



Controllable Text-to-Image Diffusion Models with Additional Conditions

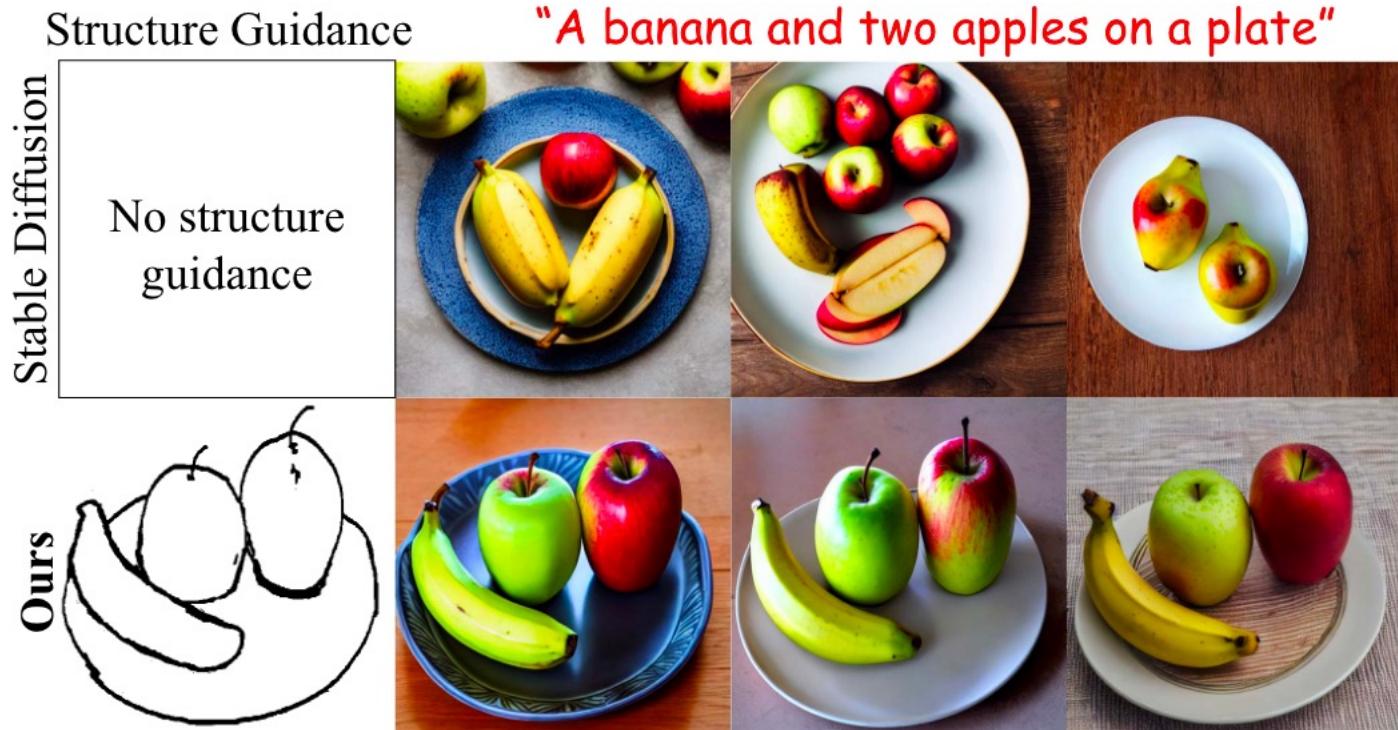
Qi Tianhao
2023/06/27



Task Introduction

2

- 扩散模型能够生成多样化的高质量图像，但是不能对生成图像提供精细化的控制





Methods

3

- Sketch-Guided Text-to-Image Diffusion Models
- Adding Conditional Control to Text-to-Image Diffusion Models
- T2I-Adapter: Learning Adapters to Dig out More Controllable Ability for Text-to-Image Diffusion Models



