

ArXiv '24

PixArt- Σ : Weak-to-Strong Training of Diffusion Transformer for 4K Text-to-Image Generation

Junsong Chen, Chongjian Ge, Enze Xie, Yue Wu, Lewei Yao, Xiaozhe Ren, Zhongdao Wang, Ping Luo, Huchuan Lu, Zhenguo Li

Jeonghoon Park, happypjh2001@unist.ac.kr
Undergraduate Research Intern,
Ubiquitous Artificial Intelligence Lab,
Department of Computer Science and Engineering,
Ulsan National Institute of Science and Technology

Reason to choose this paper

For research...

MobileDiffusion

- ◆ MobileDiffusion: Instant Text-to-Image Generation on Mobile Devices
- ◆ 512 x 512 Resolution

On-device + 4K diffusion

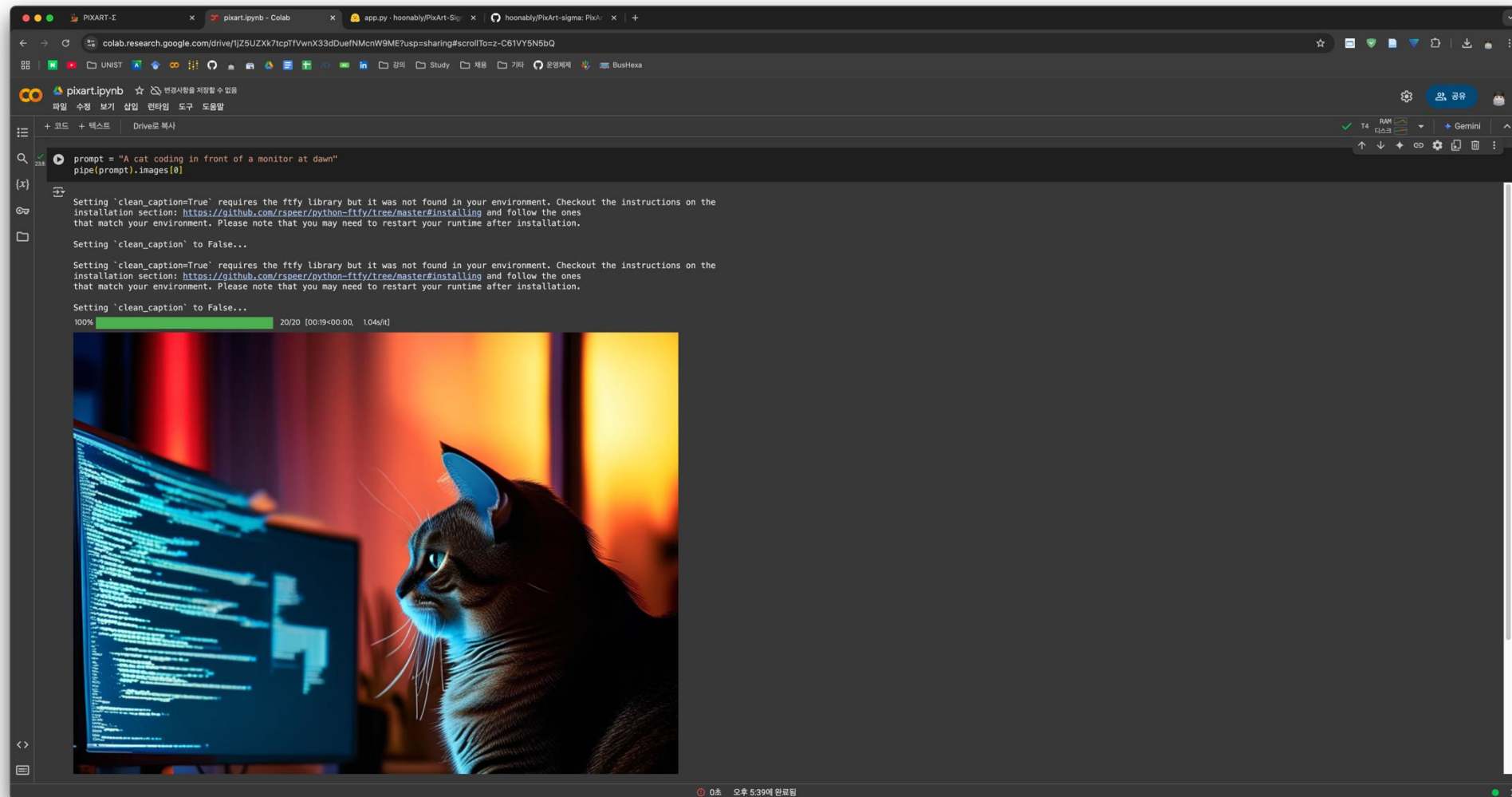
- ◆ No papers on implementing 4K at on-device
- ◆ Decide on a topic for my research

4K diffusion paper

- ◆ Exploring 4K diffusion technology
- ◆ Latest paper

Preview : Result of a run in Colab T4

PixArtAlphaPipeline.from_pretrained("PixArt-alpha/**PixArt-Sigma-XL-2-1024-MS**", torch_dtype=torch.float16)



Previous works

PixArt- α : Fast training of diffusion transformer for photorealistic text-to-image synthesis (ICLR, 2024 Spotlight)

Chen, J., Yu, J., Ge, C., Yao, L., Xie, E., Wu, Y., Wang, Z., Kwok, J., Luo, P., Lu, H., Li, Z.

➔ First Transformer-based Diffusion Model (DiT) capable of generating up to 1024×1024 resolution

Stable Diffusion : High-resolution image synthesis with latent diffusion models (CVPR, 2022)

Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.

