



# RT-Lynx: Putting **GEMM Sparsity** in the Right Place for **Diffusion Models**

**Xing Cong**<sup>1</sup>, Hanlin Tang<sup>2</sup>, Kan Liu<sup>2</sup>, Lan Tao<sup>2</sup>, Lin Qu<sup>2</sup>, **Chenhao Xie**<sup>1\*</sup>

*<sup>1</sup>School of Computer Science and Engineering, Beihang University, Beijing, China*

*<sup>2</sup>Independent Researcher, Beijing, China*



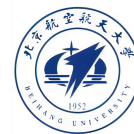
# OUTLINE

**Background & Motivation**

**RT-Lynx**

**Experimental Analysis**

**Conclusion**



# OUTLINE

**Background & Motivation**

**RT-Lynx**

**Experimental Analysis**

**Conclusion**

# 1.1 Diffusion Has Become Mainstream for Image Generation



Table 8. Tensor Shapes of Weights and Activations in the Image and Text Modules of Qwen-Image (where  $X$  denotes activations,  $W$  denotes weights,  $N_{\text{txt}}$  is the text length, and  $N_{\text{img}}$  is the image length). Text in brackets([abbr.]) indicates the abbreviation used in this paper for this type of layer in the model.

Image Module	Image X shape	Image Weight shape	Text Module	Text X shape	Text Weight shape
img_mod.1	(3072)	(3072, 18432)	txt_mod.1	(3072)	(3072, 18432)
attn.to_q [Q]	( $N_{\text{img}}$ , 3072)	(3072, 3072)	attn.add_q_proj	( $N_{\text{txt}}$ , 3072)	(3072, 3072)
attn.to_k [K]	( $N_{\text{img}}$ , 3072)	(3072, 3072)	attn.add_k_proj	( $N_{\text{txt}}$ , 3072)	(3072, 3072)
attn.to_v [V]	( $N_{\text{img}}$ , 3072)	(3072, 3072)	attn.add_v_proj	( $N_{\text{txt}}$ , 3072)	(3072, 3072)
attn.to_out.0 [Out]	( $N_{\text{img}}$ , 3072)	(3072, 3072)	attn.to_add_out	( $N_{\text{txt}}$ , 3072)	(3072, 3072)
img_mlp.net.0.proj [Up]	( $N_{\text{img}}$ , 3072)	(3072, 12288)	txt_mlp.net.0.proj	( $N_{\text{txt}}$ , 3072)	(3072, 12288)
img_mlp.net.2 [Down]	( $N_{\text{img}}$ , 12288)	(12288, 3072)	txt_mlp.net.2	( $N_{\text{txt}}$ , 12288)	(12288, 3072)

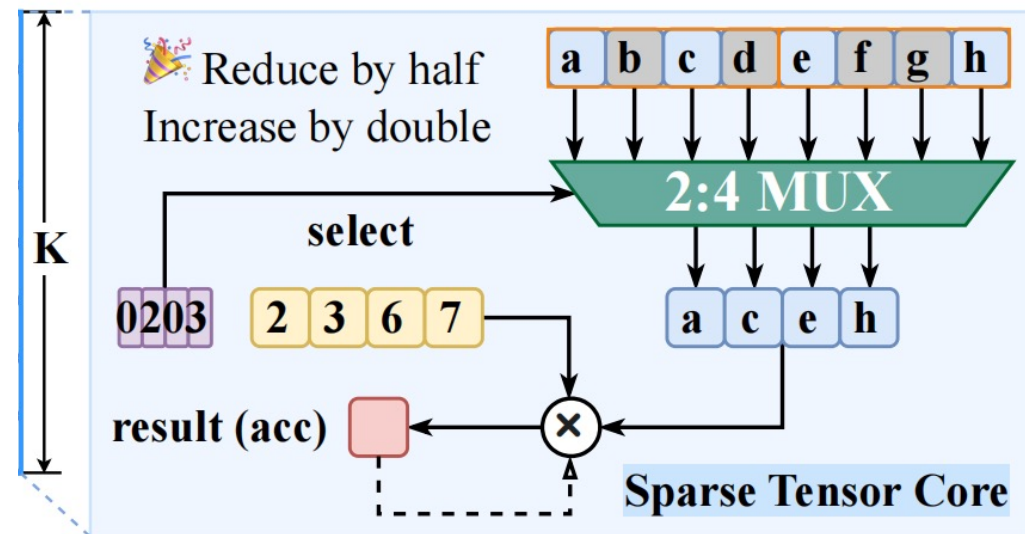
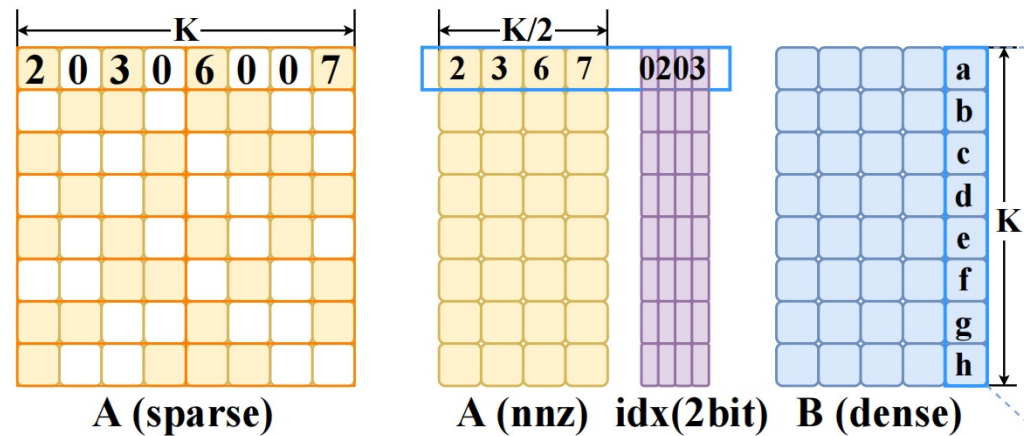
Table 10. Tensor Shapes of Weights and Activations in the Image and Text Streams of FLUX.1-dev (where  $X$  denotes activations,  $W$  denotes weights,  $N_{\text{txt}}$  is the text length, and  $N_{\text{img}}$  is the image length). Text in brackets([abbr.]) indicates the abbreviation used in this paper for this type of layer in the model.