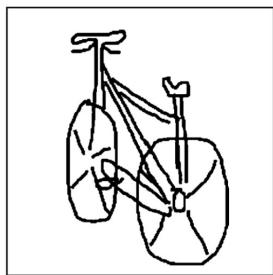
Sketch-Guided Text-to-Image Diffusion Models

4

□ Motivation

结合用户提供的**sketch**作为辅助信息，训练一个**per-pixel MLP**预测**edge map**，并通过梯度指导生成过程朝着接近给定**sketch**的方向进行

Input Sketch



“A photo of a bicycle”



“An origami bicycle”



“A bicycle in
a snowy weather”



“A macro photo of
a toy bicycle”



“A bicycle
made of wood”

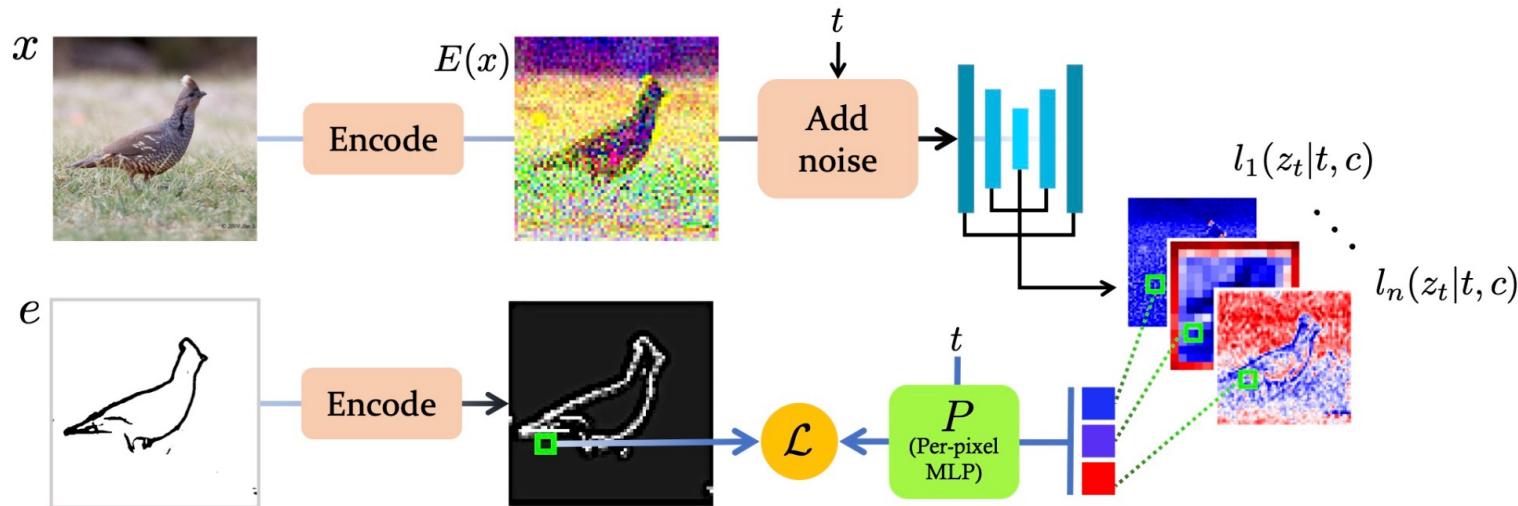


Sketch-Guided Text-to-Image Diffusion Models

5

□ Training

1. 抽取去噪过程中Unet的多个中间层特征，resize到输入尺寸后通道级联
2. 级联特征图通过per-pixel MLP预测edge map，与ground-truth计算损失