➔ Utilizes the Latent Diffusion Model (LDM) structure to generate high-resolution images beyond 1024×1024

DALL-E 3 (OpenAI, 2023)

➔ Utilizes GPT-4-based text comprehension to more accurately reflect prompts

Data Analysis : Higher Aesthetic and higher Resolution

Effective training with limited data

	Data	Data Resolution
Internal- α	14M	Only 256 ~ 1K
Internal-Σ	33M	1K~4K (33M) real photo 4K (8M)
SD v1.5 (open-source)	2B	512x512, 768x768

Data Analysis : Higher Aesthetic and higher Resolution

However, well scored!

Models	#Params (B)	FID ↓	CLIP-Score ↑
Stable 1.5	0.9	17.03	0.2748
Stable Turbo	3.1	10.91	0.2804
Stable XL	2.6	7.38	0.2913
Stable Cascade	5.1	9.96	0.2839
Playground-V2.0	2.6	8.68	0.2885
Playground-V2.5	2.6	7.64	0.2871
PIXART- α	0.6	8.65	0.2787
PIXART- Σ	0.6	8.23	0.2797

Data Analysis : Enhanced caption accuracy

PixArt- α (LLaVa) -> certain hallucination problem

PixArt- Σ (Share-Captioner) -> generate detailed and correct captions -> augmenting the collected raw prompts



	LLaVA Hallucinations	<p>The image features a large, ornate church with a tall, pointed roof and a large stained-glass window. The church has a white and gray color scheme. The style of the church is Gothic, featuring a pointed roof and the intricate details of the stained-glass window. The presence of statues further emphasizes the grand and historical nature of the structure.</p>
	Share-Captioner Correctness	<p>The image captures the grandeur of a cathedral, painted against a backdrop of a clear blue sky. The entrance to the cathedral is flanked by statues of saints, standing as silent guardians. The photo, taken from a low angle, shows the lush green trees in the foreground.</p>
	LLaVA Hallucinations	<p>The image features a woman and a man sitting on a brick walkway near a body of water, which could be a river or a lake. They are both wearing head coverings, and the woman is holding a handbag. The scene is set during the day, with the sun shining brightly, creating a warm and inviting atmosphere. The style of the image is a black and white photo, which adds a timeless and classic feel to the scene.</p>
	Share-Captioner Correctness	<p>The image captures a serene scene at a harbor. Two individuals are seated on a bench, their backs to the camera, engrossed in the view of the water. The water, a deep shade of blue, is dotted with boats of various sizes and colors, including a white boat with a green stripe and a red boat. The sky above is a light blue.</p>

Fig. 5: Comparative illustration of hallucinations: Contrasting differences in hallucination occurrences between LLaVA and Share-Captioner, with **red** indicating hallucinations and **green** denoting correctness.

LLaVa

“Visual Instruction Tuning” (2023)

➔ a study to create a visual version (Vision-Language Model, VLM) of GPT-4.

- Based on: CLIP + LLaMA (Language Model)
- Purpose: Multimodal model to view images and perform “description, question-answer (Q&A), summarization, etc.”
- Features:
 - Utilizes CLIP to convert images into linguistic representations.
 - Large Language Model (LLaMA) to generate text.
 - Performs a similar role to the traditional GPT-4V.
 - However, it can be less accurate and potentially lacks fine-grained information.

Share-Captioner

“ShareGPT4V: Improving Large Multi-Modal Models with Better Captions” (2023)

➔ Share-Captioner is a model to overcome the limitations of LLaVA and generate more sophisticated captions.

- Based on: GPT-4V (GPT-4 with Vision)
- Purpose: Generate more accurate and detailed image captions
- Features:
 - Utilizes GPT-4V to generate more accurate and detailed descriptions.
 - Longer sentences, more detail than LLaVA.

Data Analysis : Increased caption length

Internal- α : ≤ 120 tokens

Internal- Σ : ≤ 300 tokens

Share-Captioner(60%) + raw(40%) \rightarrow reduce potential biases

(Not using raw data can introduce bias!)

Table 1: Statistics of noun concepts for different datasets. VN: valid distinct nouns (appearing more than 10 times); DN: total distinct nouns; **Average**: average noun count per image; **ACL**: Average Caption length.

Dataset	Volume	Caption	VN/DN	Total Noun	ACL	Average
Internal- α	14M	Raw	187K/931K	175M	25	11.7/Img
Internal- α	14M	LLaVA	28K/215K	536M	98	29.3/Img
Internal- α	14M	Share-Captioner	51K/420K	815M	184	54.4/Img
Internal- Σ	33M	Raw	294K/1512K	485M	35	14.4/Img
Internal- Σ	33M	Share-Captioner	77K/714K	1804M	180	53.6/Img
4K- Σ	2.3M	Share-Captioner	24K/96K	115M	163	49.5/Img

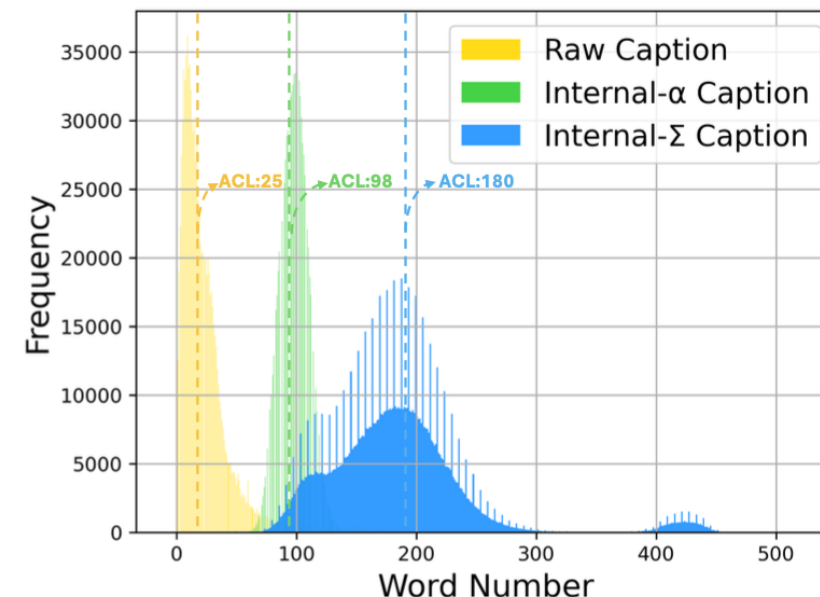


Fig. 6: Histogram Visualization of the Caption Length. We randomly select 1M captions from the **raw captions**, **Internal- α** , and **Internal- Σ** to draw the corresponding histogram. *ACL* denotes the average caption length.