Img Module	Img X shape	Img W shape	Txt Module	Txt X shape	Txt W shape
norm1_a.linear	(3072)	(3072, 18432)	norm1_b.linear	(3072)	(3072, 18432)
attn.a.to_qkv [QKV]	( $N_{\text{img}}$ , 3072)	(3072, 9216)	attn.b.to_qkv	( $N_{\text{txt}}$ , 3072)	(3072, 9216)
attn.a.to_out [Out]	( $N_{\text{img}}$ , 3072)	(3072, 3072)	attn.b.to_out	( $N_{\text{txt}}$ , 3072)	(3072, 3072)
ff.a.0 [Up]	( $N_{\text{img}}$ , 3072)	(3072, 12288)	ff.b.0	( $N_{\text{txt}}$ , 3072)	(3072, 12288)
ff.a.2 [Down]	( $N_{\text{img}}$ , 12288)	(12288, 3072)	ff.b.2	( $N_{\text{txt}}$ , 12288)	(12288, 3072)

- **DiT models are mainly constrained by two factors**
  - Generating one image requires **several or even dozens of denoising steps**
  - Each denoising step is usually **more compute-intensive** than conventional LLM inference

- Distillation:** distills multi-step generation into fewer steps, reducing denoising steps, e.g., from 50 to 8. [1]
- Quantization:** lowers weight and activation precision, e.g., from BF16/FP16 to FP8 or even FP4, to exploit high-FLOPS low-precision cores. [2]
- Caching:** uses similarity between adjacent steps to reduce computation, e.g., merging two steps into one. [3]
- Attention:** reduces redundant computation in attention to lower overall cost. [4]

**However, sparsity, the most common acceleration method in LLMs (up to 2x), remains unused here...**



[1] Zheng, K., Wang, Y., Ma, Q., Chen, H., Zhang, J., Balaji, Y., Chen, J., Liu, M., Zhu, J., & Zhang, Q. (2025). Large Scale Diffusion Distillation via Score-Regularized Continuous-Time Consistency. ArXiv, abs/2510.08431.

[2] Li, M., Lin, Y., Zhang, Z., Cai, T., Li, X., Guo, J., Xie, E., Meng, C., Zhu, J., & Han, S. (2024). SVDQuant: Absorbing Outliers by Low-Rank Components for 4-Bit Diffusion Models. ArXiv, abs/2411.05007.

[3] Liu, F., Zhang, S., Wang, X., Wei, Y., Qiu, H., Zhao, Y., Zhang, Y., Ye, Q., & Wan, F. (2024). Timestep Embedding Tells: It's Time to Cache for Video Diffusion Model. 2025 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).

[4] Zhang, J., Xiang, C., Huang, H., Wei, J., Xi, H., Zhu, J., & Chen, J. (2025). SpargeAttn: Accurate Sparse Attention Accelerating Any Model Inference. ArXiv, abs/2502.18137.

**Full ALL**  
IR(↑): 1.910



**Sparse Weight**  
IR(↑): 1.845



**Sparse Activation**  
IR(↑): 1.867



**Wanda (ICLR'24)**  
IR(↑): 1.813



**RIA (ICLR'24)**  
IR(↑): 1.769



**BaWA (ICML'25)**  
IR(↑): 1.863



**Slim (ICML'25)**  
IR(↑): 1.873



**Ours (1.5x Faster)**  
IR(↑): 1.925

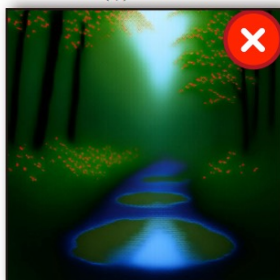


Prompt: *A capybara wearing a suit holding a sign that reads ICLM leads smart computing with strong research. Anime cute style, high quality, 4k*

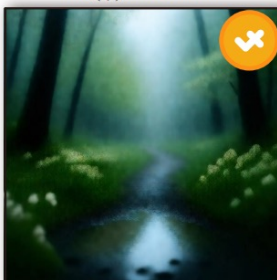
IR(↑): 1.679



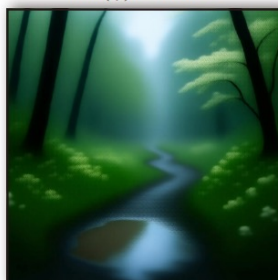
IR(↑): -2.156



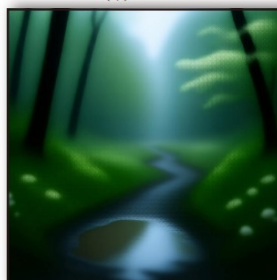
IR(↑): -0.4074



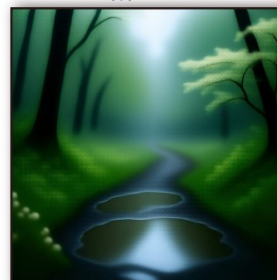
IR(↑): -1.284



IR(↑): -1.839



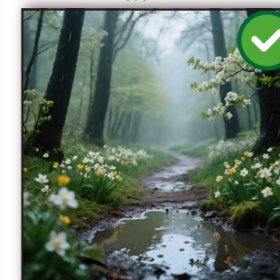
IR(↑): -1.482



IR(↑): 1.611



IR(↑): 1.763



Prompt: *Rainy woods, heavy rain, spring, flowers, day, trails, puddles of water*

**Experiment:** Comparison of sparsification strategies on Qwen-Image. Both weight and activation sparsity use Top-K selection to impose 2:4 semi-structured sparsity.

- **Sparse Weight** severely damages the original model's image-generation capability
- **Sparse Activation** shows some robustness to sparse operations, preserving the basic elements and visual quality of generated images



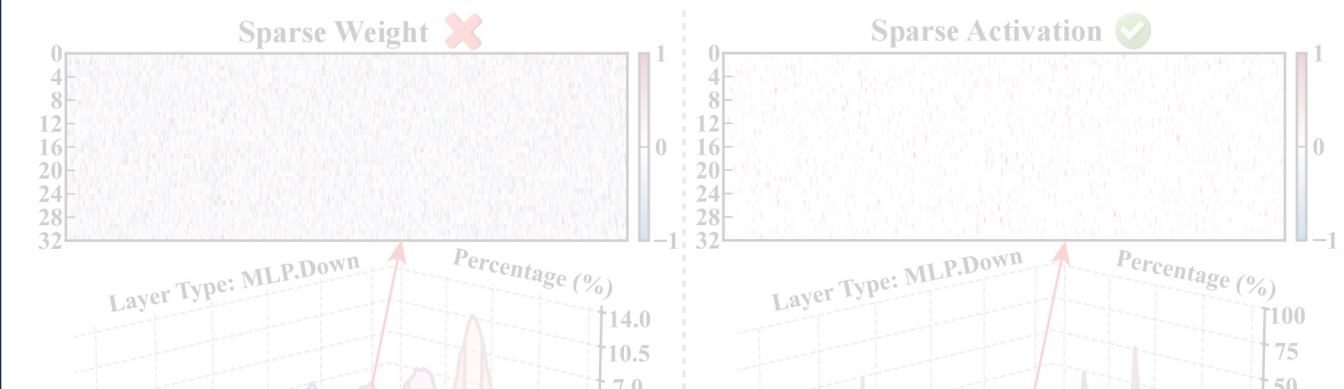
# 1.4 Why Sparse Activation is Better?