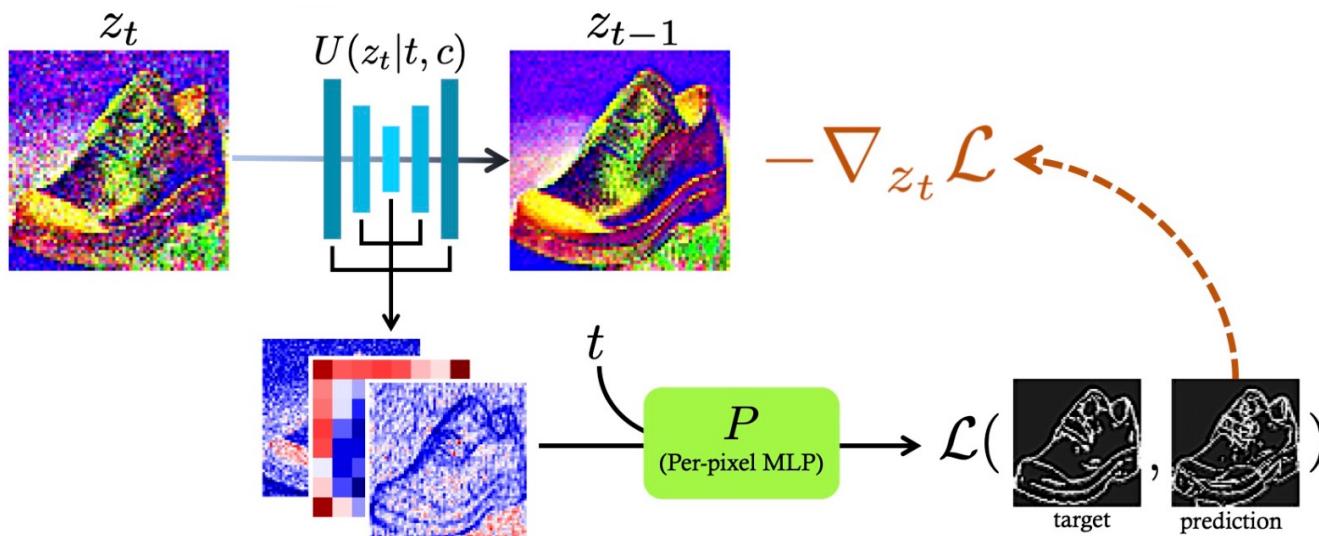
$$z_t = \alpha_t \cdot E(x) + \mu_t \cdot \xi \quad \mathbf{F}(w|c, t) = [l_1(w|c, t), \dots, l_n(w|c, t)]$$

$$\mathcal{L} = \mathbb{E}_{\substack{(x, e, c) \sim \mathcal{D} \\ \xi \sim \mathcal{N}(0, 1)}} \mathbb{E}_{t \sim \mathcal{U}([0, 1])} \sum_{i,j} \|P(\mathbf{F}(z_t|c, t)_{i,j}, t) - E(e)_{ij}\|^2$$

Sketch-Guided Text-to-Image Diffusion Models

□ Inference

1. 级联特征图通过 per-pixel MLP 预测 edge map，与 ground-truth 计算损失
2. 借鉴 classifier guidance 的思想，用损失梯度指导生成过程
3. 由于去噪过程的末端不会影响生成图像的几何布局，故仅在 $t=T, \dots, S > 1$ 使用 edge guidance，其中 $S=0.5T$



$$\mathcal{L}(\tilde{E}, E(e)) = \|\tilde{E} - E(e)\|^2$$

$$\tilde{z}_{t-1} = z_{t-1} - \alpha \cdot \nabla_{z_t} \mathcal{L}$$

$$\alpha = \frac{\|z_t - z_{t-1}\|_2}{\|\nabla_{z_t} \mathcal{L}\|_2} \cdot \beta$$



Sketch-Guided Text-to-Image Diffusion Models

7

□ Inference

2. 借鉴classifier guidance的思想，用损失梯度指导生成过程

Algorithm 2 Classifier guided DDIM sampling, given a diffusion model $\epsilon_\theta(x_t)$, classifier $p_\phi(y|x_t)$, and gradient scale s .

