

Low-Rank Adaptation / Stable Diffusion / LLaMa

24.10.31.
Mingyu Kim

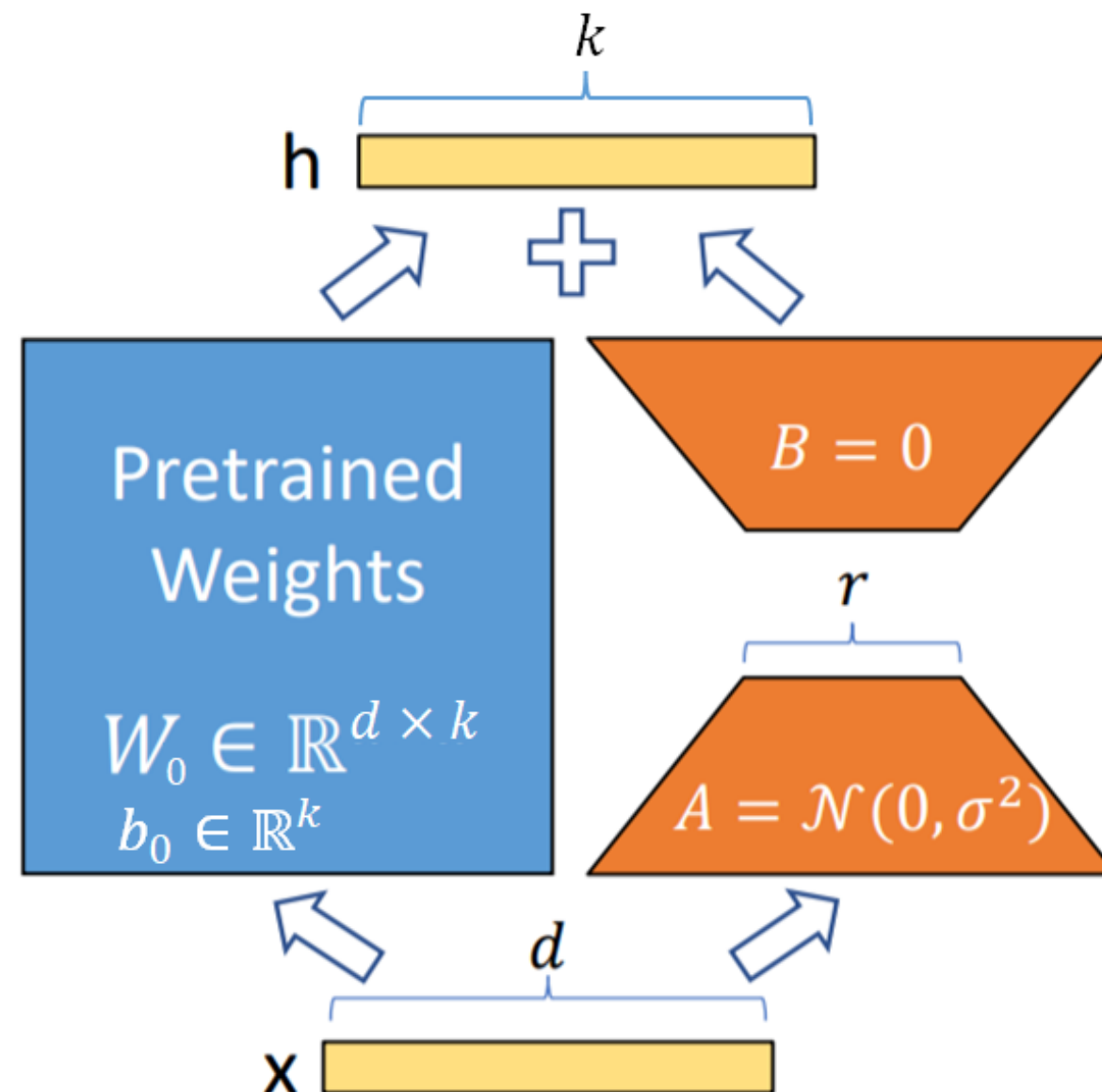
24.10.31

Repository

<https://github.com/MingyuKim87/day5.git>

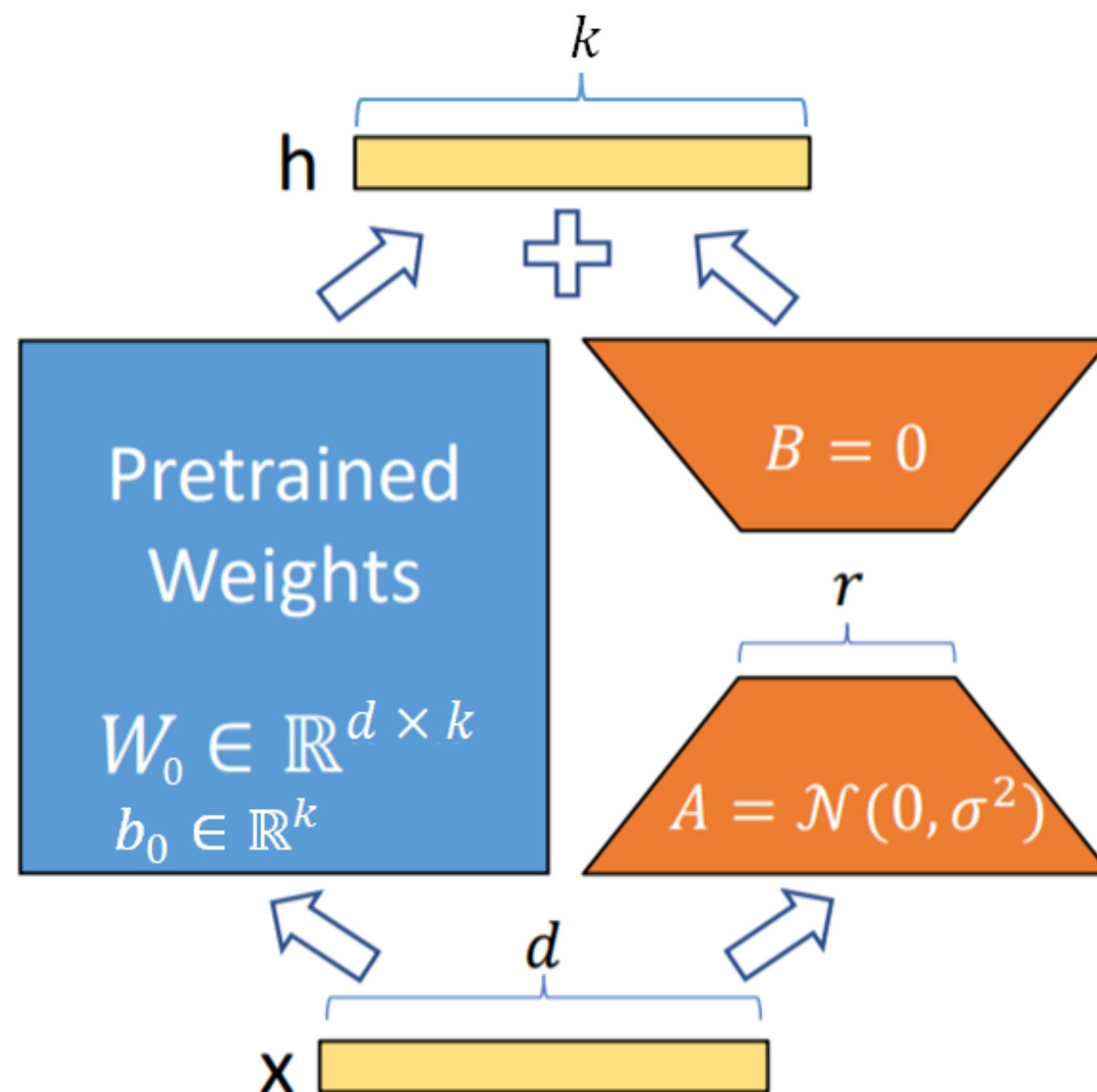
Low-Rank Adaptation

- The Pre-trained Network: $h = W_0x + b_0$ where $W_0 \in \mathbb{R}^{d \times d}$, $b_0 \in \mathbb{R}^{d \times 1}$
- The LoRA Network : A, B where $A \in \mathbb{R}^{d \times r}$, $B \in \mathbb{R}^{r \times d}$ and $d \ll \min(d, r)$
 - $A \sim N(0, \sigma^2)$, $B = 0$
- $h = (W_0x + b_0) + sBAx = W_0x + s\Delta Wx = (W_0 + s\Delta W)x + b_0$ where $s = \alpha/r$ (α : hyper-parameter)