Background Motivation

RT-Lynx

Experimental Analysis

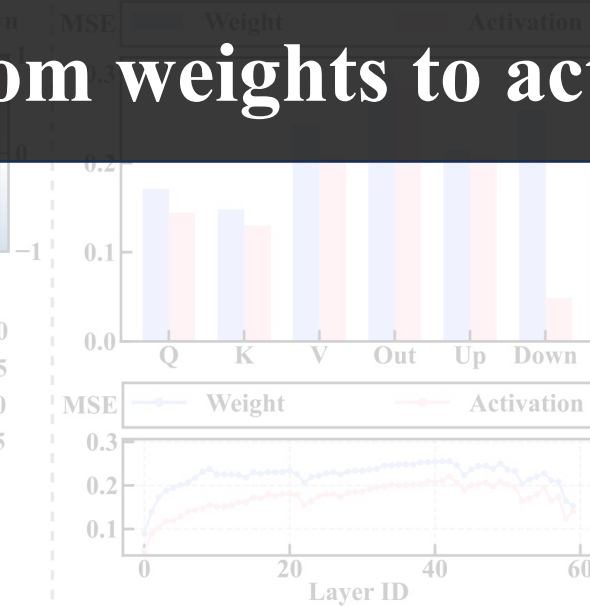
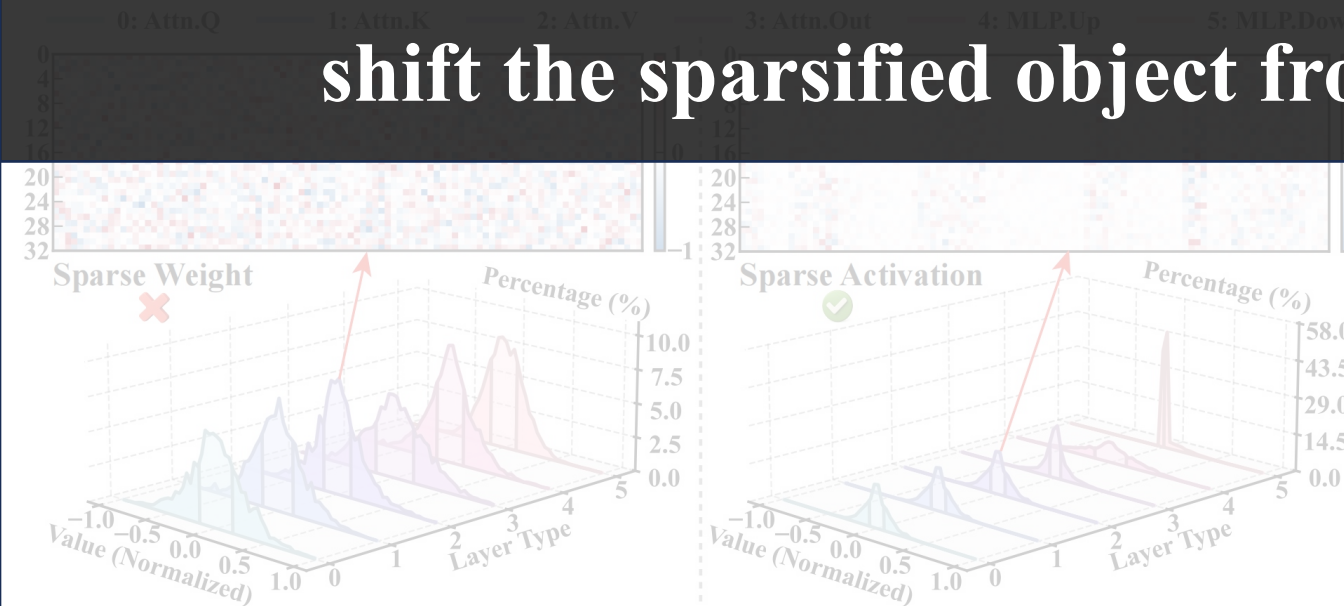
Conclusion



➤ **Experiment-1:** Distribution and error comparison of weight vs. activation sparsity in Qwen-Image MLP-down layers under Top-K 2:4 sparsification.

**We propose a new sparsity paradigm:**

**shift the sparsified object from weights to activations**



➤ **Experiment-2:** Distribution and error comparison of weight vs. activation sparsity across Qwen-Image Linear layers under Top-K 2:4 sparsification.

**So:** Activations in DiT are numerically sparser and less sensitive to errors from sparse operations!

### ➤ Accuracy

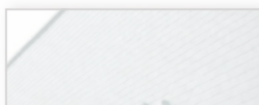
**Full ALL**  
FID(↓): 31.15  
IR(↑): 1.172



**Sparse Weight**  
FID(↓): 66.91  
IR(↑): -0.2159



**Sparse Activation**  
FID(↓): 48.59  
IR(↑): 0.4724



**Accuracy compensation and an efficient execution pipeline are needed to preserve quality and achieve speedup**

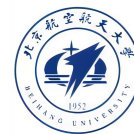
- Sparse Activation still introduces large **errors on some prompts**.

### ➤ Performance

*Table 5. Proportion of online activation sparsification overhead in total sparse execution for a typical DiT inference pipeline across methods on H20 GPUs. Gray rows indicate matrix sizes used in Owen-Image.*

$M=N$	$K$	PyTorch	CUTLASS	cuSparseLi
4096	3072	35.08%	26.98%	44.61%
8192	3072	23.14%	17.01%	26.99%
4096	12288	39.29%	28.10%	43.37%
8192	12288	23.14%	17.01%	26.99%

- Sparse Activation requires online sparsification and format conversion while also enabling efficient pipelined GEMM execution. This remains largely unexplored, **making speedup difficult to achieve**.