Input: class label y , gradient scale s

$x_T \leftarrow$ sample from $\mathcal{N}(0, \mathbf{I})$

for all t from T to 1 **do**

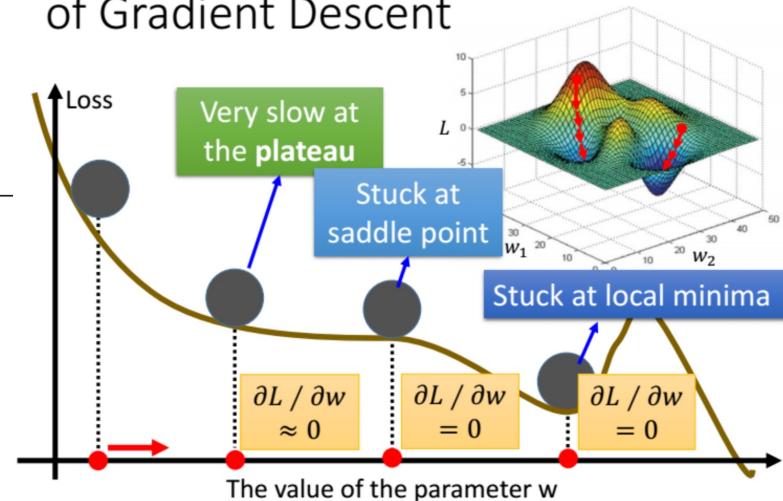
$$\hat{\epsilon} \leftarrow \epsilon_\theta(x_t) - \sqrt{1 - \bar{\alpha}_t} \nabla_{x_t} \log p_\phi(y|x_t)$$

$$x_{t-1} \leftarrow \sqrt{\bar{\alpha}_{t-1}} \left(\frac{x_t - \sqrt{1 - \bar{\alpha}_t} \hat{\epsilon}}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1 - \bar{\alpha}_{t-1}} \hat{\epsilon}$$

