### SVDQuant: Absorbing Outliers by Low-Rank Comp onents for 4-Bit Diffusion Models

ICLR 2025

Muyang Li, Yujun Lin, Zhekai Zhang, Tianle Cai, Xiuyu Li, Junxian Guo, E nze Xie, Chenlin Meng, Jun-Yan Zhu, Song Han

Presented by Changhwan Sung<sup>1</sup>, Jungyu Jin<sup>1</sup>, Junho Jeong<sup>2</sup>

<sup>1</sup>Graduate School of AI, <sup>2</sup>Electrical Engineering, POSTECH May 21, 2025





**1. Introduction** 

2. Methodology

3. Experimental results

4. Conclusion



### **Motivation**

- Diffusion models generate high-quality images but are computationally intensive.
- As model size scales (e.g., SD 1.4 (0.8B) → SDXL (2.6B) → FLUX.1 (12B)): Better quality. Heavier memory.

Higher inference latency.

 Low-bit quantization (4-bit) is promising Goal: W4A4 Quantization. FLUX.1-dev BF16 (25 Steps) DiT Memory: 22.7 GiB **E2E Latency: 111.7 s** 



SVDQuant INT4 (W4A4) LPIPS: 0.223 DiT Memory: 6.5 GiB (3.5× Less) E2E Latency: 12.9 s (8.7× Faster)



Prompt: a cyberpunk cat holding a huge neon sign that says "SVDQuant is lite and fast", wearing fancy goggles and a black leather jacket.



## **Limits of Standard PTQ in Diffusion Models**

■ As model size scales (e.g., SD 1.4 (0.8B)  $\rightarrow$  SDXL (2.6B)  $\rightarrow$  FLUX.1 (12B)):

Better quality. Heavier memory.

Higher inference latency.

- Unlike LLMs:
  - LLMs are bandwidth-bound
  - Diffusion models are compute-bound
  - FLUX.1 (12B) vs Llama2 (13B):
     FLUX.1 is over 3x heavier than Llama2





### **Prior Works**

- Q-Diffusion, PTQ4DM: 8-bit quantization.
  - For U-Net backbone.
- MixDQ, ViDiT-Q:
  - Use low-bit inference engines.
  - Show actual speedups on GPUs.
  - ViDiT-Q: W8A8, MixDQ: W4A8
- Most prior work:
  - Targets only weights (W4A16)
  - No support for both weights & activations in 4-bit



# **Existing Strategy: Smoothing**

- Shift outliers from activations to weights
  - Scaling each channel in X by smoothing factor  $\lambda$ .
  - $\hat{X} = X \cdot \operatorname{diag}(\lambda)^{-1}, \ \widehat{W} = W \cdot \operatorname{diag}(\lambda).$
  - $Y = (X \cdot \operatorname{diag}(\lambda)^{-1}) \cdot (W \cdot \operatorname{diag}(\lambda))$
- For Diffusion Models, weight outlier is also critical.
  - Looks like this:





### **4bit Quantization + Low-Rank Branch**

- High-precision branch for outliers.
  - 16-bit low rank branch + 4-bit residual branch.
  - $Y = (X \cdot \operatorname{diag}(\lambda)^{-1}) \cdot (W \cdot \operatorname{diag}(\lambda)) = \hat{X} \cdot \hat{W} = \hat{X} \cdot (L_1 L_2 + R) = \hat{X} L_1 L_2 + \hat{X} R$
  - $Y \approx \hat{X}L_1L_2 + s_X s_R Q(\hat{X})Q(R)$





## **Diffusion model (U-Net)**





## **Diffusion model (DiT)**





# **Quantization Preliminary**

• Given a tensor X, the quantization process is defined as:

$$Q_X = round\left(\frac{X}{s_X}\right)$$
,  $s_X = \frac{\max(|X|)}{q_{max}}$ 

- For signed k-bit integer quantization,  $q_{max} = 2^{k-1} 1$
- For 4-bit floating-point quantization (E2M1),  $q_{max} = 6 (1.1_2 * 2^{3-1})$
- Dequantized tensor  $Q(X) = s_X \cdot Q_X$
- For a linear layer with input *X* and weight *W*, its computation:

$$XW \approx Q(X)Q(W) = s_X s_W \cdot Q_X Q_W$$

- we denote x-bit weight, y-bit activation as  $W_x A_y$ .
  - This work focus on W4A4 quantization for acceleration.





# **Quantization Preliminary**

- Traditional methods to suppress outliers include QAT and rotation.
  - QAT requires massive computing resources.
- Rotation is inapplicable due to the usage of adaptive normalization layers.
  - Online rotation of both activations and weights incurs significant runtime overhead.