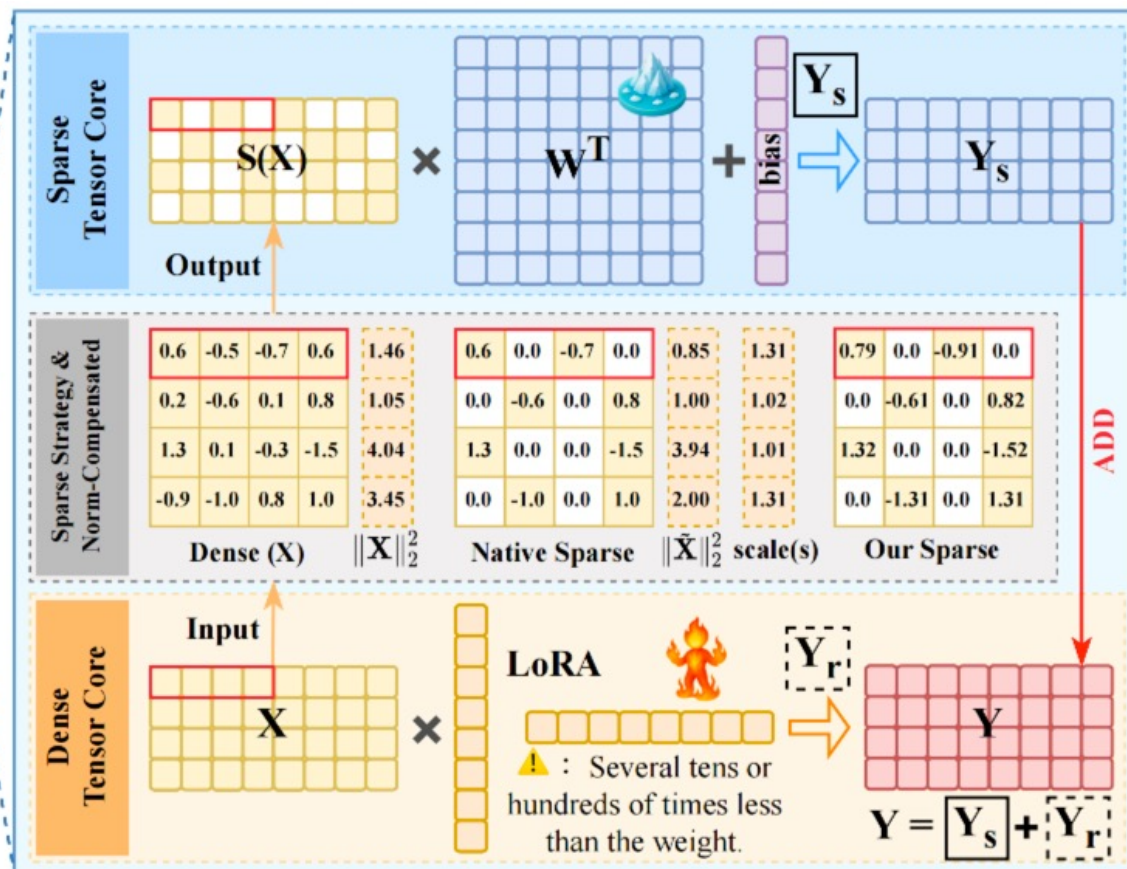
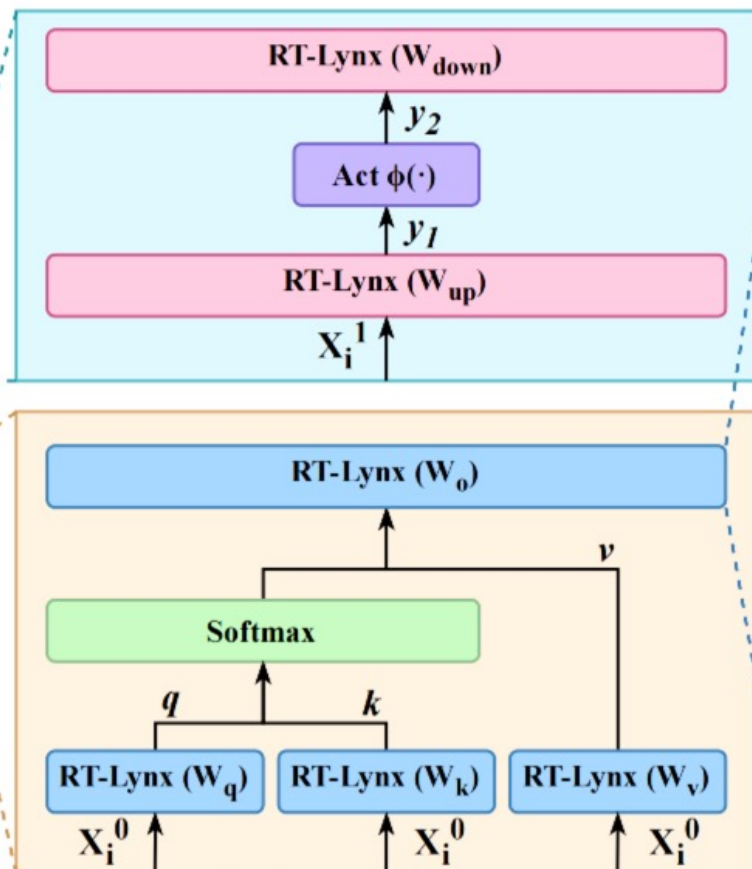
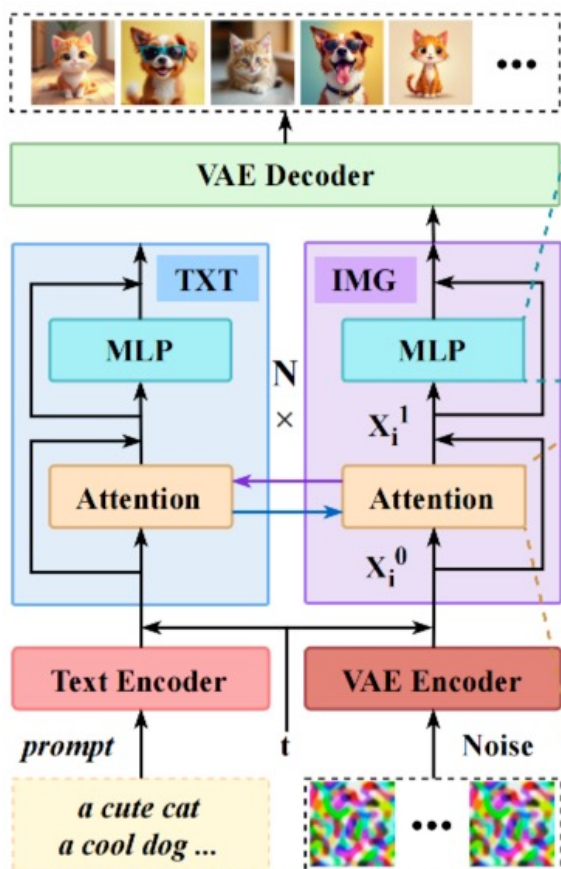
# OUTLINE

