# Towards Understanding the Working Mechanism of Text-to-Image Diffusion Model

Mingyang Yi
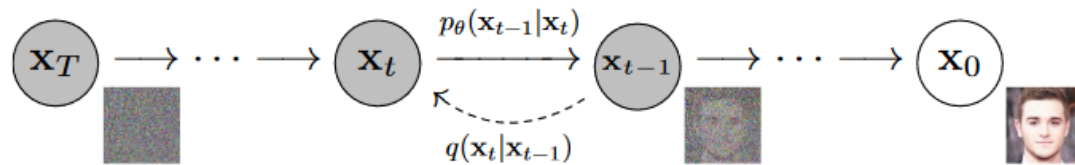
Huawei Noah's Ark Lab

Joint work with Aoxue Li, Yi Xin, Zhenguo Li
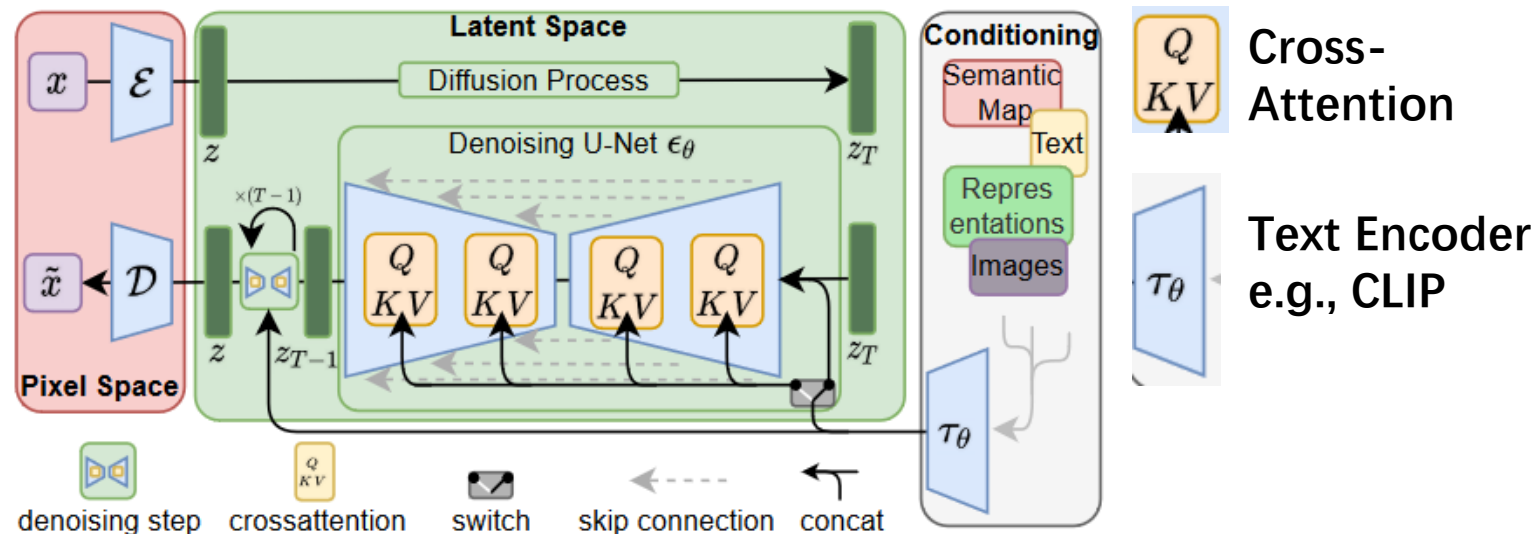
# Conditional Diffusion for T2I Generation

'A painting of a squirrel eating a burger'

- The Process of Diffusion Model

$$\mathbf{x}_T \longrightarrow \cdots \longrightarrow \mathbf{x}_t \xrightarrow{\ p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)\ } \mathbf{x}_{t-1} \longrightarrow \cdots \longrightarrow \mathbf{x}_0$$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

- Condition Diffusion Model for Text-to-Image Generation

**Cross-Attention** $Q$, $K$, $V$

**Text Encoder e.g., CLIP** $\tau_\theta$

How does T2I generation diffusion model works in practice?