Data Analysis : High-Quality Evaluation Dataset

Most SoTA T2I models chose **MSCOCO** (MobileDiffusion too)

SoTA : State of the Art

-> Not enough to evaluate aesthetics and text-image alignment





Image	Prompt	Image	Prompt
	A red apple sitting on a wooden table, remote control aerial photography.		A photographic work capturing a polar bear walking through icy and snowy terrain.
	A serene beach with palm trees, turquoise water, and a hammock between two trees, star trail.		A bird known for its distinctive blue and orange plumage. The kingfisher is perched on a branch, its body angled slightly to the left as if poised to take flight at any moment.

Fig. 12: Samples in our proposed High-Quality Evaluation Dataset. The evaluation dataset presented in this paper contains samples of superior visual quality compared to those in COCO-30K.

This paper use 30,000 high quality dataset

Efficient DiT Design : Previous problems

in PixArt- α

Self-Attention computation increases proportional to **the square of the number of tokens**

→ **$O(N^2)$**

4K resolution needs higher number of tokens → **Model execution = slow**

Memory usage spikes when generating 4K → **GPU costs = increase**

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{Q \cdot f_c(K)^T}{\sqrt{d_k}} \right) f_c(V)$$

Efficient DiT Design : Key-Value Token Compression

- Compress Key (K) and Value (V) using **Group Convolution on Stride 2**

➔ reduce the number of tokens by $N \rightarrow N/R^2$

- Using $1 \leq R \leq 4$ without losing too much accuracy

➔ reduce computation by about 34%

$$\text{from } O(N^2) \text{ to } O\left(\frac{N^2}{R^2}\right)$$

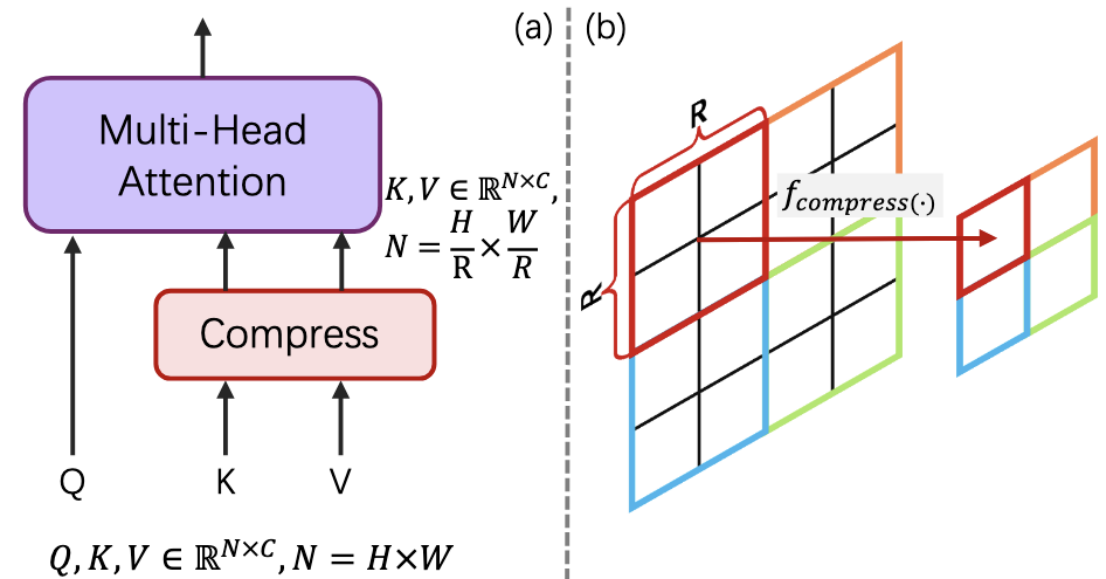
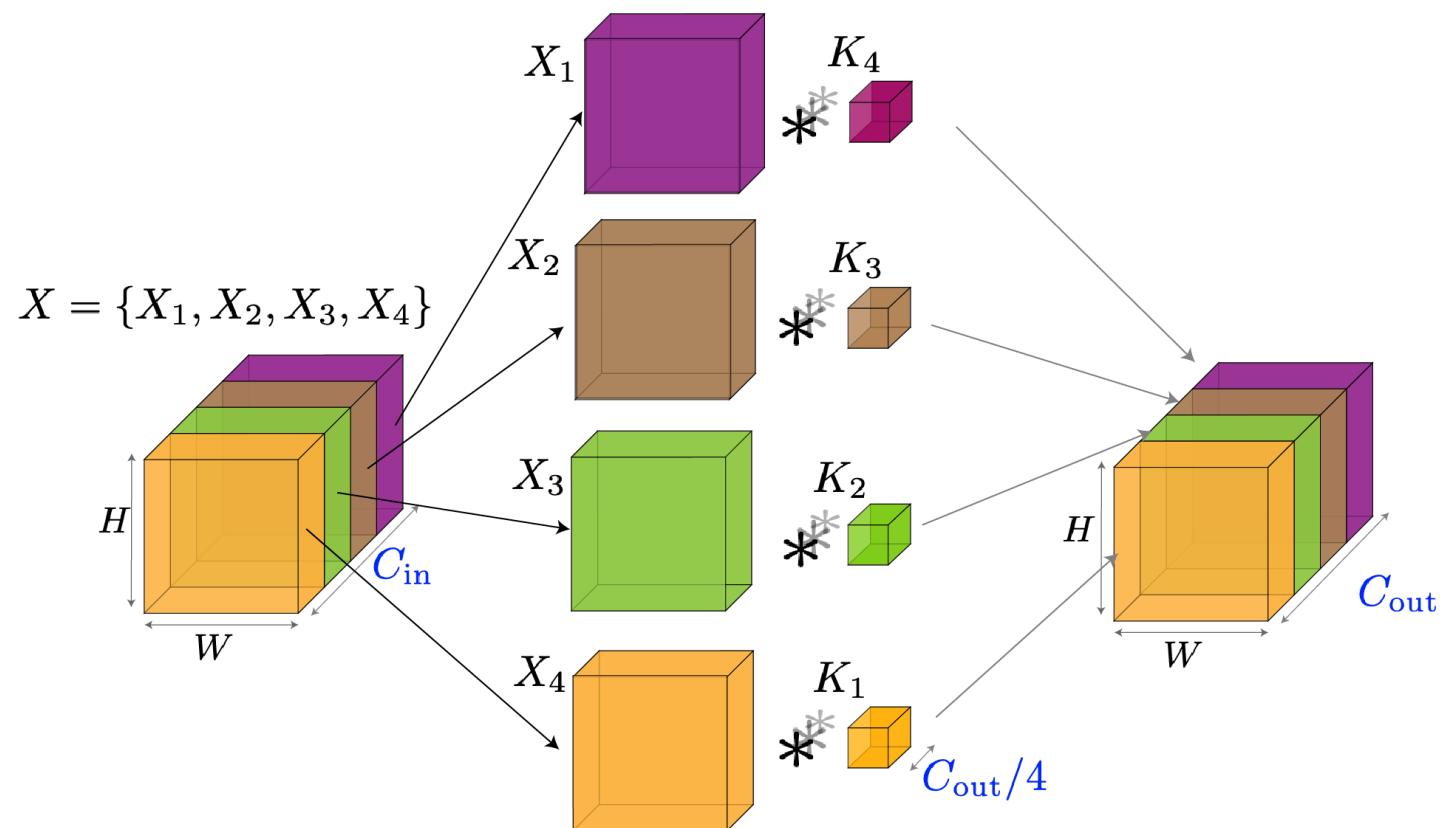