**Background & Motivation**

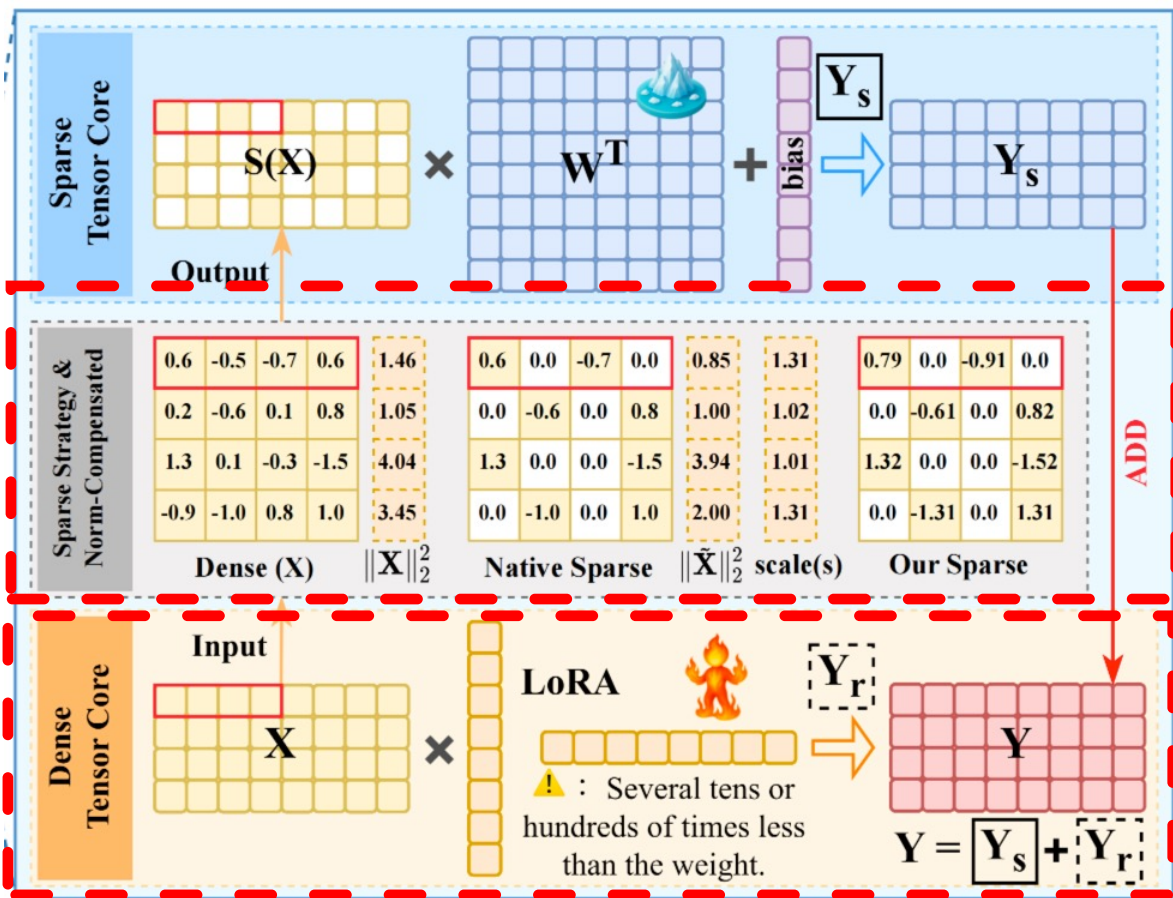
**RT-Lynx**

**Experimental Analysis**

**Conclusion**



**Replace all linear layers in the model with RT-Lynx (Plug-and-Play).**



■ Pruning **reduces the overall magnitude.**

$$S(\mathbf{X}) = s \cdot \tilde{\mathbf{X}}, \quad s = \sqrt{\frac{\|\mathbf{X}\|_2^2}{\|\tilde{\mathbf{X}}\|_2^2 + \epsilon}}$$

✓ **Rescale sparse activations:** so their L2 norm aligns with the original dense activations.

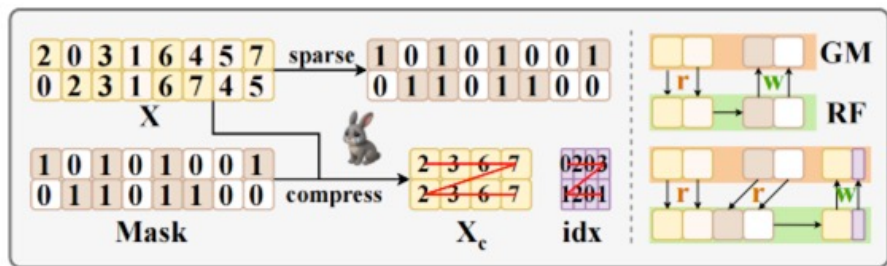
■ Although activations are close to zero, they still **contain fine-grained information;**

$$\mathbf{Y} = \mathbf{Y}_s + \mathbf{Y}_r = S(\mathbf{X}) \cdot \mathbf{W}^\top + \mathbf{X} \cdot (\mathbf{L}_A \mathbf{L}_B)^\top$$

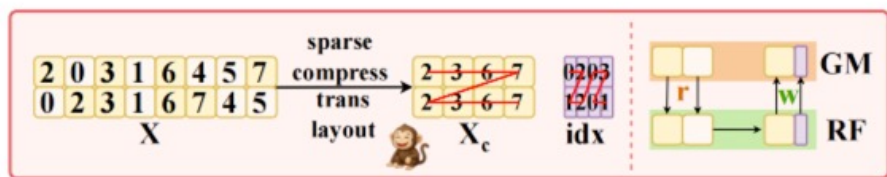
$$Loss = \|\mathbf{X} \cdot \mathbf{W}^\top - (S(\mathbf{X}) \cdot \mathbf{W}^\top + \mathbf{X} \cdot (\mathbf{L}_A \mathbf{L}_B)^\top)\|^2$$

✓ A lightweight **LoRA branch** recovers sparsification residuals by aligning sparse outputs with dense outputs through self-distillation.

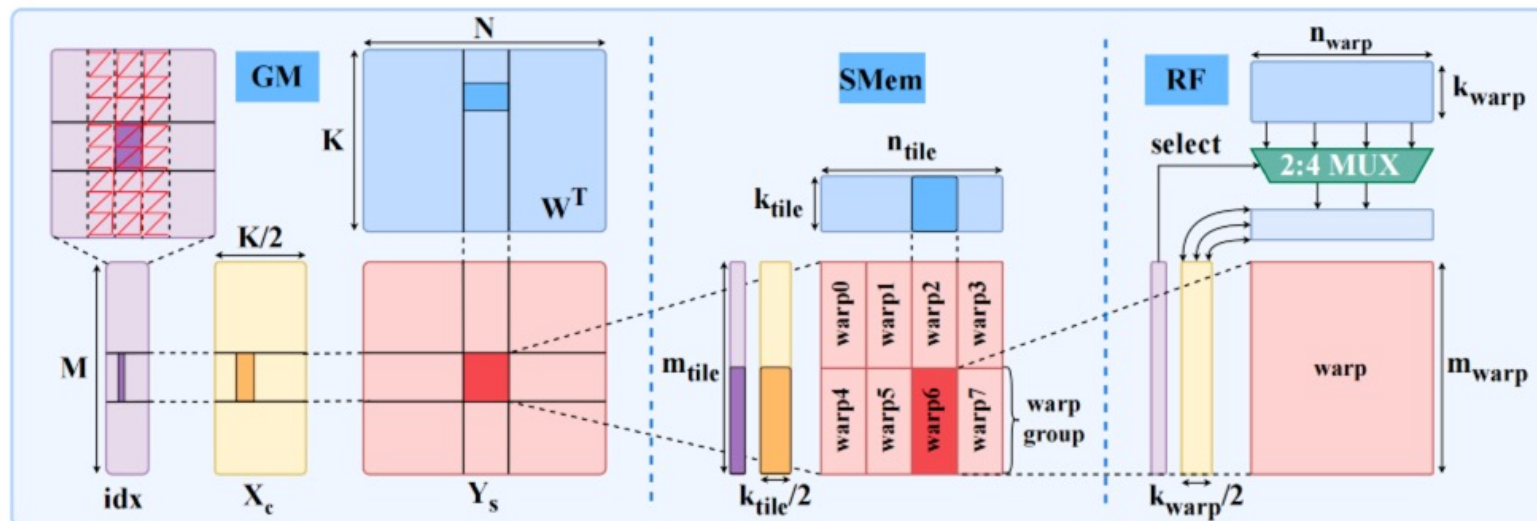
✓ For sensitive layers that cannot be fully compensated, sparsification is **skipped** to balance accuracy and speed.