Latent Space — Diffusion Process — Denoising U-Net $\epsilon_\theta$ — Conditioning: Semantic Map, Text, Representations, Images

Pixel Space — $x$ — $\mathcal{E}$ — $z$ — $\times(T-1)$ — $z_{T-1}$ — $z$ — $\tilde{x}$ — $\mathcal{D}$ — $z_T$ — $Q$ $KV$ — $\tau_\theta$

denoising step   crossattention   switch   skip connection   concat

# Quickly Appeared Shape

- Cross-Attention is Weighted Sum over Tokens

$$Q = W_Q \phi(x_t); K = W_K \mathcal{C}; V = W_V \mathcal{C}.$$

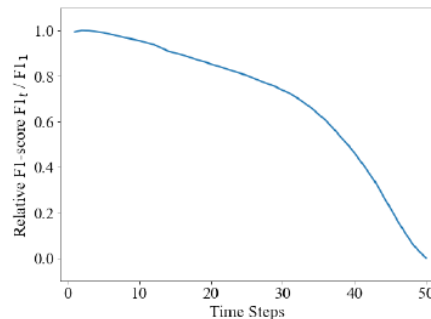$$\text{Attention}(Q, K, V) = \text{Softmax}(QK^\top/\sqrt{d})V$$

**Cross-Attention Map: Image-Token Correlation**

$\phi(x_t)$: **Pixel**    $\mathcal{C}$: **Textual Prompt (Embedding)**



Generate Image

Text Prompt: The square coaster was next to the circular mug.

- The Shape is Quickly Recovered

**The shape of image has been decided in the first few diffusion steps.**



(b) Convergence of Cross-Attention Map

# A Frequency Explain

**high-freq -> shape**　　　**low-freq -> details**

- Noisy data and its Frequency

frequency
$$F_{\boldsymbol{x}_t}(u,v) = \frac{1}{MN} \sum_{k=0}^{M-1} \sum_{l=0}^{N-1} x_t^{kl} \exp\left(-2\pi i \left(\frac{ku}{M} + \frac{lv}{N}\right)\right)$$
$$= \sqrt{\bar{\alpha}_t} F_{\boldsymbol{x}_0}(u,v) + \sqrt{1-\bar{\alpha}_t} F_{\boldsymbol{\epsilon}_t}(u,v),$$

noisy data $\quad \boldsymbol{x}_t = \sqrt{\bar{\alpha}_t}\boldsymbol{x}_0 + \sqrt{1-\bar{\alpha}_t}\boldsymbol{\epsilon}_t,$
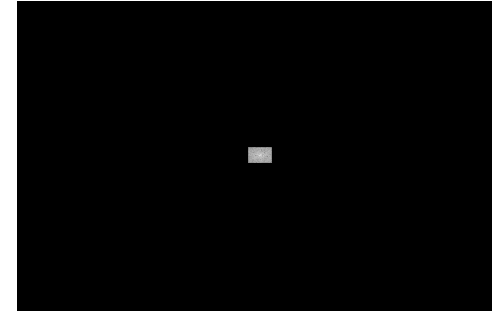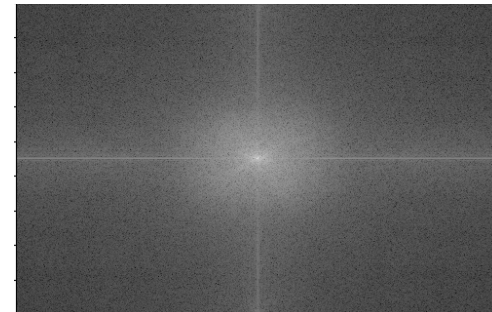
- Energy of High-Freq v.s. Low-Freq

**Proposition 1.** *For all* $u \in [M], v \in [N]$, *with high probability, we have*

$$\|F_{\boldsymbol{\epsilon}_t}(u,v)\|^2 \approx \mathcal{O}\left(\frac{1}{\sqrt{MN}}\right).$$