Low-Rank Adaptation

- For instance, $d = 768, k = 256 \rightarrow$ the number of parameters $W_0 : 196,608$
 - $A : 767 \times 4 = 3,068$
 - $B : 4 \times 256 = 1,024$
- $A + B = 4,092$



Stable Diffusion VI

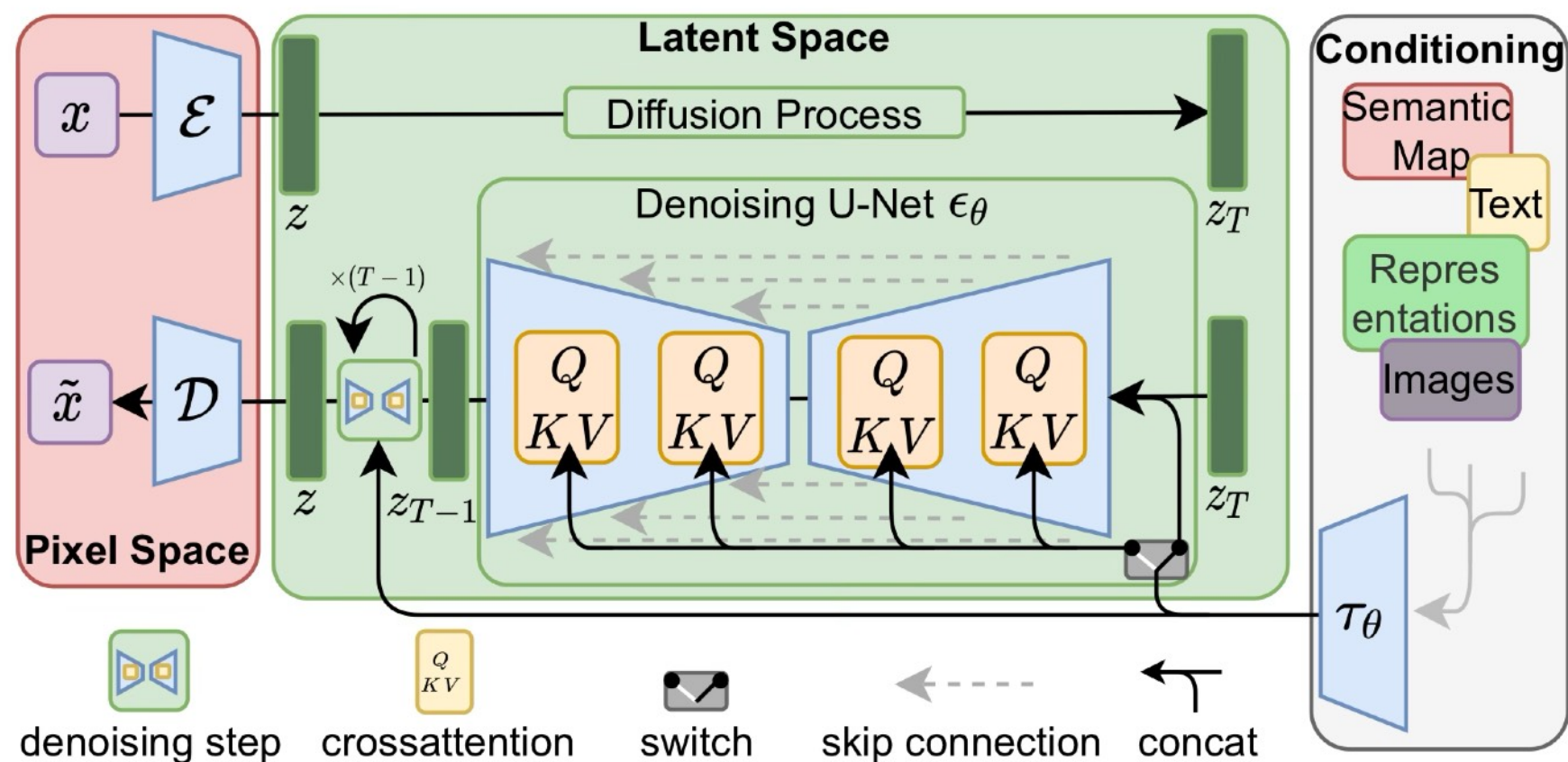
- Limitations of Previous Diffusion Models
 - Earlier diffusion models used a pixel-space UNet framework.
 - Generating high-resolution images with these models faces significant computational challenges due to the large amount of pixel data.