Fig. 7: Design of KV Token Compression. We merge KV tokens in spatial space to reduce the computation complexity.

Group Convolution?

Divide input channels into groups and perform convolution on each independently

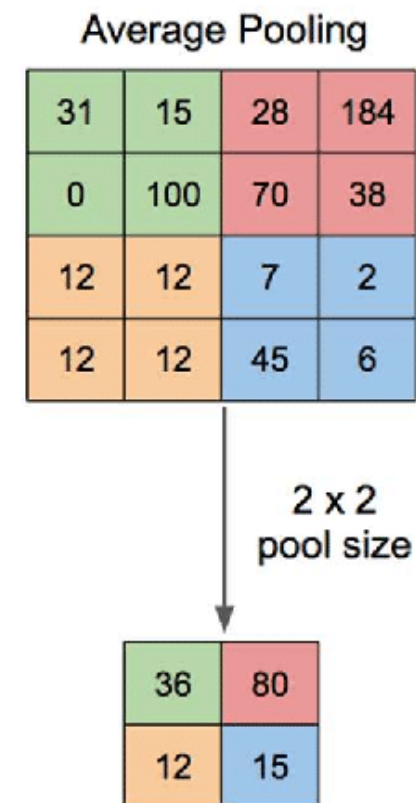
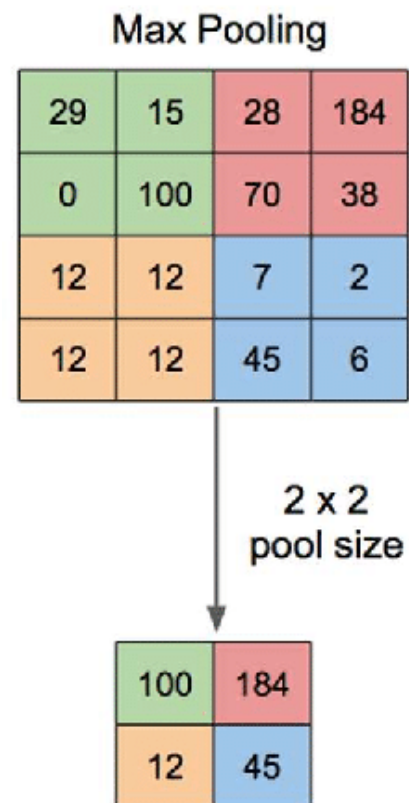
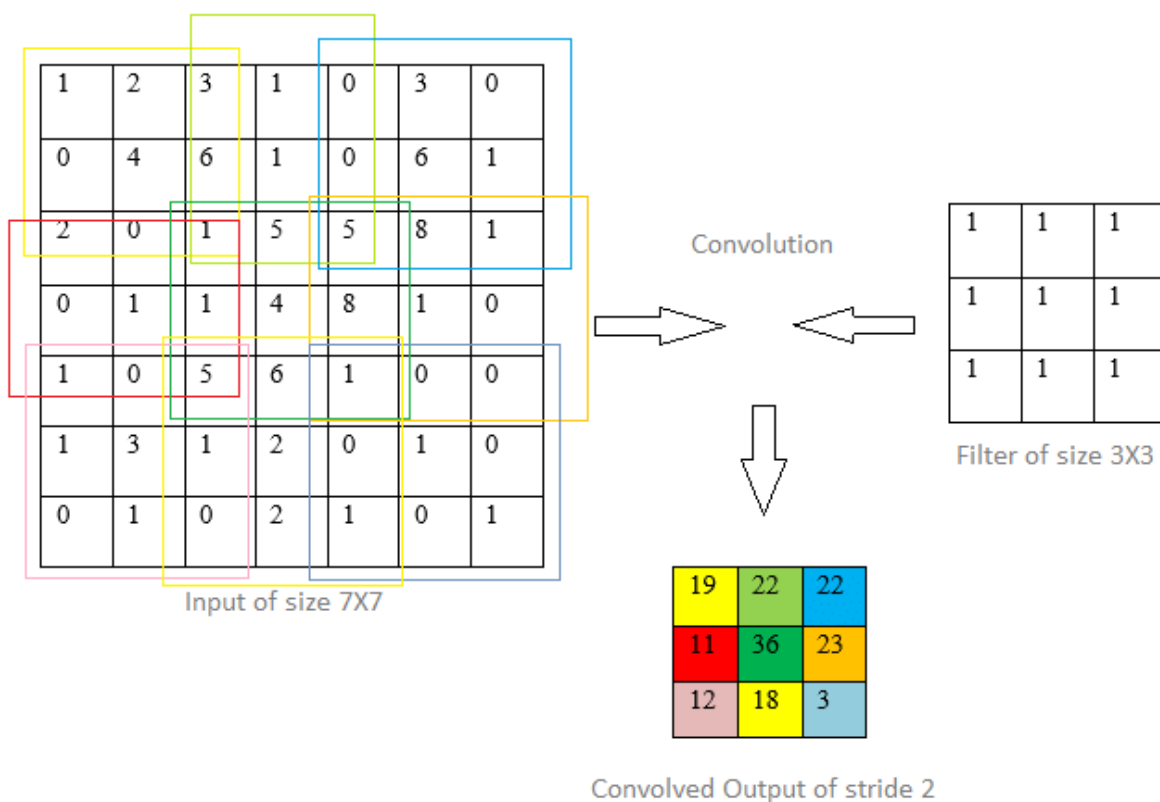
- ➔ Less computation
- ➔ More speed