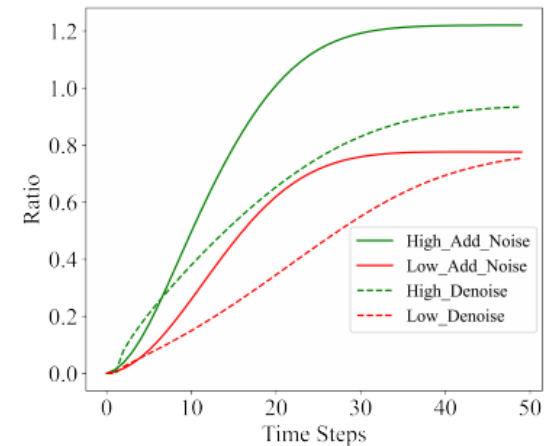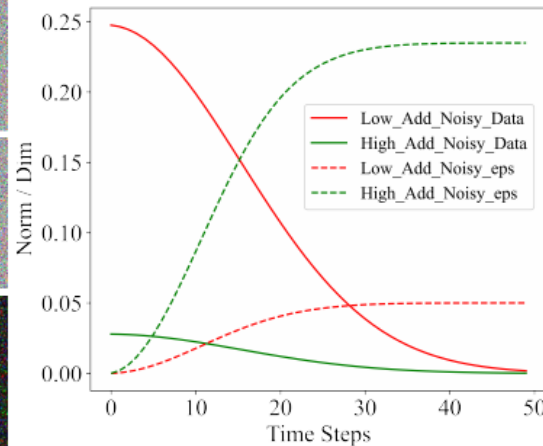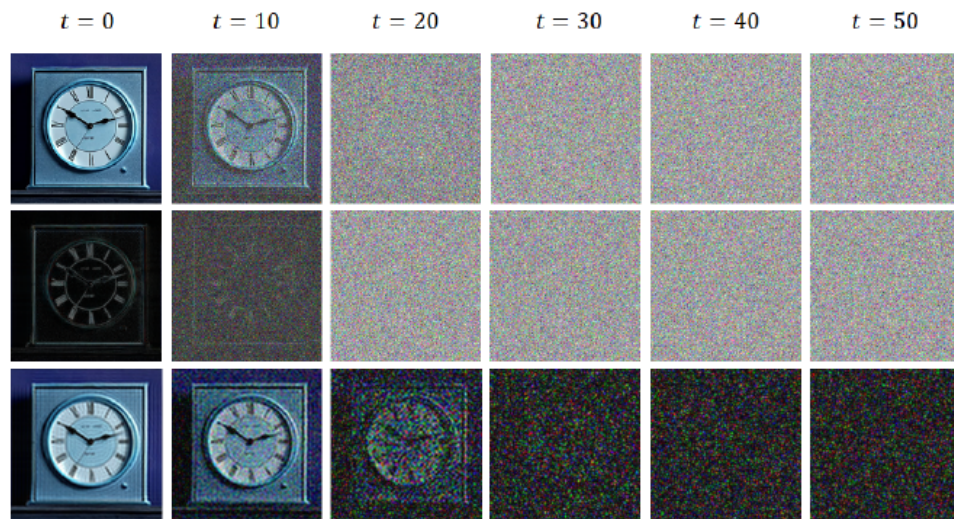white noise has more energy on high-freq.　　natural noise has more energy on low-freq.