(a) Other : sparse kernel and compress kernel



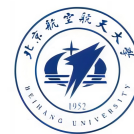
(b) RT-Lynx : fused sparse kernel



(c) RT-Lynx : compute sparse-dense matrix multiplication with tiling and pipeline

■ The online computation of sparse matrix multiplication using existing kernels **takes too long**.

- ✓ **Fuse pattern** determination, Top-K, and compression into a single CUDA kernel, directly generating a 2:4 SpTC-compatible layout at the register level.
- ✓ Build a **execution pipeline** for sparse computation + dense LoRA computation, combining  $Y_r$  and  $Y_s$  on chip. Avoid materializing intermediate LoRA results, reducing synchronization and memory access overhead.



# OUTLINE

**Background & Motivation**

**RT-Lynx**

**Experimental Analysis**

**Conclusion**

### ■ Model

- Qwen-Image
- FLUX.1-Dev
- Z-Image
- Qwen-Image-2512

### ■ Datasets (5000 prompt)

- MJHQ
- sDCI

### ■ Baseline

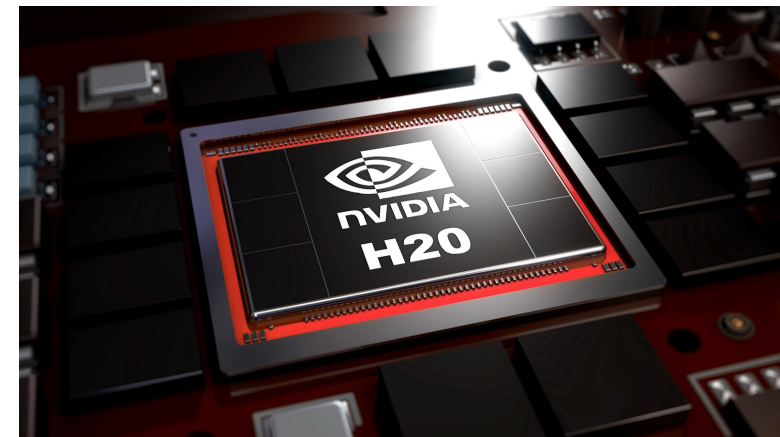
- Wanda
- RIA
- BaWA
- Slim

### ■ Metrics

- FID
- Image Reward
- CLIP-IQA
- CLIP-Score

### ■ Details

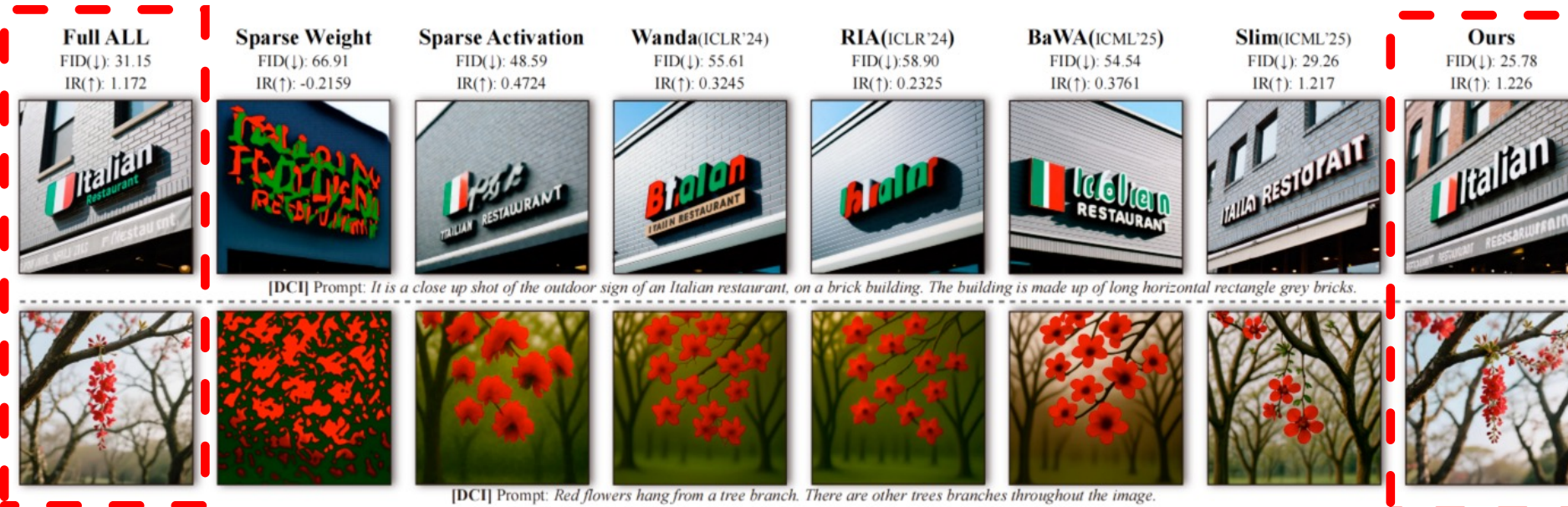
- H20 GPUs
- Driver: 580.82.07
- CUDA 13.0



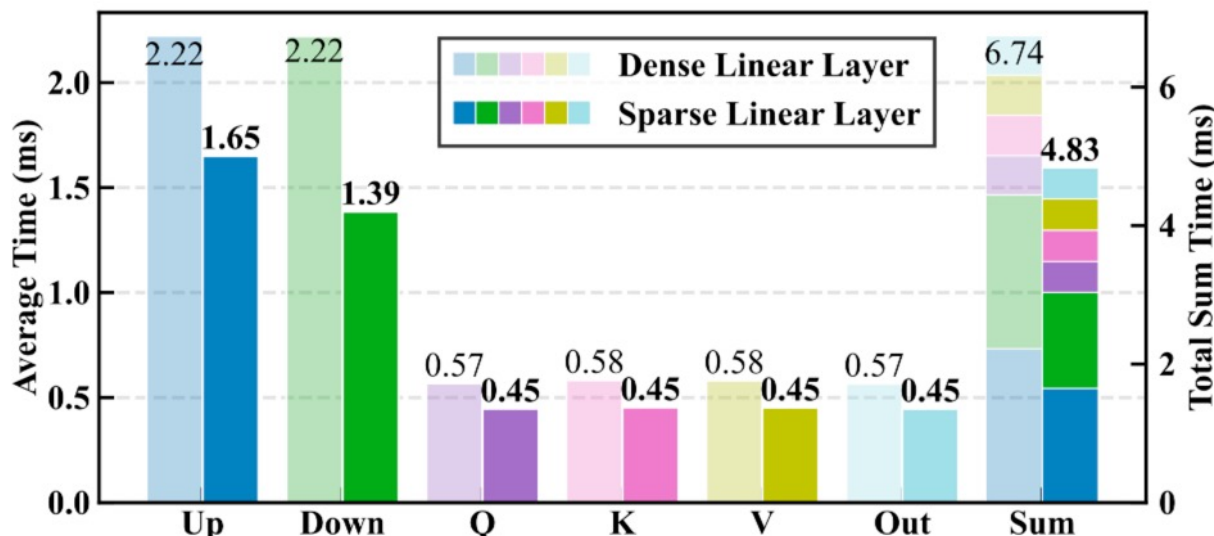
Method	MJHQ				sDCI			
	FID(↓)	IR(↑)	C.SCR(↑)	C.IQA(↑)	FID(↓)	IR(↑)	C.SCR(↑)	C.IQA(↑)
Full	21.98	1.219	27.12	0.9317	31.15	1.172	26.49	0.9492
<del>Sparse Weight</del>	<del>51.63</del>	<del>-0.1605</del>	<del>24.22</del>	<del>0.5783</del>	<del>66.91</del>	<del>-0.2159</del>	<del>24.93</del>	<del>0.5792</del>
<del>Sparse Activation</del>	<del>35.85</del>	<del>0.5994</del>	<del>26.28</del>	<del>0.7473</del>	<del>48.59</del>	<del>0.4724</del>	<del>26.13</del>	<del>0.7550</del>
Wanda (ICLR'24)	40.81	0.5356	26.06	0.7344	55.61	0.3245	26.01	0.7268
RIA (ICLR'24)	43.02	0.4627	25.90	0.7150	58.90	0.2325	25.90	0.7124
BaWA (ICML'25)	39.68	0.5888	26.15	0.7660	54.54	0.3761	26.15	0.7692
Slim (ICML'25)	22.25	1.278	27.18	0.9242	29.26	1.217	26.23	0.9350
RT-Lynx (Ours)	<b>21.25</b>	<b>1.304</b>	<b>27.33</b>	<b>0.9263</b>	<b>25.78</b>	<b>1.226</b>	<b>26.61</b>	<b>0.9366</b>

