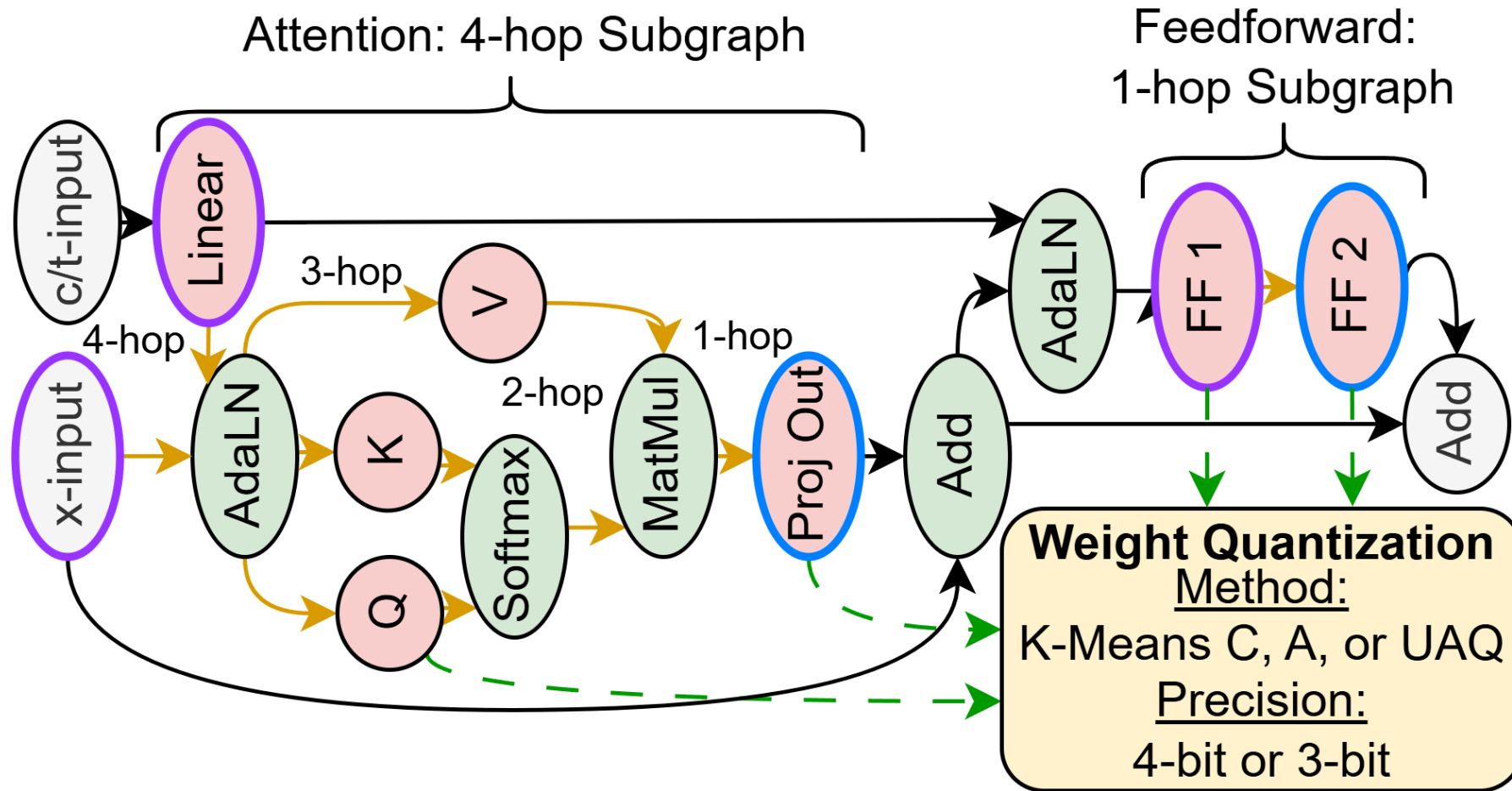
# Visual Example



# Visual Example



# Visual Example



# Predictor with Hop-Level Ranking Loss

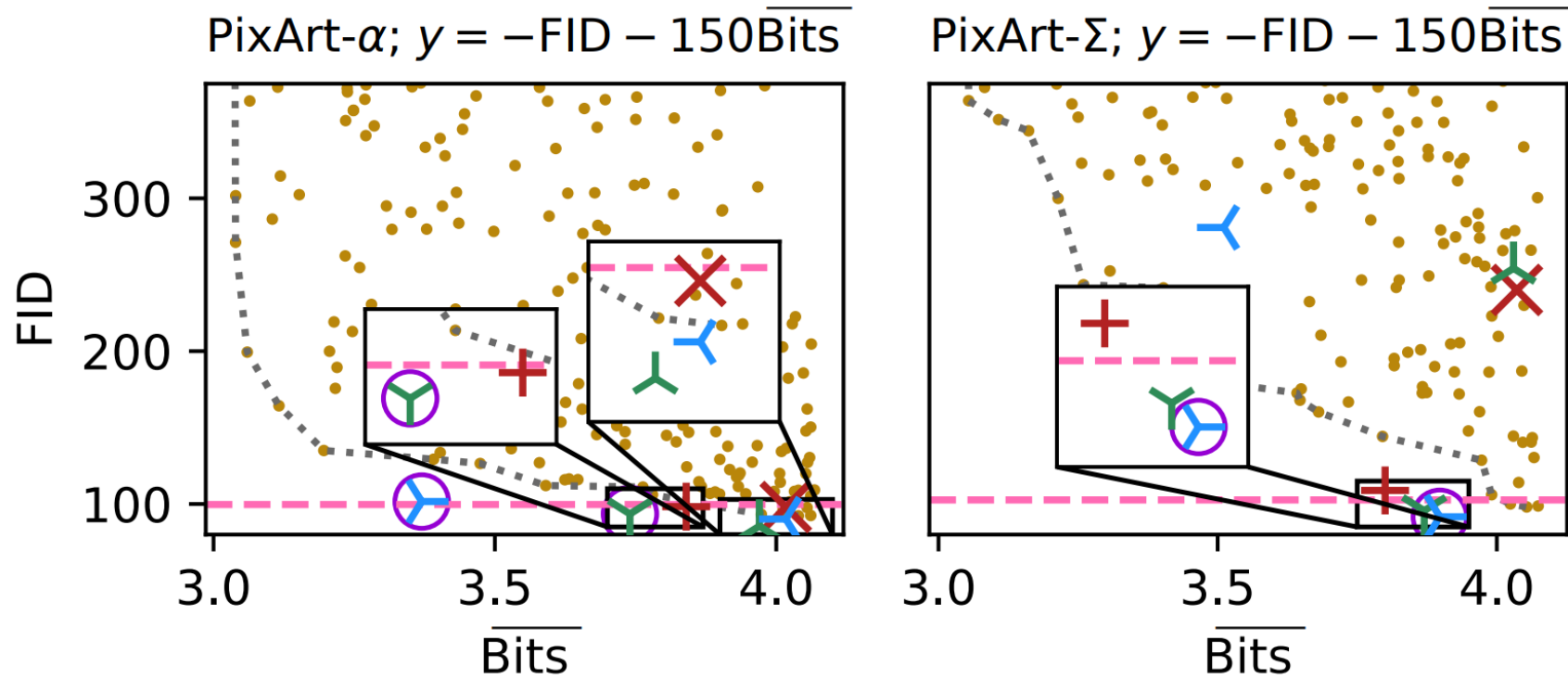
Preliminary: Graphs and GNNs

- $(arch, perf) = (G_1, y_1)$
- Learn  $y'_1 = GNN(G_1)$

Optimize  $L_{orig}(y, y') + \frac{1}{M+1} \sum_{m=0}^M L_{rank}(y, \|h_G^m\|_1)$

- $L_{orig}$  is traditional predictor loss, like MSE
- $L_{rank}$  is SRCC, LambdaRank, or both.