end for

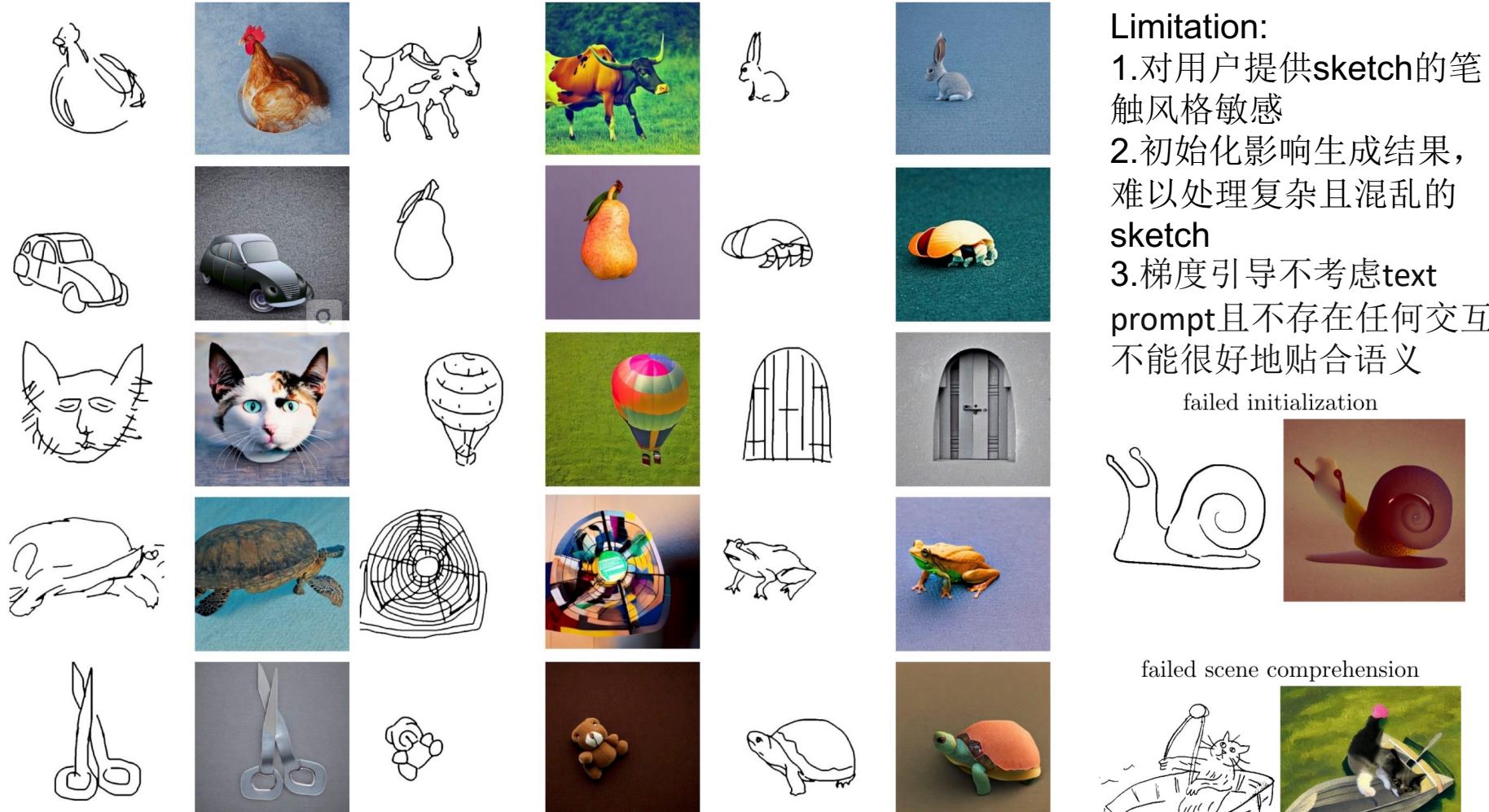
return x_0

More Limitation
of Gradient Descent



Sketch-Guided Text-to-Image Diffusion Models

8





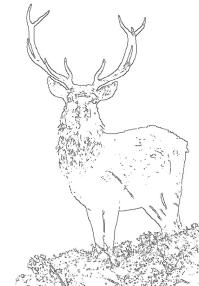
Adding Conditional Control to Text-to-Image Diffusion Models

Lvmin Zhang and Maneesh Agrawala
Stanford University

- 2021 苏州大学本科毕业
- 2021-2022 香港中文大学科研助理
- 2022至今 Stanford University
- [Style2Paints](#) 17.1k star
- [ControlNet](#) 27.3k star



Source image
(for canny edge detection)



Canny edge (input)

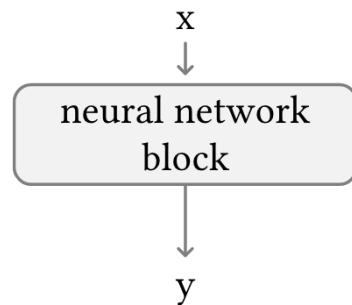
Generated images (output)