Stride 2?

Moves 2 pixels at a time, so the output is half the size of the input

- ➔ Downsampling
- ➔ Less computation
- ➔ Keep important information than Max Pooling or Average Pooling (general downsampling)



Efficient DiT Design : Key-Value Token Compression

Effect of Compression Ratio

Res.	Ratio	FID ↓	CLIP-Score ↑	Train Latency ↓
512	1	8.244	0.276	2.3
512	2	9.063	0.276	2.2 (-4%)
512	4	9.606	0.276	2.1 (-9%)
1024	1	5.685	0.277	27.5
1024	2	5.512	0.273	22.5 (-18%)
1024	4	5.644	0.276	20.0 (-27%)
1024	9	5.712	0.275	17.8 (-35%)

(c) Compression ratios on different resolutions.

Res.	Ratio	Train Latency ↓ (s/Iter@32BS)	Test Latency ↓ (s/Img)
2K	1	56	58
2K	4	37 (-34%)	38 (-34%)
4K	1	191	91
4K	4	125 (-35%)	60 (-34%)

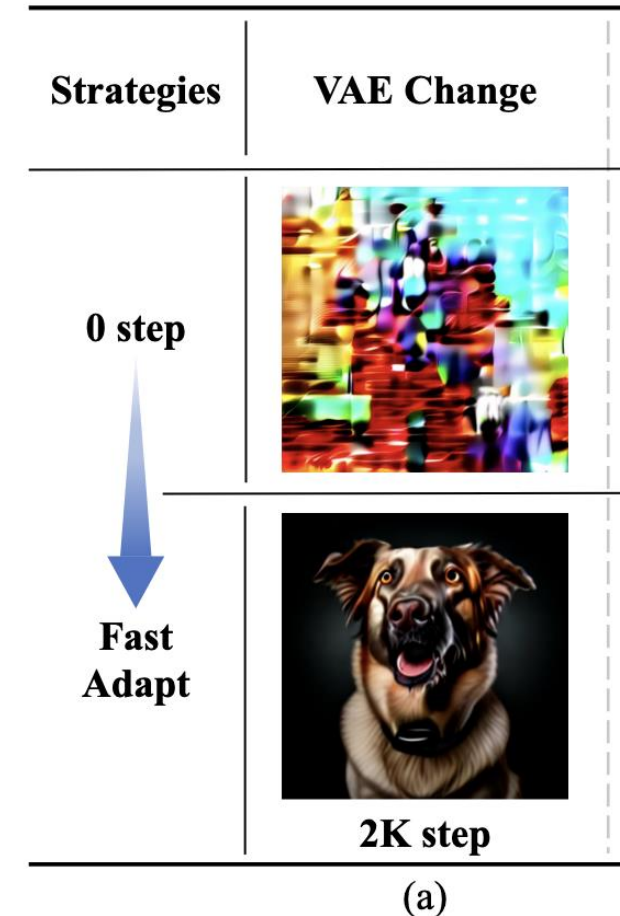
(d) Speed of different resolutions.

Table 3: KV-Token Compression Settings in Image Generation. This study employs FID, CMMD, and CLIP-Score metrics to assess the impact of various token compression components, such as compression ratio, positions, operators, and varying resolutions. Speed calculation in Tab. 3c is Second/Iteration/384 Batch-size.

Weak-to-Strong Training Strategy

Adapting model to new VAEs

- ◆ PixArt- α : VAE (8x downsampling)
 - ➔ PixArt- Σ : **Stable Diffusion XL(SDXL) VAE (4x downsampling)**
 - ➔ Preserve details
- ◆ If training T2I models from scratch = resource-intensive
 - ➔ choosing fine-tuning
- ◆ How?
 - ➔ fine-tuning quickly converges at **2K training steps**



Weak-to-Strong Training Strategy

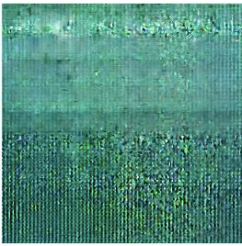

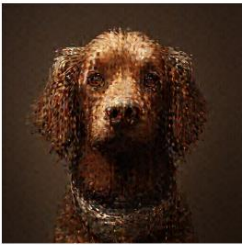

Adapting to Higher-Resolution

fine-tune from a low-resolution (LR) model to a high-resolution (HR) model
 → observe a performance degradation 😞

Using Positional Embedding Interpolation (PE Interpolation)

- Adapt quickly to new resolutions with fewer training steps (1000 steps)
- Create high-resolution images without learning from scratch

Resolution	Iterations	FID ↓	CLIP ↑
256	20K	16.56	0.270
256 → 512	1K	9.75	0.272
256 → 512	100K	8.91	0.276

Strategies	512 → 1024 RA	512 → 1024 RA + PE Interp.
0 step		
Fast Adapt		
	100 step	100 step

(b)

PE Interpolation?