(1) Naive weight sparsity severely degrades image quality, while **activation sparsity provides a stronger starting point.**

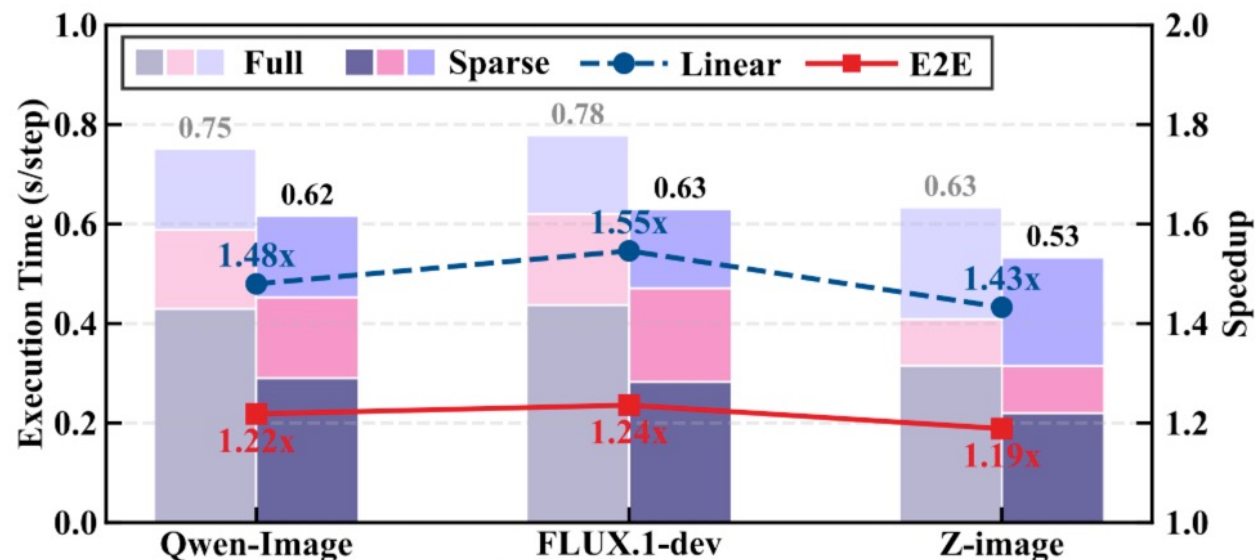
(2) RT-Lynx further restores quality, achieving the best FID and IR among sparse methods on both MJHQ and sDCI.



$M=N$	$K$	Pytorch-GEMM	Pytorch-SpMM	CUTLASS	cuSparseLt	RT-Lynx-Kernel	Sparse-Cost(%)
2048	3072	0.199	0.221 (0.90 $\times$ )	0.216 (0.92 $\times$ )	0.282 (0.71 $\times$ )	0.135 ( <b>1.47<math>\times</math></b> )	<b>8.33%</b>
4096	3072	0.781	0.715 (1.09 $\times$ )	0.616 (1.26 $\times$ )	0.709 (1.10 $\times$ )	0.465 ( <b>1.68<math>\times</math></b> )	<b>4.60%</b>
8192	3072	2.994	2.087 (1.43 $\times$ )	2.044 (1.46 $\times$ )	2.066 (1.50 $\times$ )	1.802 ( <b>1.66<math>\times</math></b> )	<b>2.28%</b>
2048	12288	0.781	1.259 (0.62 $\times$ )	0.801 (0.97 $\times$ )	1.137 (0.69 $\times$ )	0.454 ( <b>1.72<math>\times</math></b> )	<b>9.39%</b>
4096	12288	3.099	2.648 (1.17 $\times$ )	2.287 (1.35 $\times$ )	2.669 (1.16 $\times$ )	1.652 ( <b>1.88<math>\times</math></b> )	<b>4.83%</b>
8192	12288	11.95	8.217 (1.45 $\times$ )	7.497 (1.59 $\times$ )	8.202 (1.45 $\times$ )	6.754 ( <b>1.77<math>\times</math></b> )	<b>2.37%</b>



(a) Single-Step Inference Average Latency of Different Linear Layers



(b) End-to-End Inference Average Latency Analysis

**Kernel:** RT-Lynx accelerates GEMM with low sparsification overhead, achieving up to **1.88 $\times$  kernel speedup**.