Latent Diffusion Transformer

scaling ( $\alpha$ ) has been absorbed into the weight matrices

DiT Block with adaLN-Zero



### **Problem Formulation**

• Consider a linear layer with input  $X \in \mathbb{R}^{b \times m}$  and weight  $W \in \mathbb{R}^{m \times n}$ , Quantization error can be defined as:

$$E(X,W) = \|XW - Q(X)Q(W)\|_F$$

Proposition 4.1 (Error decomposition).
 The quantization error can be decomposed as follows:

 $E(X,W) \le \|X\|_F \|W - Q(W)\|_F + \|X - Q(X)\|_F (\|W\|_F + \|W - Q(W)\|_F)$ 

• To minimize the overall quantization error, we aim to optimize four terms.



#### POSTECH

### **Proof of Proposition 4.1**

- The proof of proposition use two inequalities.
  - 1) Triangle inequality:  $||A + B||_F \le ||A||_F + ||B||_F$
  - 2) Cauchy-Schwarz Inequality:  $\|AB\|_F \leq \|A\|_F \|B\|_F$

$$\begin{split} \| \boldsymbol{X} \boldsymbol{W} - Q(\boldsymbol{X})Q(\boldsymbol{W}) \|_{F} \\ &= \| \boldsymbol{X} \boldsymbol{W} - \boldsymbol{X}Q(\boldsymbol{W}) + \boldsymbol{X}Q(\boldsymbol{W}) - Q(\boldsymbol{X})Q(\boldsymbol{W}) \|_{F} \\ &\leq \| \boldsymbol{X} (\boldsymbol{W} - Q(\boldsymbol{W})) \|_{F} + \| (\boldsymbol{X} - Q(\boldsymbol{X}))Q(\boldsymbol{W}) \|_{F} \\ &\leq \| \boldsymbol{X} \|_{F} \| \boldsymbol{W} - Q(\boldsymbol{W}) \|_{F} + \| \boldsymbol{X} - Q(\boldsymbol{X}) \|_{F} \| Q(\boldsymbol{W}) \|_{F} \\ &\leq \| \boldsymbol{X} \|_{F} \| \boldsymbol{W} - Q(\boldsymbol{W}) \|_{F} + \| \boldsymbol{X} - Q(\boldsymbol{X}) \|_{F} \| \boldsymbol{W} - (\boldsymbol{W} - Q(\boldsymbol{W})) \|_{F} \\ &\leq \| \boldsymbol{X} \|_{F} \| \boldsymbol{W} - Q(\boldsymbol{W}) \|_{F} + \| \boldsymbol{X} - Q(\boldsymbol{X}) \|_{F} \| \boldsymbol{W} - (\boldsymbol{W} - Q(\boldsymbol{W})) \|_{F} \\ &\leq \| \boldsymbol{X} \|_{F} \| \boldsymbol{W} - Q(\boldsymbol{W}) \|_{F} + \| \boldsymbol{X} - Q(\boldsymbol{X}) \|_{F} \| \boldsymbol{W} - (\boldsymbol{W} - Q(\boldsymbol{W})) \|_{F} \end{split}$$

- SVDQuant 1) first migrates outliers from activation to weight by smoothing, and then
   2) absorbs the migrated weight outliers into the low-rank branch.
- Smooth outliers in activations using a per-channel smoothing factor  $\lambda \in \mathbb{R}^m$ :



• While the smoothed input  $\hat{X}$  exhibits fewer outliers,  $\hat{W}$  has a significant increase in both magnitude and the presence of outliers.  $\rightarrow$  raise  $\|W - Q(W)\|_F$ 



- Core insight is to introduce a 16-bit low-rank branch to further migrate the weight quantization difficulty.
- Specifically, we decompose  $\widehat{W} = L_1L_2 + R$ , where  $L_1 \in R^{m \times r}$ ,  $L_2 \in R^{r \times n}$ . Then XW can be approximated as:

$$XW = \hat{X}\hat{W} = \hat{X}L_1L_2 + \hat{X}R \approx \hat{X}L_1L_2 + Q(\hat{X})Q(R)$$
  
16-bit low-rank path 4-bit residual path

- r is typically set to 16 or 32, so additional parameters and computation for the lowrank branch are negligible.
  - However, it still requires careful system design to eliminate redundant memory access



From reformulated equation, quantization error can be expressed as

$$E(X,W) = \|XW - Q(X)Q(W)\|_{F}$$
  
=  $\|\hat{X}\widehat{W} - (\hat{X}L_{1}L_{2} + Q(\hat{X})Q(R))\|_{F} = \|\hat{X}R - Q(\hat{X})Q(R)\|_{F} = E(\hat{X},R)$ 

Since  $\hat{X}$  is already free from outliers, we only need to focus on optimizing the  $\|R\|_F$  and  $\|R - Q(R)\|_F$ .