80% spectrum are high-freq



Low-Freq part of Image

# Why First Shape then Details

- The high-freq part is quickly destroyed and will not be recovered until the end of reverse diffusion process. (vice-versa for low-freq)



(a) Noisy data and its high, low frequency parts (b) Norm of features $\sqrt{\bar{\alpha}_t}\boldsymbol{x}_0$ and $\sqrt{1-\bar{\alpha}_t}\boldsymbol{\epsilon}_t$ (c) Ratio of high / low frequency parts variation
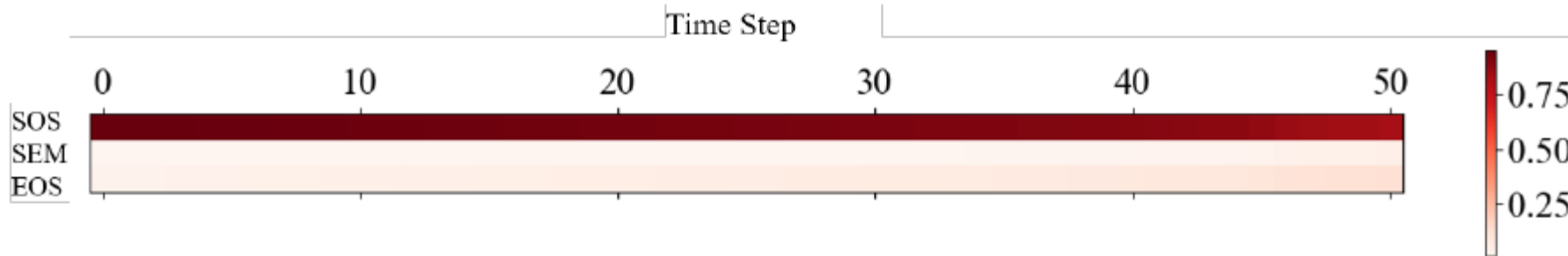
Focusing Shape and Details at beginning and end of diffusion, respectively.

# Text Prompts Related to the Phenomenon

- Three Classes of Tokens

    Prompt: [SOS] a white vase [EOS]  →  [SOS] + Sem + [EOS]

- Auto-regressive encoder makes [SOS] contains no information



Weights on tokens, [SOS] adjust
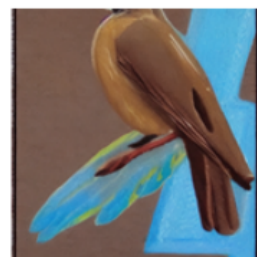weights on cross-attention map.

# [EOS] Decides Generation

- ## Generation Under Switched [EOS]

Prompt $A$ + [EOS]$_A$    Prompt $B$ + [EOS]$_B$    Prompt $A$ + [EOS]$_B$    Prompt $B$ + [EOS]$_A$

**Prompt A:** A blue bird
**Prompt B:** A brown chair



**Prompt A:** The sharp, angular edges of the city skyline pierced the clouds, a symbol of human innovation and progress.
**Prompt B:** The delicate, fluttering wings of the butterfly signaled the arrival of spring, a natural symbol of rebirth and renewal.



Observation I: [EOS] decides the overall T2I generation

Observation II: Slighter information in SEM is conveyed.

Table 1: The alignment of generated image with its source and target prompts. The prompts are constructed with switched [EOS].

| Alignment \ Prompt | Source | Target |
|---|---|---|
| Text-CLIPScore ↑ | 0.2363 | **0.2758** |
| BLIP-VQA↑ | 0.3325 | **0.4441** |
| MiniGPT-CoT↑ | 0.6473 | **0.7213** |