**E2E:** The optimized sparse Linear layers reduce single-step latency and bring consistent end-to-end gains.

Model	Method	MJHQ				sDCI			
		FID(↓)	IR(↑)	C.SCR(↑)	C.IQA(↑)	FID(↓)	IR(↑)	C.SCR(↑)	C.IQA(↑)
Qwen-Image	Full	21.98	1.219	27.12	0.9317	31.15	1.172	26.49	0.9492
	SA-Native	35.85	0.5994	26.28	0.7473	48.59	0.4724	26.13	0.7550
	SA-NC	25.28	0.9393	26.76	0.8349	37.56	0.8399	26.34	0.8321
	SA-NC-LoRA	<b>21.25</b>	<b>1.304</b>	<b>27.33</b>	<b>0.9263</b>	<b>25.78</b>	<b>1.226</b>	<b>26.61</b>	<b>0.9366</b>
Qwen-Image (2512)	Full	20.45	1.290	26.45	0.9569	24.42	1.130	25.91	0.9435
	SA-Native	39.13	0.6575	25.60	0.7379	51.41	0.3205	<b>26.48</b>	0.6593
	SA-NC	26.46	1.064	26.08	0.8891	33.29	0.8747	26.32	0.8251
	SA-NC-LoRA	<b>20.84</b>	<b>1.302</b>	<b>26.54</b>	<b>0.9484</b>	<b>24.10</b>	<b>1.132</b>	26.07	<b>0.9285</b>
Flux.1-dev	Full	22.05	1.010	26.51	0.9401	25.96	1.074	26.01	0.9487
	SA-Native	49.87	0.4520	24.83	0.7881	52.88	0.590	25.74	0.7471
	SA-NC	38.67	0.7823	25.56	0.8565	39.67	0.8704	25.86	0.8340
	SA-NC-LoRA	22.61	0.9783	26.30	0.9341	26.35	1.050	25.90	0.9415
	SA-NC-LoRA-SL	<b>21.17</b>	<b>1.011</b>	<b>26.42</b>	<b>0.9381</b>	<b>24.41</b>	<b>1.091</b>	<b>25.90</b>	<b>0.9445</b>
Z-image	Full	25.70	0.9928	25.57	0.9307	25.40	0.9974	25.91	0.9467
	SA-Native	42.65	0.5982	25.42	0.7917	36.18	0.7306	<b>26.72</b>	0.7508
	SA-NC	31.65	0.8054	<b>25.61</b>	0.8525	28.48	0.8818	26.45	0.8504
	SA-NC-LoRA	27.39	0.9292	25.48	0.9016	25.98	0.9788	26.08	0.9137
	SA-NC-LoRA-SL	<b>26.17</b>	<b>0.9673</b>	25.39	<b>0.9250</b>	<b>24.93</b>	<b>0.9803</b>	26.22	<b>0.9396</b>



**Each component contributes to stable sparse inference:** norm compensation reduces error, LoRA restores output quality, and sensitive-layer skipping protects fragile layers.

Together, they deliver the best visual quality and metric performance across models.

Model	Method	MJHQ				sDCI			
		FID(↓)	IR(↑)	C.SCR(↑)	C.IQA(↑)	FID(↓)	IR(↑)	C.SCR(↑)	C.IQA(↑)
Distillation (Z-Image)	Baseline(step-50)	22.31	0.8790	26.25	0.9046	18.60	0.9545	26.23	0.9196
	<i>RT-Lynx</i>	22.68	0.9301	26.25	0.9058	18.80	1.006	26.25	0.9167
	Turbo(step-8)	25.70	0.9928	25.57	0.9307	25.40	0.9974	25.91	0.9467
	<b>RT-Lynx+Turbo</b>	<b>26.17</b>	<b>0.9673</b>	<b>25.39</b>	<b>0.9250</b>	<b>24.93</b>	<b>0.9803</b>	<b>26.22</b>	<b>0.9396</b>
Quant (Qwen-Image)	Baseline	21.98	1.219	27.12	0.9317	31.15	1.172	26.49	0.9492
	<i>RT-Lynx</i>	21.25	1.304	27.33	0.9263	25.78	1.226	26.61	0.9366
	W8A8	21.78	1.220	27.10	0.9298	30.59	1.173	26.48	0.9474
	<b>RT-Lynx+W8A8</b>	<b>21.43</b>	<b>1.299</b>	<b>27.33</b>	<b>0.9250</b>	<b>25.92</b>	<b>1.227</b>	<b>26.65</b>	<b>0.9341</b>
Cache (FLUX.1-dev)	Baseline	22.05	1.010	26.51	0.9401	25.96	1.074	26.01	0.9487
	<i>RT-Lynx</i>	21.17	1.011	26.42	0.9381	24.41	1.091	25.90	0.9445
	Teacache( $l=0.6$ )	20.67	0.937	25.83	0.9373	25.76	1.024	25.46	0.9487
	<b>RT-Lynx+Teacache</b>	<b>20.57</b>	<b>0.941</b>	<b>25.73</b>	<b>0.9414</b>	<b>24.98</b>	<b>1.026</b>	<b>25.32</b>	<b>0.9443</b>
Attention (Qwen-Image-2512)	Baseline	20.45	1.290	26.45	0.9569	24.42	1.130	25.91	0.9435
	<i>RT-Lynx</i>	20.84	1.302	26.54	0.9484	24.10	1.132	26.07	0.9285
	SparseAttn( $k=0.5$ )	21.17	1.192	26.18	0.9557	25.46	1.067	26.16	0.9390
	<b>RT-Lynx+SparseAttn</b>	<b>21.39</b>	<b>1.212</b>	<b>26.37</b>	<b>0.9420</b>	<b>25.22</b>	<b>1.075</b>	<b>26.42</b>	<b>0.9267</b>

Model	Method	End2End(s)
Distillation (Z-Image)	Baseline(step-50)	49.44
	<i>RT-Lynx</i>	40.76 (1.21×)
	Turbo(step-8)	4.99 (9.91×)
	<b>RT-Lynx+Turbo</b>	<b>4.17 (11.86×)</b>
Quant (Qwen-Image)	Baseline	60.66
	<i>RT-Lynx</i>	50.60 (1.20×)
	W8A8	54.45 (1.11×)
	<b>RT-Lynx+W8A8</b>	<b>46.04 (1.32×)</b>
Cache (FLUX.1-dev)	Baseline	77.99
	<i>RT-Lynx</i>	63.00 (1.24×)
	TeaCache( $l=0.6$ )	29.56 (2.64×)
	<b>RT-Lynx+TeaCache</b>	<b>24.89 (3.13×)</b>
Attention (Qwen-Image-2512)	Baseline	60.42
	<i>RT-Lynx</i>	49.35 (1.22×)
	SparseAttn( $k=0.5$ )	54.49 (1.11×)
	<b>RT-Lynx+SparseAttn</b>	<b>44.35 (1.36×)</b>

(1) RT-Lynx can be seamlessly combined with Distillation, Quantization, Cache, and Sparse Attention.

(2) These combinations further **improve end-to-end speedup**, reaching up to **11.86×**, **while preserving competitive** FID, IR, C.SCR, and C.IQA scores.



# OUTLINE

**Background & Motivation**

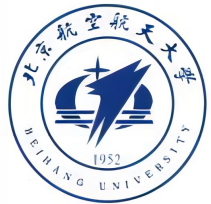
**RT-Lynx**

**Experimental Analysis**

**Conclusion**

- We explore **shifting the sparsity paradigm from weights to activations** to enable efficient inference acceleration for DiT models.
- We **reveal that activations exhibit inherent sparsity** and are more robust than weights under N:M sparsity.
- We propose RT-Lynx, introducing norm compensation, a LoRA branch, and a layer-skipping mechanism to **reduce sparsification error**.
- We design **an efficient CUDA kernel** that fuses online activation sparsity generation with sparse computation.
- Across multiple diffusion models, RT-Lynx achieves lossless acceleration with significant inference speedups.





**北京航空航天大学**  
BEIHANG UNIVERSITY



**Thanks!**



Thank you, and have a  
great day!

Email: [congxingcs@163.com](mailto:congxingcs@163.com)