• Then, how to minimize terms related to  $R = \widehat{W} - L_1 L_2$ ?



Proposition 4.2 (Quantization error bound).

For any tensor **R** and quantization method  $Q(\mathbf{R}) = s_{\mathbf{R}}Q_{\mathbf{R}}$ . Assuming the elements of **R** follow a distribution that satisfies the following regularity condition: There exists a constant c such that

$$E[max(|\mathbf{R}|)] \le c \cdot E[||\mathbf{R}||_F].$$

Then, we have

$$E[\|\boldsymbol{R} - Q(\boldsymbol{R})\|_F] \leq \frac{c\sqrt{size(R)}}{q_{max}} \cdot E[\|\boldsymbol{R}\|_F].$$

- From this proposition, the quantization error is bounded by the magnitude of **R**.
  - Thus, our goal is to find the optimal  $L_1L_2$  that minimizes  $\|\mathbf{R}\|_F = \|\widehat{W} L_1L_2\|_F$ .



### **Proof of Proposition 4.2**

**Proposition 4.2** (Quantization error bound). For any tensor  $\mathbf{R}$  and quantization method described in Equation 1 as  $Q(\mathbf{R}) = s_{\mathbf{R}} \cdot \mathbf{Q}_{\mathbf{R}}$ . Assuming the elements of  $\mathbf{R}$  follow a distribution that satisfies the following regularity condition: There exists a constant c such that

$$\mathbb{E}\left[\max(|\boldsymbol{R}|)\right] \le c \cdot \mathbb{E}\left[\|\boldsymbol{R}\|_{F}\right].$$
(7)

Then, we have

$$\mathbb{E}\left[\left\|\boldsymbol{R} - Q(\boldsymbol{R})\right\|_{F}\right] \leq \frac{c\sqrt{size(\boldsymbol{R})}}{q_{\max}} \cdot \mathbb{E}\left[\left\|\boldsymbol{R}\right\|_{F}\right]$$
(8)

where  $size(\mathbf{R})$  denotes the number of elements in  $\mathbf{R}$ . Especially if the elements of  $\mathbf{R}$  follow a normal distribution, Equation 7 holds for  $c = \sqrt{\frac{\log(size(\mathbf{R}))\pi}{size(\mathbf{R})}}$ .



$$\begin{aligned} \|\boldsymbol{R} - Q(\boldsymbol{R})\|_{F} \\ &= \|\boldsymbol{R} - s_{\boldsymbol{R}} \cdot \boldsymbol{Q}_{\boldsymbol{R}}\|_{F} \\ &= \left\|s_{\boldsymbol{R}} \cdot \frac{\boldsymbol{R}}{s_{\boldsymbol{R}}} - s_{\boldsymbol{R}} \cdot \operatorname{round}\left(\frac{\boldsymbol{R}}{s_{\boldsymbol{R}}}\right)\right\|_{F} \\ &= |s_{\boldsymbol{R}}| \left\|\frac{\boldsymbol{R}}{s_{\boldsymbol{R}}} - \operatorname{round}\left(\frac{\boldsymbol{R}}{s_{\boldsymbol{R}}}\right)\right\|_{F} \\ &\leq \mathbb{E}\left[\|\boldsymbol{R} - Q(\boldsymbol{R})\|_{F}\right] \\ &\leq \mathbb{E}\left[|s_{\boldsymbol{R}}|\right] \sqrt{\operatorname{size}(\boldsymbol{R})} \\ &= \frac{\sqrt{\operatorname{size}(\boldsymbol{R})}}{q_{\max}} \cdot \mathbb{E}\left[\max(|\boldsymbol{R}|)\right] \\ &\leq \frac{c\sqrt{\operatorname{size}(\boldsymbol{R})}}{q_{\max}} \cdot \mathbb{E}\left[\|\boldsymbol{R}\|_{F}\right] \end{aligned}$$