Previous limitation

If learned location embedding is 512×512 in size (LR Model),
Directly applying this embedding to a higher resolution (1024×1024) will result in mismatch
→ poor performance 😞

Apply PE Interpolation

Interpolate the existing position embedding to the 1024×1024 size.
This means that 512 values are **naturally converted to smooth values**
in the process of scaling to 1024.

Weak-to-Strong Training Strategy

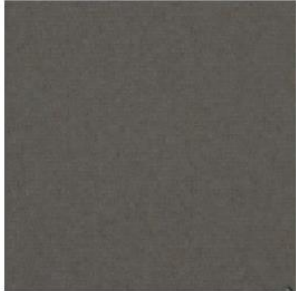
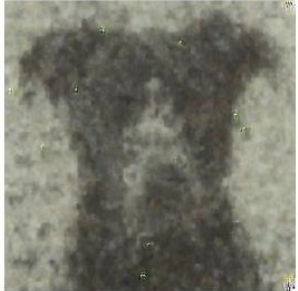
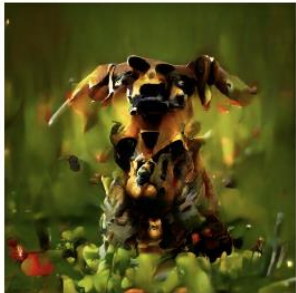

Adapting model to KV compression

Using KV Token Compression

- Risk of different structure 😞
- Difficult to use the trained weights of PixArt- α

Using Conv Avg Init

- Set the weighting value to $1/R^2$ to smooth the transition
- Preserving as much of the existing spatial information as possible

RA + PE Adjust + KV Compress	RA + PE Interp. + KV Compress + Conv Avg Init
	
	
100 step	100 step

(c)

Experiment : Implementation Details

Model

Text-Incoder	Flan-T5-XXL (= PixArt- α)
VAE	Stable Diffusion XL(SDXL)
Base model	PixArt- α

Hardware

Training GPU ($\leq 1K$ model)	32 NVIDIA Tesla V100
Training GPU (2K, 4K model)	16 NVIDIA A100
Optimization algorithms	CAME Optimizer

Evaluation Metrics

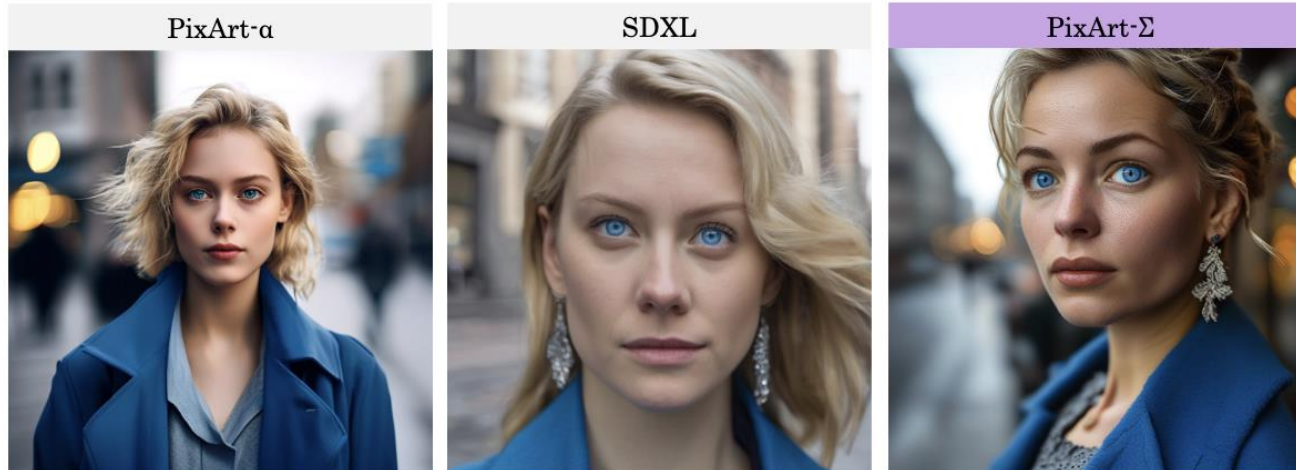
- 30,000 high quality dataset
- benchmark the most powerful T2I models.

Dataset	Volume	Caption	VN/DN	Total Noun	ACL	Average
Internal- α	14M	Raw	187K/931K	175M	25	11.7/Img
Internal- α	14M	LLaVA	28K/215K	536M	98	29.3/Img
Internal- α	14M	Share-Captioner	51K/420K	815M	184	54.4/Img
Internal- Σ	33M	Raw	294K/1512K	485M	35	14.4/Img
Internal- Σ	33M	Share-Captioner	77K/714K	1804M	180	53.6/Img
4K- Σ	2.3M	Share-Captioner	24K/96K	115M	163	49.5/Img