ControlNet

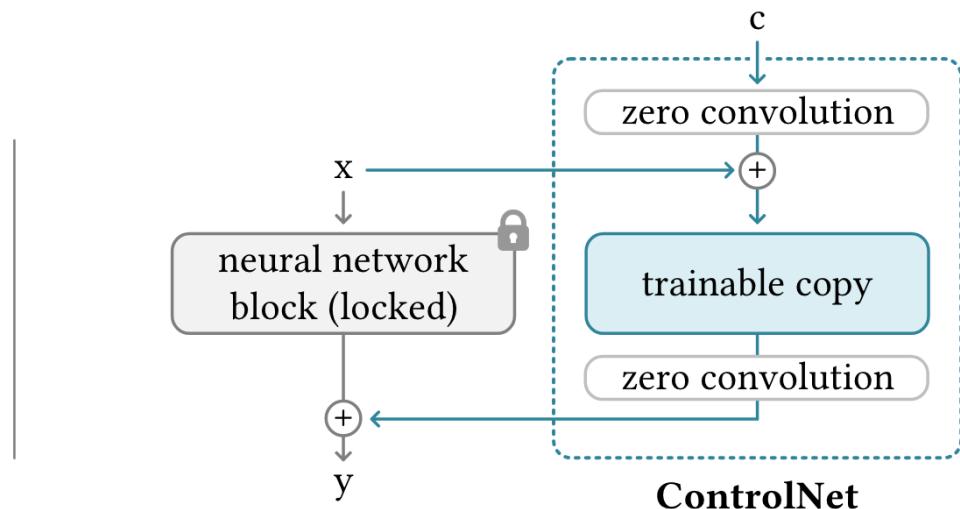
10

□ Motivation

1. 冻结原文生图模型参数，以保留在数十亿数据上训练得到的文生图能力
2. 引入额外分支，以实现额外条件的控制作用
3. zero convolution保证模型在优化前与原模型的容量、功能和结果保持一致