- The problem of minimizing  $\|\widehat{W} L_1 L_2\|_F$  can be solved by Singular Value Decomposition (SVD).
  - Minimize  $\|\widehat{W} L_1 L_2\|_F$  => Maximize  $\|L_1 L_2\|_F$  ·  $\|A\|_F = \sqrt{\sum_i^m \sum_j^n A_{ij}^2} = \sqrt{\operatorname{trace}(A^T A)} = \sqrt{\sum_{i=1}^{\min\{m,n\}} \sigma_i^2(A)},$ where  $\sigma_i(A)$  is i-th singular value of A.
- Given SVD of  $\widehat{W} = U\Sigma V$ , the optimal solution is  $L_1 = U\Sigma_{:,:r}$  and  $L_2 = V_{:r,:}$ 
  - By removing dominant values, the magnitude of the residual R is dramatically reduced,



### Nunchaku: Fusing Low-rank And Low-bit Branch Kernels

 Although the low-rank branch introduces negligible computation in theory, running it as a separate branch incurs large latency overhead.



Figure 6: (a) Naïvely running low-rank branch with rank 32 will introduce 57% latency overhead due to extra read of 16-bit inputs in *Down Projection* and extra write of 16-bit outputs in *Up Projection*.

- Because the input and output activation sizes remain unchanged, shifting **bottleneck** from computation to **memory access.**
- For example, Up projection for QKV projection is much slower since its output exceeds the available L2 cache, resulting in the extra DRAM load and store operations.

### POSTECH

### Nunchaku: Fusing Low-rank And Low-bit Branch Kernels

- Author eliminate these extra memory access and halving the number of kernel calls by fusing.
  - 1. Down projection( $L_1$ ) shares the **same input** as quantization kernel => Fusing
  - 2. Up projection( $L_2$ ) shares the **same output** as the 4-bit computational kernel => Fusing



Figure 6: (a) Naïvely running low-rank branch with rank 32 will introduce 57% latency overhead due to extra read of 16-bit inputs in *Down Projection* and extra write of 16-bit outputs in *Up Projection*.

■ As a result, low-rank branch adds only 5~10% latency.



### **Experimental Setup**

٠	Models :	FLUX.1	PixArt-Σ	SANA	SDXL	SDXL-Turbo
•	Datasets :	COCO Capt MJHQ-30K sDCI	ions 2024 → fo → 5 → $\xi$	or Calibration K prompts for be 5K prompts for be	nchmarking enchmarking	
•	Comparison :	NF4	ViDiT-Q	MixDQ	TensorRT	SVDQuant (ours)



### **Baselines Comparison**

### PTQ-based Quantization Methods

Category	NF4	ViDiT-Q	MixDQ	TensorRT
Precision	W4	W8A8	W4A8	W8A8
Core Technique	NormalFloat: nonlinear normal dis tribution-based 4-bit	Token-level quantization + s moothing	16-bit for initial tokens + 4-bi t for the rest	Percentile-based calibration + smoothing
Strengths	Very low memory usage, strong w eight preservation	Strong for text alignment and token outliers	Minimal performance loss, maximizing compression efficiency	Optimized for production, fas t inference, stable
Weaknesses	Fixed activation precision (needs F P16 etc.)	Limited memory savings due to 8-bit usage	Performance varies by sente nce length, complex logic	Hard to customize, only activ ation is quantized
Use Cases	FLUX.1	PixArt-Σ	SDXL-Turbo	SDXL, LLaMA



### **Models Comparison**

### Diffusion Transformer & UNet based Models

Model	Parameters (B)	Architecture	Step Type	Estimated Memory (FP16)	Use Case
FLUX.1	12B	DiT (Large)	multi-step	22.2 GiB	High-fidelity generatio n, LoRA support
PixArt-Σ	0.6B	DiT (Compact PixArt variant)	multi-step	2~3 GiB	Compact model, alignment-sensitive
SANA	0.6B	Linear DiT	1-step	<4 GiB	Real-time generation, low-latency
SDXL	2.6B	LDM (multi-step)	multi-step (~30)	6~7 GiB	High-quality baseline diffusion
SDXL-Turbo	~1B	Optimized SDXL (1~2 steps)	1~2 steps	4~5 GiB	Latency-critical, real-time image gen



### **Visual Quality Results**



Figure 7: Oualitative visual results on MJHO. Image Reward is calculated over the entire dataset. On FLUX.1 models, our 4-bit models outperform the NF4 W4A16 baselines, demonstrating superior text alignment and closer similarity to the 16-bit models. For instance, NF4 misses the swinging chair in the top right example. On PixArt- $\Sigma$  and SDXL-Turbo, our 4-bit results demonstrate noticeably better visual quality than ViDiT-Q's

and MixDQ's W4A8 results.