**Pay More Attention on [EOS]**

# When Does [EOS] Works

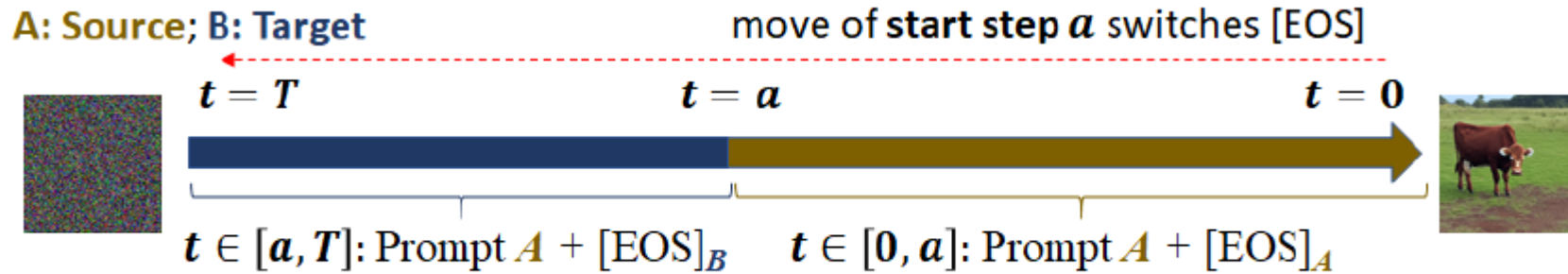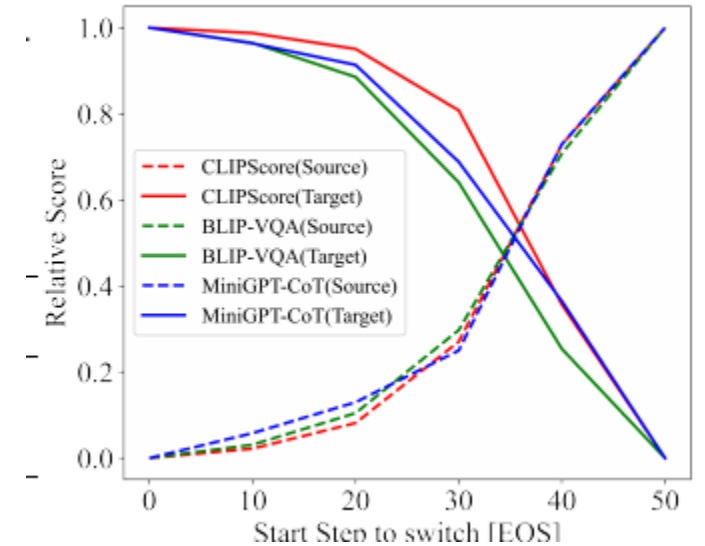- **The [EOS] works on the first shape reconstruction stage.**



A: Source; B: Target

move of **start step $a$** switches [EOS]

$t = T$         $t = a$         $t = 0$

$t \in [a, T]$: Prompt $A$ + $[EOS]_B$     $t \in [0, a]$: Prompt $A$ + $[EOS]_A$

Figure 5: Desnoising process under text prompt with switched [EOS] in $[a, 50]$.

CLIPScore(Source)
CLIPScore(Target)
BLIP-VQA(Source)
BLIP-VQA(Target)
MiniGPT-CoT(Source)
MiniGPT-CoT(Target)

Relative Score

Start Step to switch [EOS]

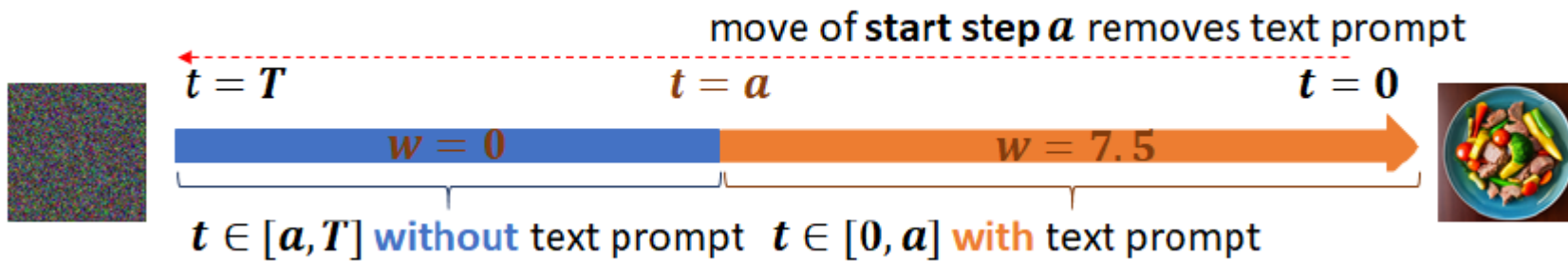The effect of [EOS] is not disappeared until the removing it at the beginning of denoising process.