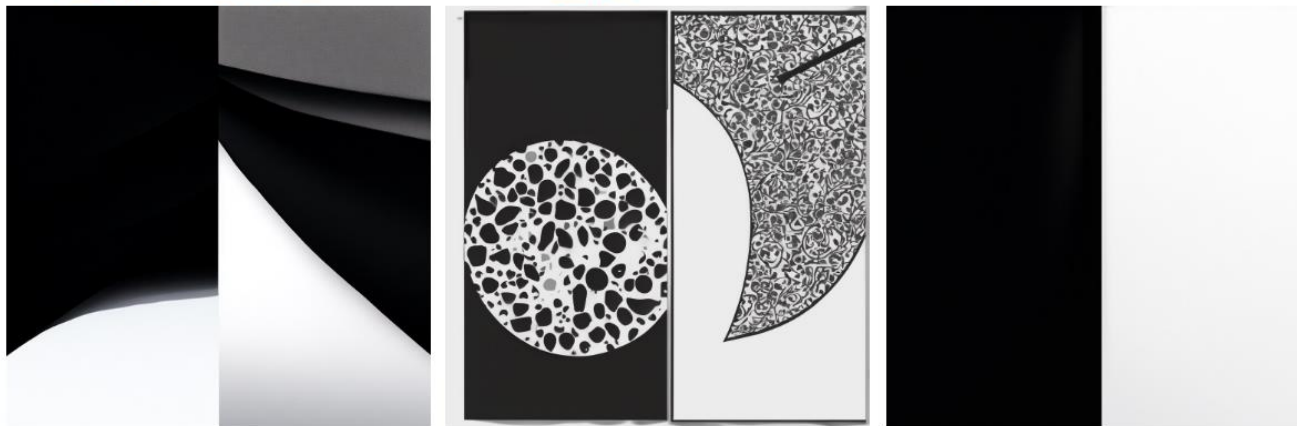
Experiment

Image Quality Assessment

Compared with open-source models



Prompt: A close-up photo of a person. The subject is a woman. She wore a **blue coat** with a **gray dress underneath**. She has **blue eyes** and **blond hair**, and wears a pair of **earrings**. Behind are blurred city buildings and streets.



Prompt: half a **solid black** background and half a **solid white** background



Experiment

Image Quality Assessment

competitive with these commercial products



Prompt: a small cactus with a happy face in the Sahara desert

Experiment

High-resolution Generation

A stylish woman walks down a Tokyo street filled with warm glowing neon and animated city signage. She wears a black leather jacket, a long red dress, and black boots, and carries a black purse. She wears sunglasses and red lipstick. She walks confidently and casually. The street is damp and reflective, creating a mirror effect of the colorful lights. Many pedestrians walk about.



Fig. 2: 4K image generation with complex dense instructions. PIXART- Σ can directly generate 4K resolution images without post-processing, and accurately respond to the given prompt.

Experiment

Human/AI (GPT4V) Preference Study

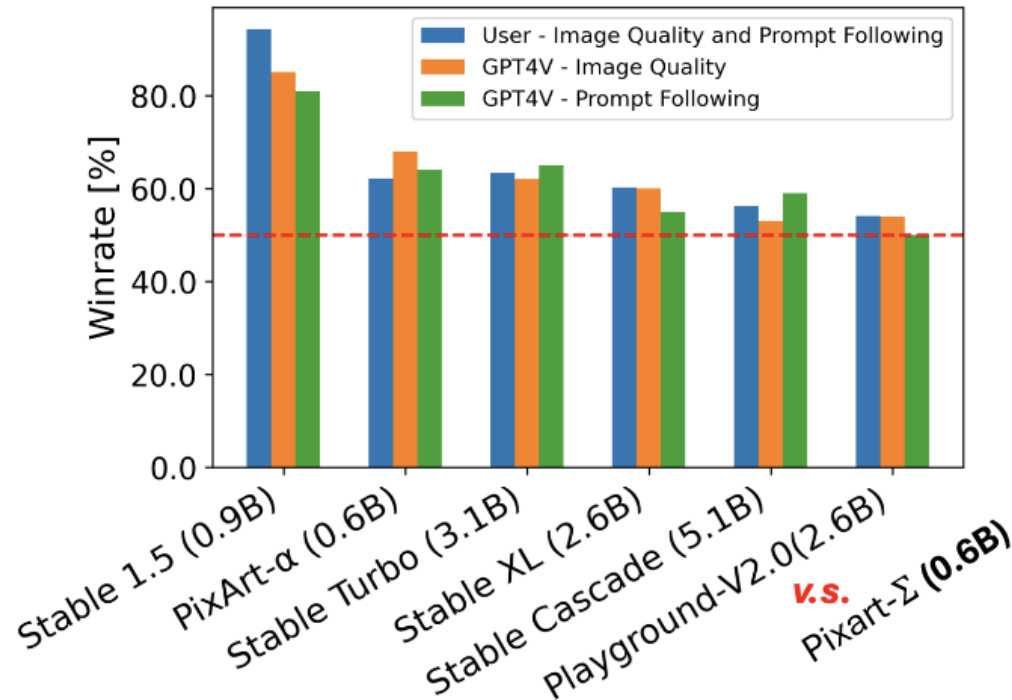


Fig. 9: Human(blue)/AI(orange and green) preference evaluation against current open T2I models. PIXART-Σ compares favorably against current state-of-the-art T2I models in both image quality and prompt-following.

(a) Designed prompts for Image-Quality Comparisons

Received with two images.

As an AI visual assistant, you are analyzing two specific images. Given a specific caption, you need to judge which image aligns with the caption more closely. Please pay attention to the key information, including object identities, properties, spatial relationships, object numbers and image style, etc.

The caption for the two images is: An old castle with a Gothic architectural style, soft focus photography, painted.

Please respond me strictly in the following format: <the first image is better> or <the second image is better> or <The two images are tied>. The reason is <give your reason here>.

<The first image is better>. The reason is that it more closely aligns with the description of "an old castle with a Gothic architectural style" due to its characteristic Gothic elements such as pointed arches, ribbed vaults, and flying buttresses. Furthermore, the image appears to have a soft focus effect that gives it a painted texture, matching the caption's requirements. The second image, while it does depict an old castle, it has a clearer focus and lacks the soft, painted quality described in the caption.

(b) Designed prompts for Prompt-Following Comparisons

Received with two images.