Model (<i>reg.-type</i>)	train throughput samples/sec.	sampling throughput [†]		train+val hours/epoch	FID@2k epoch 6
		@256	@512		
<i>LDM-1</i> (no first stage)	0.11	0.26	0.07	20.66	24.74
<i>LDM-4</i> (<i>KL</i> , w/ attn)	0.32	0.97	0.34	7.66	15.21
<i>LDM-4</i> (<i>VQ</i> , w/ attn)	0.33	0.97	0.34	7.04	14.99
<i>LDM-4</i> (<i>VQ</i> , w/o attn)	0.35	0.99	0.36	6.66	15.95

Table 6. Assessing inpainting efficiency. [†]: Deviations from Fig. 7 due to varying GPU settings/batch sizes *cf.* the supplement.

Stable Diffusion VI

- Advancements with Stable Diffusion
 - Stable Diffusion introduces diffusion models that operate within *latent spaces*, rather than direct pixel-space.
 - This approach is made possible by employing a Variational Autoencoder (VAE) architecture, which enables efficient high-resolution image synthesis by encoding images into compact latent representations.



'A street sign that reads
"Latent Diffusion" '

'A zombie in the
style of Picasso'



Stable Diffusion VI

- Diversity in Conditional Information Sources
 - Text Information: Allows for semantic input, guiding image generation with descriptive text.
 - Sketch Information: Adds structure or outlines to provide foundational visual guidance.
 - Scene Generation: Combines various conditional inputs to produce cohesive and context-rich scenes.