# Text Information is Quickly Conveyed

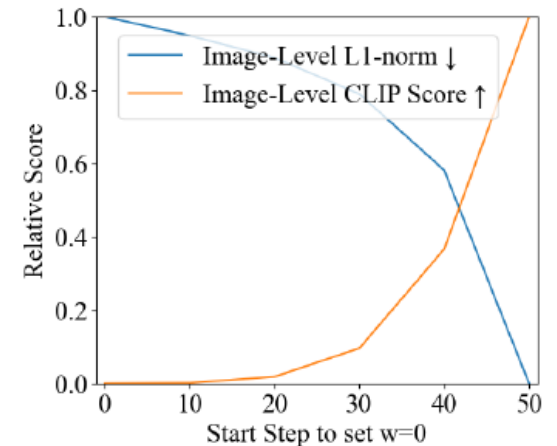- Noise Prediction

$$\epsilon_{\boldsymbol{\theta}}(t, \boldsymbol{x}_t, \mathcal{C}, \emptyset) = \epsilon_{\boldsymbol{\theta}}(t, \boldsymbol{x}_t, \emptyset) + w\left(\epsilon_{\boldsymbol{\theta}}(t, \boldsymbol{x}_t, \mathcal{C}) - \epsilon_{\boldsymbol{\theta}}(t, \boldsymbol{x}_t, \emptyset)\right),$$

- **Text Prompt Working on the Shape Reconstruction Stage**



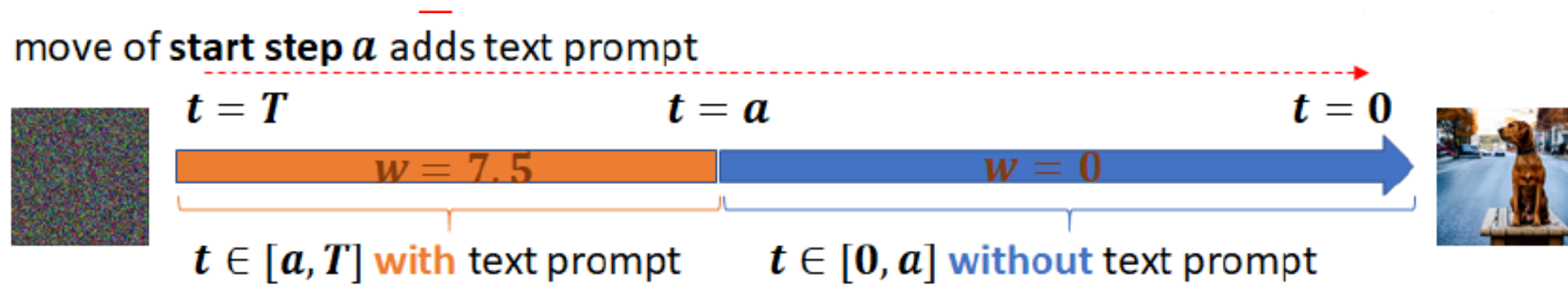Textual prompt is useless adding it at the beginning of diffusion process

# Conclusions

- The T2I Generation "First Overall Shape then Details".
- The [EOS] Has More Impact.
- The Mainly Text Prompt Works in the First Stage.

# Application

- Accelerating Sampling with Removing Text Information

$$\epsilon_{\theta}(t, x_t, C, \emptyset) = \begin{cases} \epsilon_{\theta}(t, x_t, \emptyset) + w\left(\epsilon_{\theta}(t, x_t, C) - \epsilon_{\theta}(t, x_t, \emptyset)\right) & a \leq t; \\ \epsilon_{\theta}(t, x_t, \emptyset) & 0 \leq t < a. \end{cases}$$

move of **start step** $a$ adds text prompt



$t = T$          $t = a$          $t = 0$

$w = 7.5$       $w = 0$

$t \in [a, T]$ **with** text prompt      $t \in [0, a]$ **without** text prompt

# Results

# Thanks!