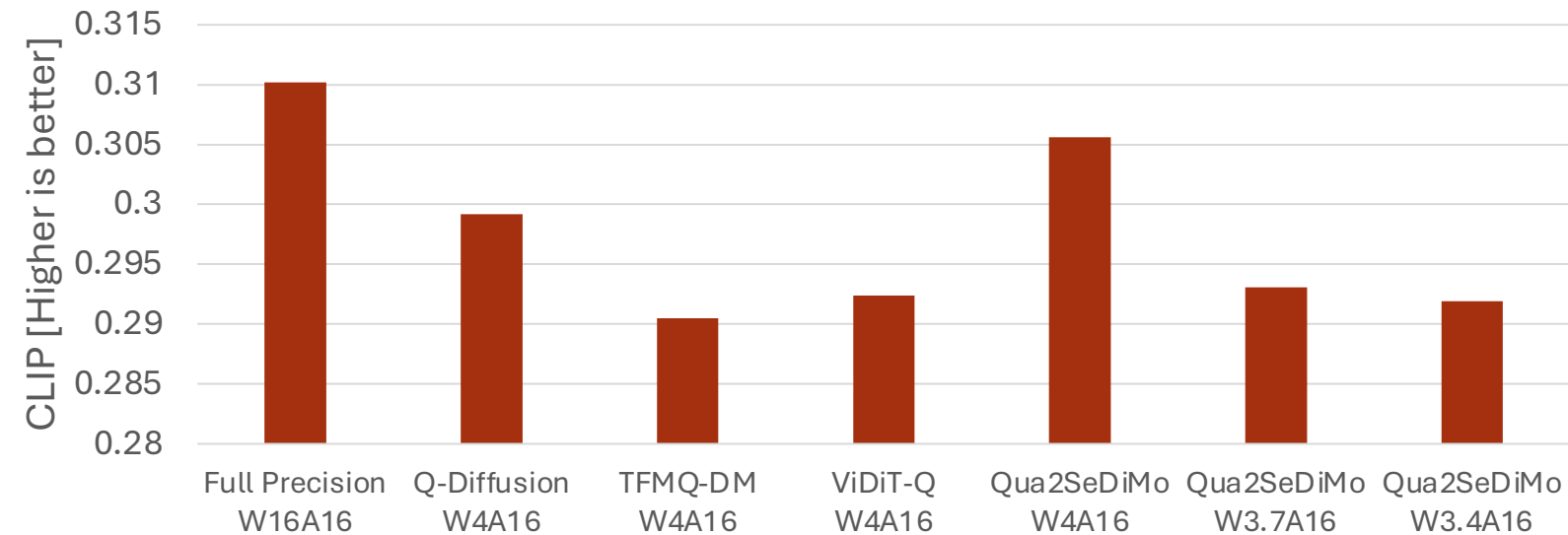
# Pareto Frontier Results



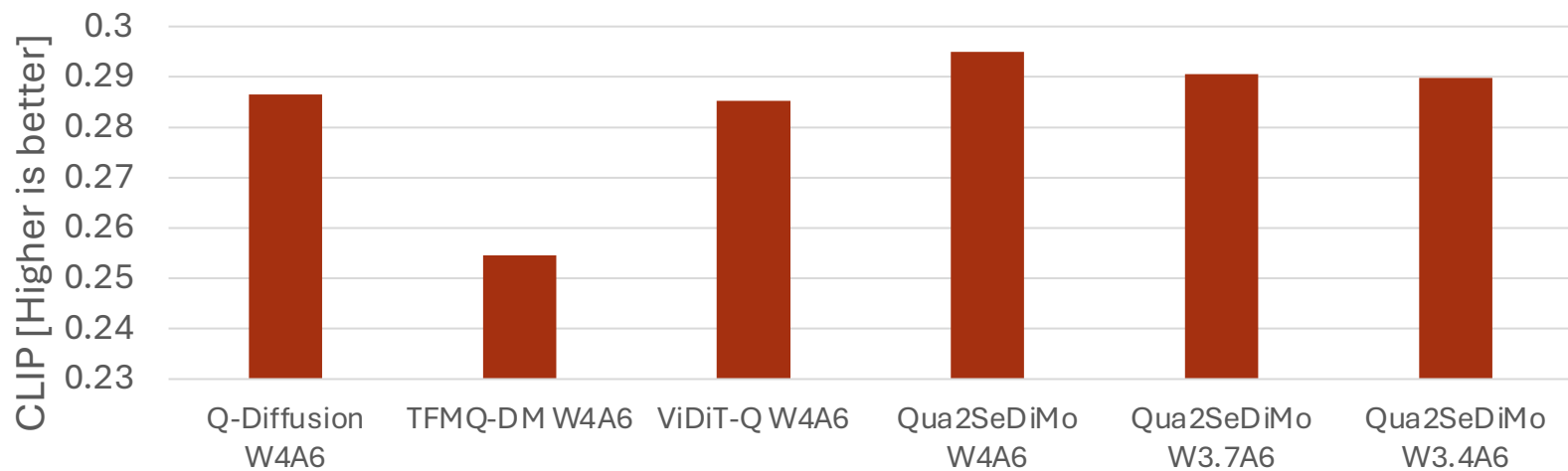
× SRCC Op-level   + SRCC Block-level   × NDCG Op-level   y NDCG Block-level   x Hybrid Op-level   y Hybrid Block-level