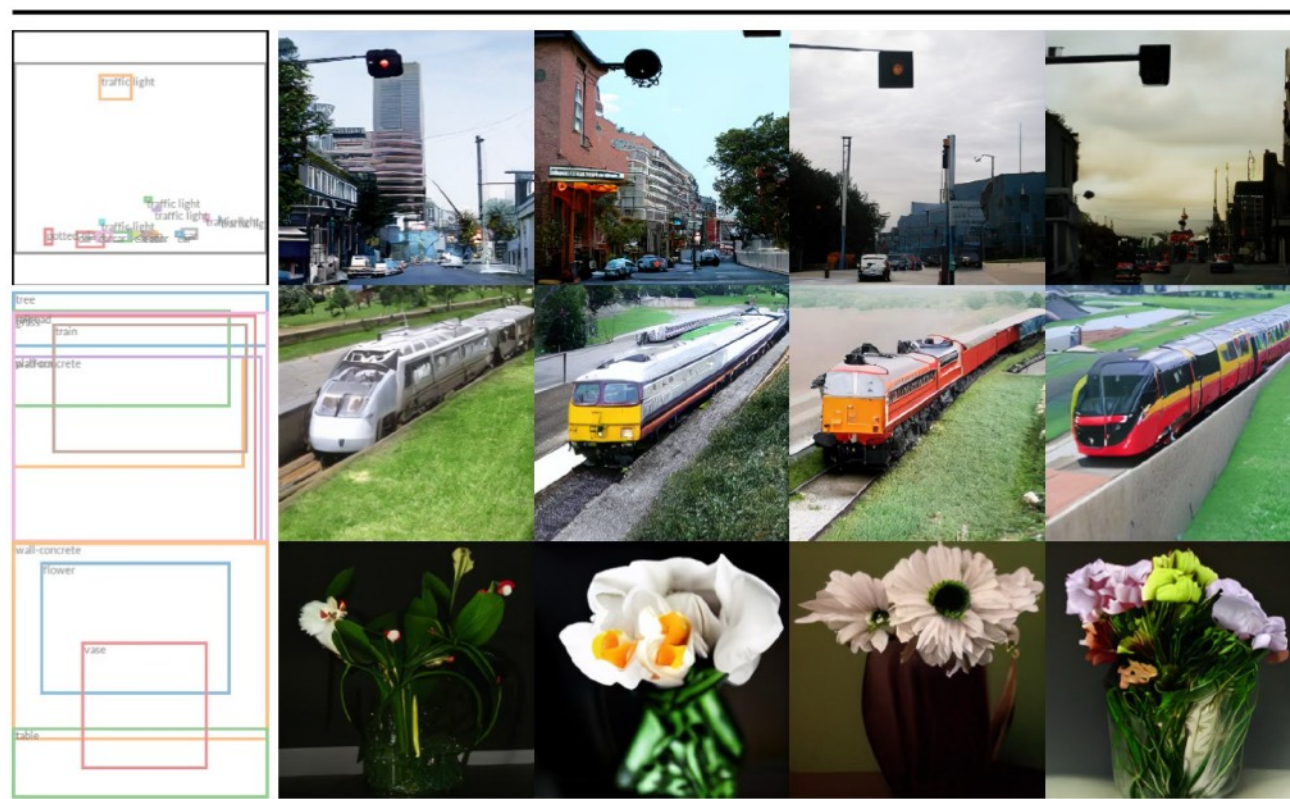


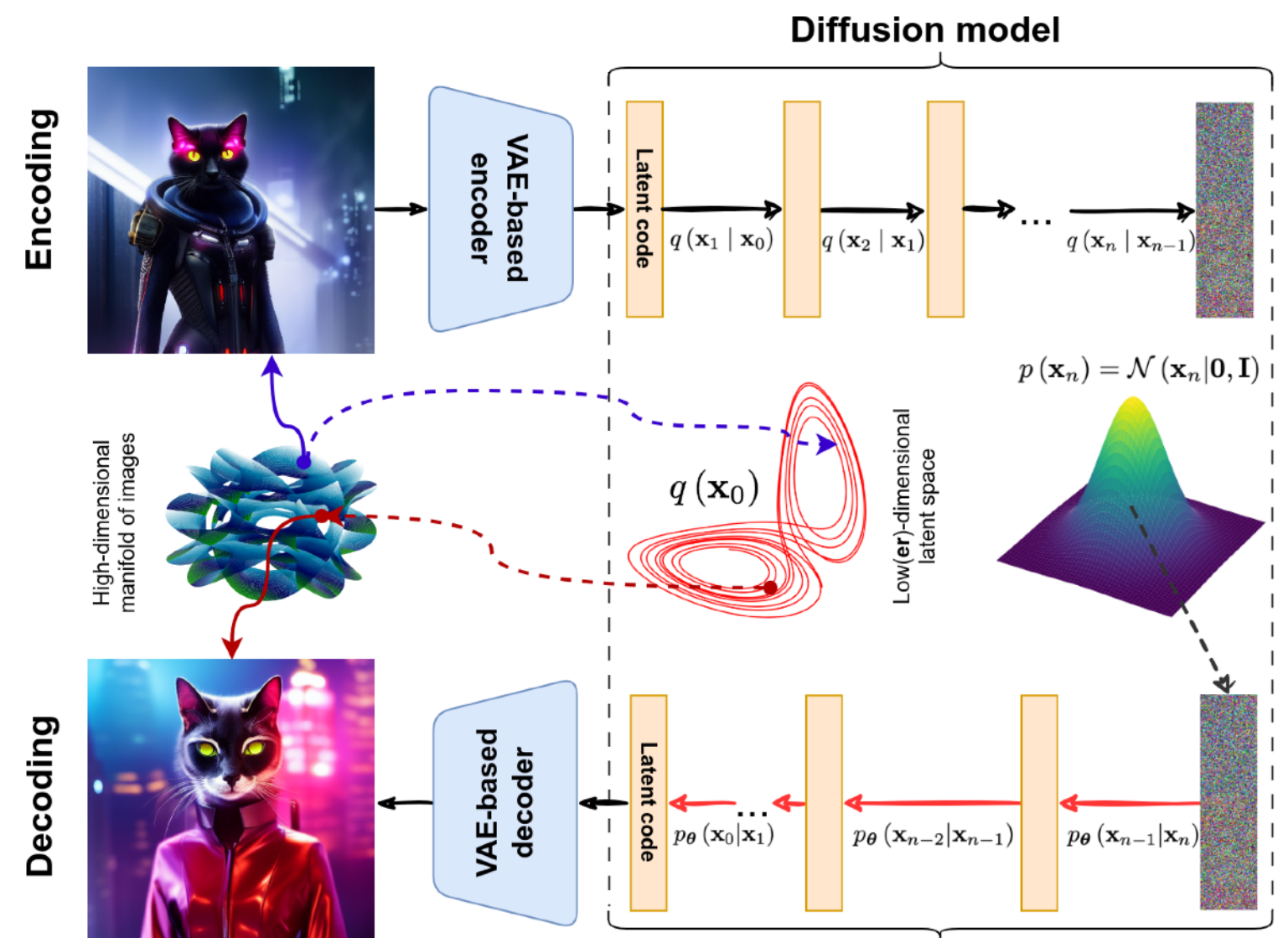
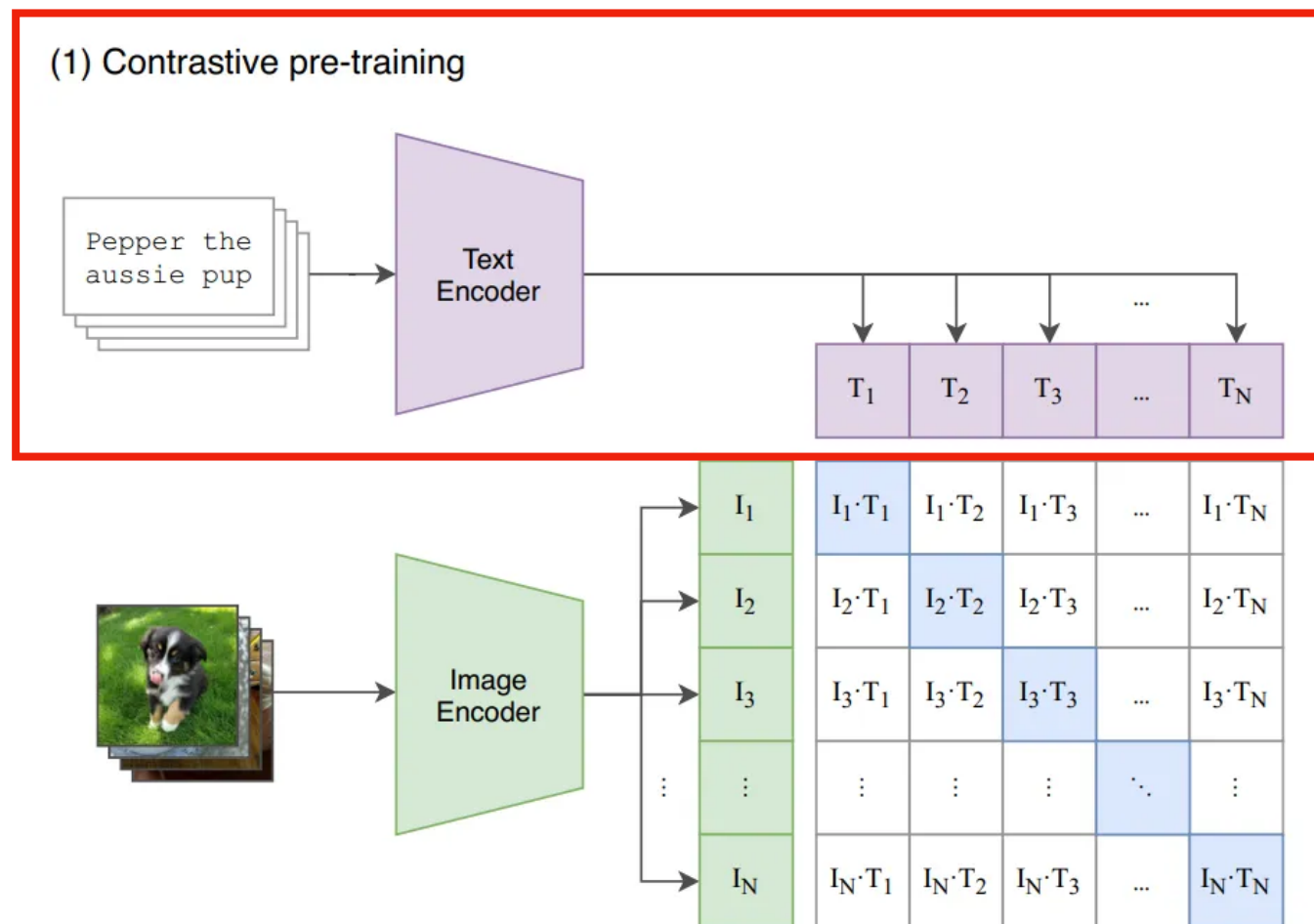
Figure 8. Layout-to-image synthesis with an *LDM* on COCO [4], see Sec. 4.3.1. Quantitative evaluation in the supplement D.3.