As an AI visual assistant, you are analyzing two specific images. When presented with a specific caption, it is required to evaluate and determine which image exhibits greater photorealism and aesthetical. Please pay attention to the key factors, including image style, the artistic quality of the image, realism, etc.

The caption for the two images is: An old castle with a Gothic architectural style, soft focus photography, painted.

Please respond me strictly in the following format: <the first image is better> or <the second image is better> or <The two images are tied>. The reason is <give your reason here>.

<The first image is better>. The reason is that the first image displays a higher level of photorealism with a more detailed portrayal of the castle and the surrounding environment. The lighting is more dynamic, with contrasts that suggest depth and texture, giving the castle a three-dimensional appearance. The mist and the overall composition add to the gothic aesthetic and the soft focus is effectively used to create an atmospheric quality, which aligns well with the concept of a painted old castle in the Gothic architectural style.

User
 GPT-4Vision

Conclusion

Key : Weak-to-Strong Training + KV Token Compression

Category	PixArt- α	PixArt- Σ
Maximum Resolution	1K (1024×1024)	4K (3840×2160) Supported
Computation Optimization	X	KV Token Compression (34% Reduction in Computation)
VAE Model	Basic VAE	SDXL VAE (Higher Quality Image Generation)
Text Token Length	120 Tokens	300 Tokens (More Precise Text-Image Alignment)
Training Strategy	Standard Training	Weak-to-Strong Training (Utilizing Pre-trained Model for Faster Training)
Training Cost	High	9% GPU Cost

Limitations

Not on-device

- ◆ Can't run on mobile and edge devices
→ Privacy concerns

Insufficient dataset

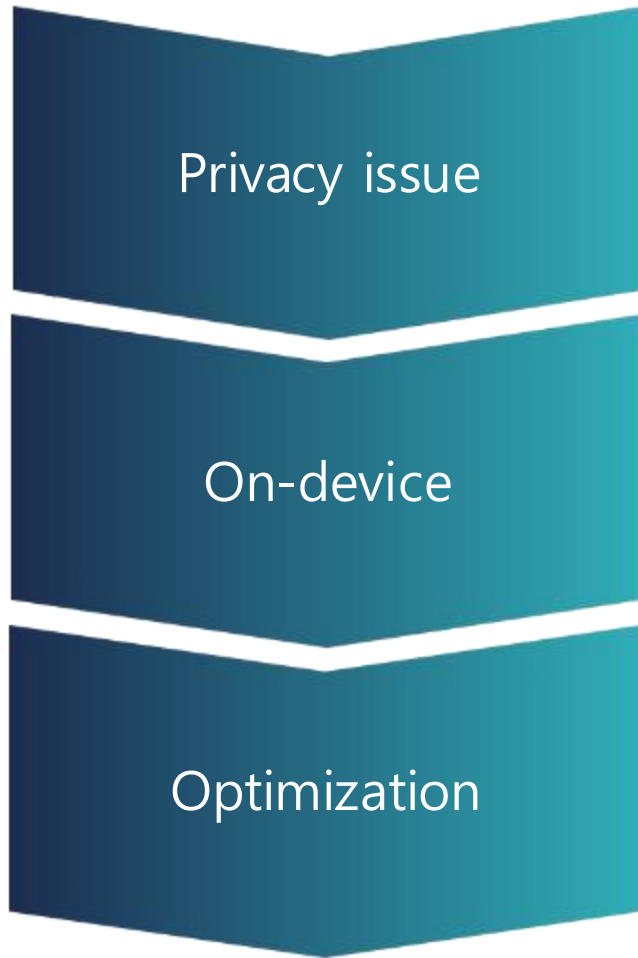
- ◆ Using 33M data = less than Stable Diffusion v1.5 (2B data)
→ Quality degradation

speed issues

- ◆ 4K creation is possible but **not optimized for speed**

Future work

4K On-device diffusion



- ◆ Handling photos is always a privacy risk
 - ◆ There is no on-device 4K diffusion paper now
 - ◆ Cloud can use your photo 😡
-
- ◆ Experimenting on-device with this model
 - ◆ Identify issues on-device (latency, battery, memory etc.)
-
- ◆ Optimized to work on **smartphones**
 - ◆ Optimize by applying modern paper techniques like '*MobileDiffusion*'

vs MobileDiffusion

Category	PixArt-Σ	MobileDiffusion
Model Architecture	Diffusion Transformer (DiT) based	Latent Diffusion + Optimized UNet
Text Encoder	Flan-T5-XXL	CLIP-ViT/L14
Image Resolution	Direct 4K (3840×2160)	512×512
KV Token Compression	✓	✗
Model Size	0.6B	386M
VAE (Autoencoder)	SDXL VAE	Lightweight VAE
Resolution Upscaling Method	PE Interpolation	Fixed at 512px (No upscaling)
Computation Optimization	Weak-to-Strong Training (Reuses pre-trained models)	Transformer block removal + Convolution-based optimization
On-Device Execution	✗ Requires high-performance GPU	✓ iPhone 15 Pro, Samsung S24 etc.
Training Dataset Size	33M (Includes 4K)	150M
Image Quality Evaluation (FID Score)	8.23	11.67 (1-step) / 8.65 (50-step DDIM)
Text-Image Alignment (CLIP Score)	0.2797	0.320 (1-step) / 0.325 (50-step DDIM)
Generation Speed	Slow on high-end GPU for 4K	0.2s on iPhone 15 Pro

PixArt- Σ : Weak-to-Strong Training of Diffusion Transformer for 4K Text-to-Image Generation

Junsong Chen, Chongjian Ge, Enze Xie, Yue Wu, Lewei Yao, Xiaozhe Ren, Zhongdao Wang, Ping Luo, Huchuan Lu, Zhenguo Li

Jeonghoon Park, happypjh2001@unist.ac.kr
Undergraduate Research Intern,
Ubiquitous Artificial Intelligence Lab,
Department of Computer Science and Engineering,
Ulsan National Institute of Science and Technology