# Pareto Frontier Results

Quantitative CLIP on PixArt- $\alpha$

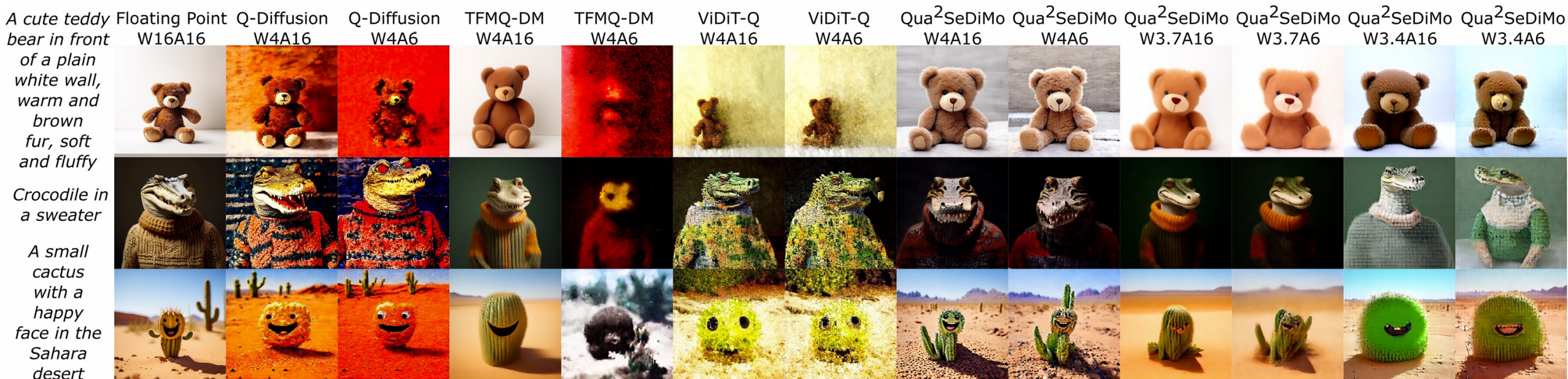


Quantitative CLIP on PixArt- $\alpha$  W\*A6

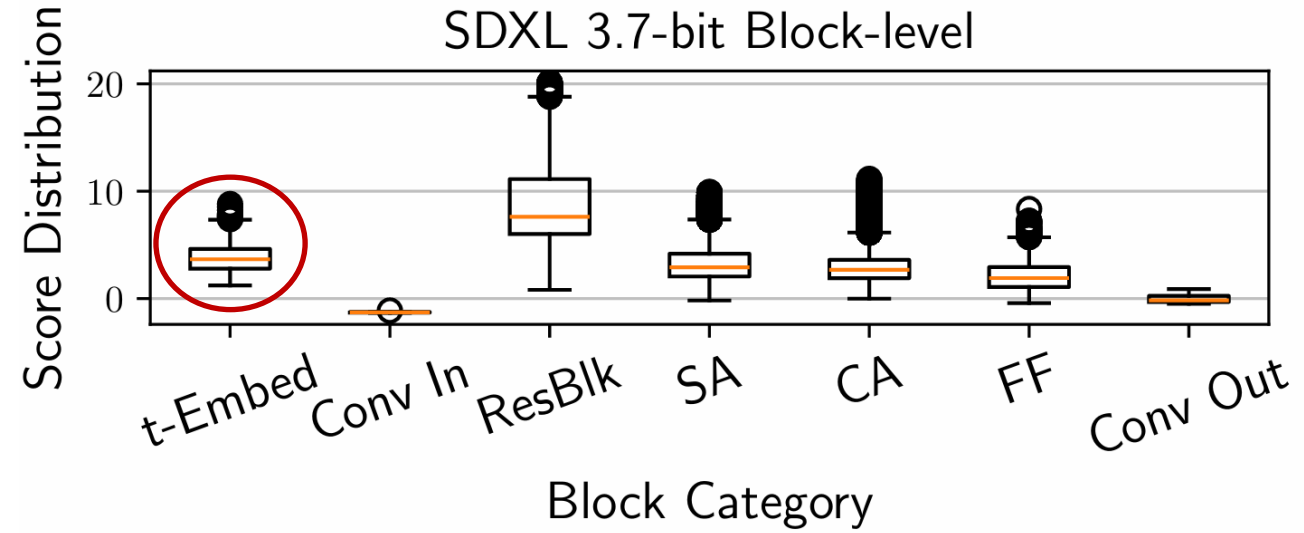
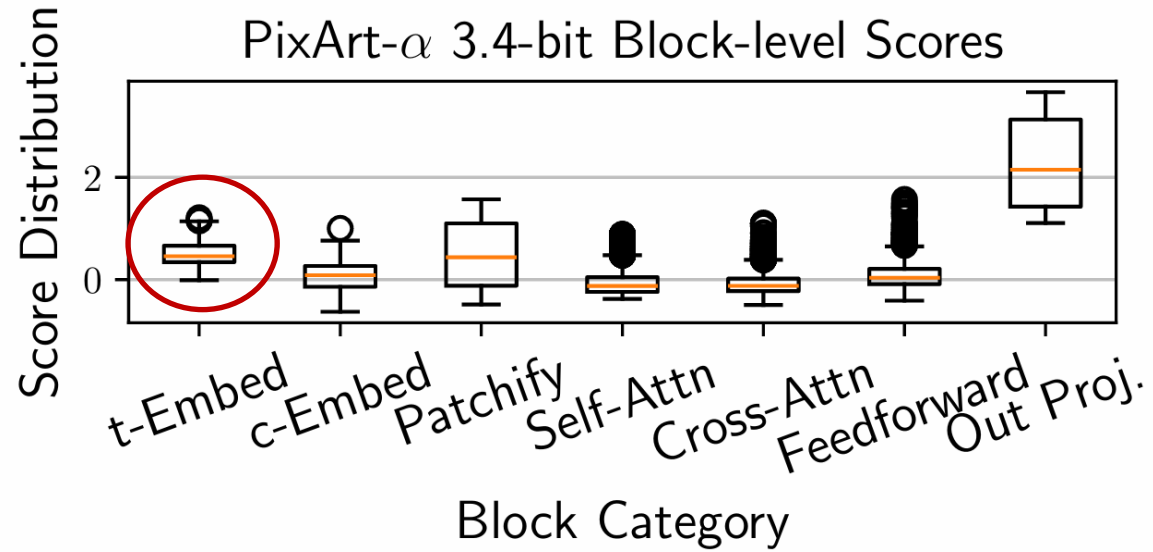




# Qualitative Visual Results



# Sample Insights





**AAAI-25 / IAAI-25 / EAAI-25**

**FEBRUARY 25 – MARCH 4, 2025 | PHILADELPHIA, USA**



**Qua<sup>2</sup>SeDiMo:  
Quantifiable Quantization Sensitivity of  
Diffusion Models**

Thank you for watching 'till the end!

See you in Philly!



**UNIVERSITY  
OF ALBERTA**



**HUAWEI**

**ALBERTA**   
**INNOVATES**