POSTECH

25/31

\_

24.9

18.2

19.7

20.2

\_

21.3

17.1

17.2

18.1

20.4

21.8

11.2

6.70

16.5

17.2

\_

13.8

16.2

15.6

17.3

\_

21.6

23.7

15.7

11.3

18.0

18.9

\_

21.7

25.9

20.7

21.7

0.265

0.129

0.307

0.261

22.0

26.4

21.0

21.8

25.4

22.4

26.2

23.9

0.453

0.574

0.477

0.502

INT W8A8

INT W8A8

INT W4A4

NVFP W4A4

SDXL

(30 Steps)

TensorRT

Ours

Ours

Ours

20.2

16.6

20.6

18.3

0.591

0.718

0.601

0.640

0.247

0.119

0.288

0.250

## **Integrate with LoRA**

### Integrate with LoRA (Low-Rank Adaptation)

- No re-quantization required for LoRA integration.
- Maintains 16-bit image quality across five LoRA styles, from Realism to Anime and Yarn Art.



W'	$= W + \Delta W = W + AB$ (LoRA)
	$A \in \mathbb{R}^{m \times r_{LORA}}, B \in \mathbb{R}^{r_{LORA} \times n} \cdots r_{LORA} \ll \min(m, n)$
<i>Ŵ</i> =	= $L_1L_2 + R$ (SVDQuant)
A	$r \in L_1^{m  imes r}$ , $L_2 \in R^{r  imes n}$ , R = residual
XW	$\approx XL_1L_2 + Q(X)Q(R)$
$\rightarrow \lambda$	$XW \approx X(L_1L_2 + AB) + Q(X)Q(R)$
$\rightarrow r$	$rank r' = r + r_{LORA}$

Figure 9: Our 4-bit model seamlessly integrates with off-the-shelf LoRAs without requiring requantization. When applying LoRAs, it matches the image quality of the original 16-bit FLUX.1-dev. See Appendix F for the text prompts.

### POSTECH

### **Ablation Study**

### **Ablation Study**

- Naïve quantization and SVD only show severe degradation in 4-bit.
- SVDQuant improves quality by decomposing weights and quantizing only the residual Smoothing
- Low-rank branch absorbs outliers, achieving near-FP16 quality



Prompt: recipe image, angry crab sallad, in salvador dali style photographed by david lachapelle, eerie, rennaisance colors, award winning recipe on white background

Figure 10: Ablation study of SVDQuant on PixArt- $\Sigma$ . The rank of the low-rank branch is 64. Image Reward is measured over 1K samples from MJHQ. Our results significantly outperform the others, achieving the highest image quality by a wide margin.



Figure 5: First 64 singular values of W,  $\hat{W}$ , and R. The first 32 singular values of  $\hat{W}$  exhibit a steep drop, while the remaining values are much more gradual.

#### POSTECH

## **Increasing Rank**

### **Increasing Rank**

- Higher rank(r) improves image quality, with Image Reward increasing from 0.787(r=16) to 0.859(r=64).
- However, model size and latency overhead grow significantly, reaching 11.3% size and 8.8% latency at rank 64.
- Rank 32 was chosen as a middle-ground trade-off in the paper



Prompt: award winning photography of a beautiful medic smiling

Figure 18: Increasing the rank r of the low-rank branch in SVDQuant can enhance image quality, but it also leads to higher parameter and latency overhead.



### Memory save and speedup

Memory save and speedup

- Up to 3.6x model size and 3.5x memory reduction compared to BF16 models.
- **3.0~10.1x latency speedup** across RTX 4090 and RTX 5090.
- Outperforms NF4 (W4A16) in both efficiency and speed, while preserving quality.



Figure 8: SVDQuant reduces the 12B FLUX.1 model size by 3.6× and cuts the 16-bit model's memory usage by 3.5×. With Nunchaku, our INT4 model runs 3.0× faster than the NF4 W4A16 baseline on both desktop and laptop NVIDIA RTX 4090 GPUs. Notably, on the laptop 4090, it achieves a total 10.1× speedup by eliminating CPU offloading. Our NVFP4 model is also 3.1× faster than both BF16 and NF4 on the RTX 5090 GPU.



### Conclusion

- **SVDQuant** : A 4-bit PTQ method for diffusion models
- Introduces a low-rank branch to absorb outliers in weight and activations
- The Nunchaku engine fuses low-rank and 4-bit branches to reduce memory and overhead
- Achieves 3.5x memory reduction and 3.0 speedup on RTX4090/5090



### Limitation

- Lack of fine-grained control over quantization granularity
- Limited generality and stability of activation quantization



### Thank you