Figure 9. A *LDM* trained on 256^2 resolution can generalize to larger resolution (here: 512×1024) for spatially conditioned tasks such as semantic synthesis of landscape images. See Sec. 4.3.2.

Stable Diffusion VI

- Text-to-Image Architecture
 - CLIP Text Encoder (+ Tokenizer)
 - VAE: Encoder / Decoder
 - UNET: Diffusion Models



24.10.31

Guidance

- Three Equivalent Interpretations

24.10.31

Guidance

- Tweedie's Formula

24.10.31

Guidance

- Score-Based Generative Models

Guidance

- Classifier Guidance

$$\nabla \log p(\mathbf{x}_t|y) = \nabla \log \left(\frac{p(\mathbf{x}_t)p(y|\mathbf{x}_t)}{p(y)} \right) \quad (163)$$

$$= \nabla \log p(\mathbf{x}_t) + \nabla \log p(y|\mathbf{x}_t) - \nabla \log p(y) \quad (164)$$

$$= \underbrace{\nabla \log p(\mathbf{x}_t)}_{\text{unconditional score}} + \underbrace{\nabla \log p(y|\mathbf{x}_t)}_{\text{adversarial gradient}} \quad (165)$$

- Implementation

$$\nabla \log p(\mathbf{x}_t|y) = \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(y|\mathbf{x}_t) \quad (166)$$

Guidance

- Classifier Free Guidance
 - We assume that $\log p(x_t|y)$ is accessible

$$\nabla \log p(\mathbf{x}_t|y) = \nabla \log \left(\frac{p(\mathbf{x}_t)p(y|\mathbf{x}_t)}{p(y)} \right) \quad (163)$$

$$= \nabla \log p(\mathbf{x}_t) + \nabla \log p(y|\mathbf{x}_t) - \nabla \log p(y) \quad (164)$$

$$= \underbrace{\nabla \log p(\mathbf{x}_t)}_{\text{unconditional score}} + \underbrace{\nabla \log p(y|\mathbf{x}_t)}_{\text{adversarial gradient}} \quad (165)$$

- Adversarial gradient breaks down:

$$\nabla \log p(y|\mathbf{x}_t) = \nabla \log p(\mathbf{x}_t|y) - \nabla \log p(\mathbf{x}_t) \quad (167)$$

- Substituting this into:

$$\nabla \log p(\mathbf{x}_t|y) = \nabla \log p(\mathbf{x}_t) + \gamma (\nabla \log p(\mathbf{x}_t|y) - \nabla \log p(\mathbf{x}_t)) \quad (168)$$

$$= \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(\mathbf{x}_t|y) - \gamma \nabla \log p(\mathbf{x}_t) \quad (169)$$

$$= \underbrace{\gamma \nabla \log p(\mathbf{x}_t|y)}_{\text{conditional score}} + \underbrace{(1 - \gamma) \nabla \log p(\mathbf{x}_t)}_{\text{unconditional score}} \quad (170)$$

Stable Diffusion VI

- Classifier Free Guidance outperforms other methods and its simple generations.
 - FID (Fréchet Inception Distance): how close the generated images are to real images by comparing the distributions of features
 - IS (Inception Score): Represent Objects + Diverse Across Classes

Text-Conditional Image Synthesis				
Method	FID ↓	IS↑	N_{params}	
CogView [†] [17]	27.10	18.20	4B	self-ranking, rejection rate 0.017
LAFITE [†] [109]	26.94	<u>26.02</u>	75M	
GLIDE* [59]	<u>12.24</u>	-	6B	277 DDIM steps, c.f.g. [32] $s = 3$
Make-A-Scene* [26]	11.84	-	<u>4B</u>	c.f.g for AR models [98] $s = 5$
<i>LDM-KL-8</i>	23.31	20.03 ± 0.33	1.45B	250 DDIM steps
<i>LDM-KL-8-G*</i>	12.63	30.29 ± 0.42	<u>1.45B</u>	250 DDIM steps, c.f.g. [32] $s = 1.5$

24.10.31

Q&A

End