(a) Before

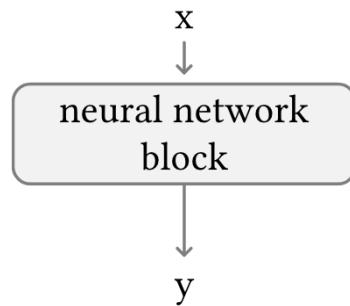


(b) After

ControlNet

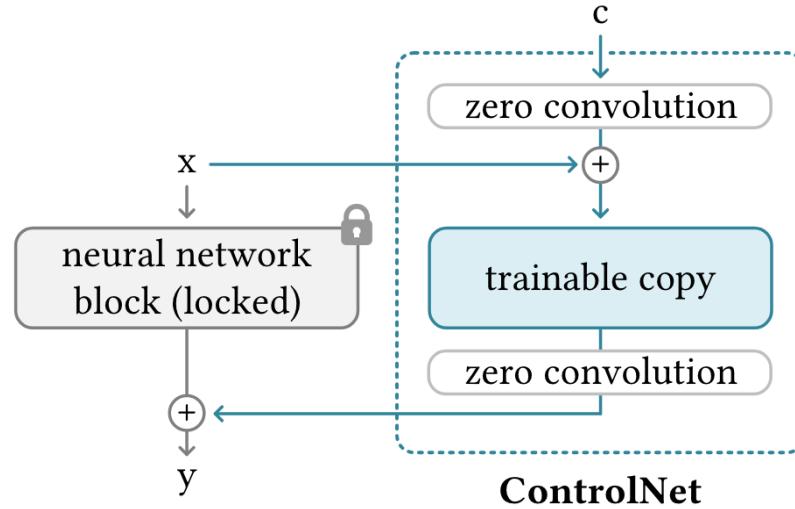


11



(a) Before

$$y = \mathcal{F}(x; \Theta)$$



(b) After

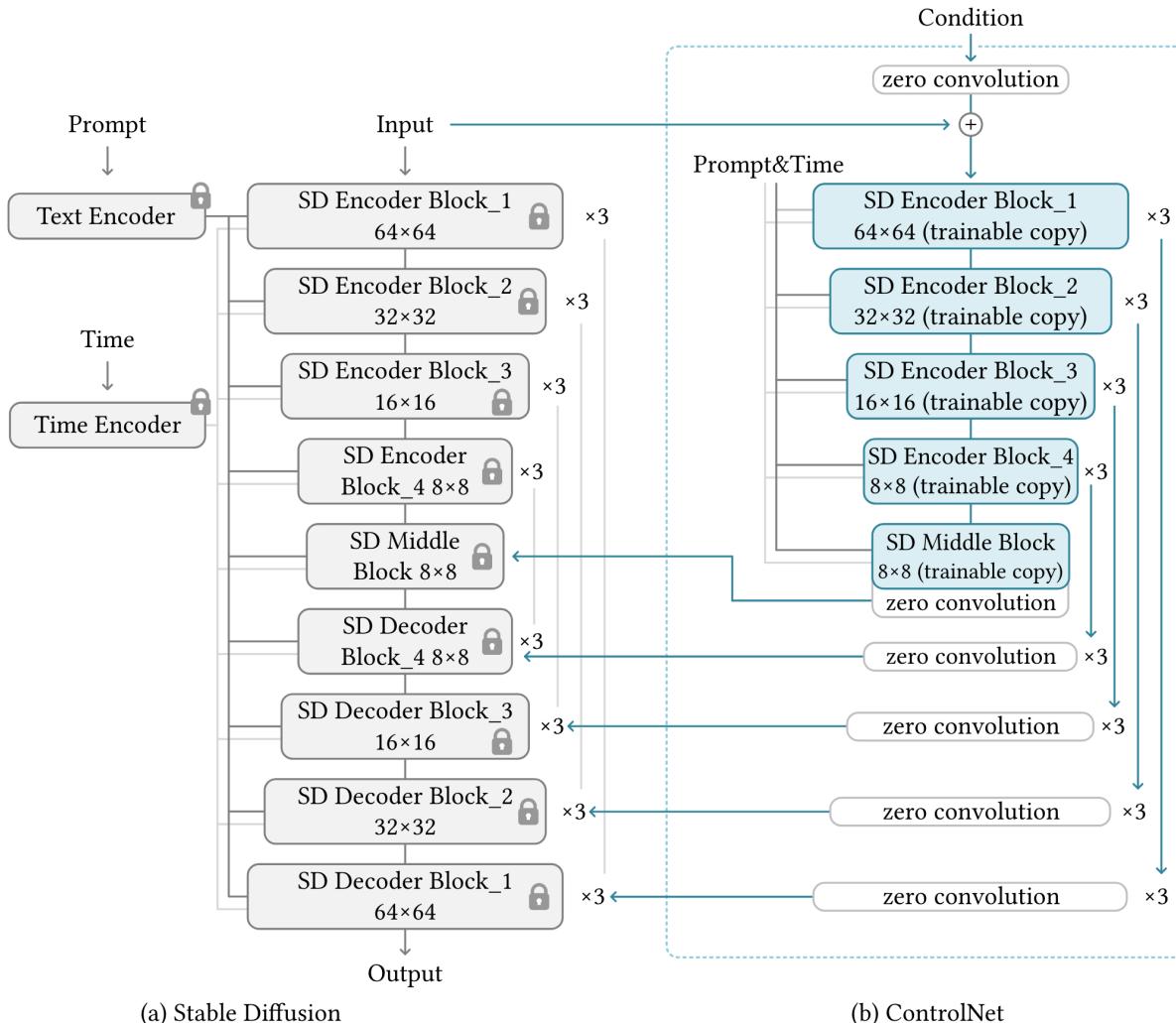
$$\mathbf{y}_c = \mathcal{F}(\mathbf{x}; \Theta) + \mathcal{Z}(\mathcal{F}(\mathbf{x} + \mathcal{Z}(\mathbf{c}; \Theta_{z1}); \Theta_c); \Theta_{z2})$$

$$\begin{cases} \mathcal{Z}(\mathbf{c}; \Theta_{z1}) = \mathbf{0} \\ \mathcal{F}(\mathbf{x} + \mathcal{Z}(\mathbf{c}; \Theta_{z1}); \Theta_c) = \mathcal{F}(\mathbf{x}; \Theta_c) = \mathcal{F}(\mathbf{x}; \Theta) \\ \mathcal{Z}(\mathcal{F}(\mathbf{x} + \mathcal{Z}(\mathbf{c}; \Theta_{z1}); \Theta_c); \Theta_{z2}) = \mathcal{Z}(\mathcal{F}(\mathbf{x}; \Theta_c); \Theta_{z2}) = \mathbf{0} \end{cases}$$

$$y_c = y$$

ControlNet

12



1. Condition 通过一个 tiny network 将图像空间的条件映射到隐空间 (8 conv layers in total, 3 conv layers with stride 2)

2. 随机 mask 50% 的 text prompt, 使得 controlnet 从额外条件中学到更丰富的语义信息

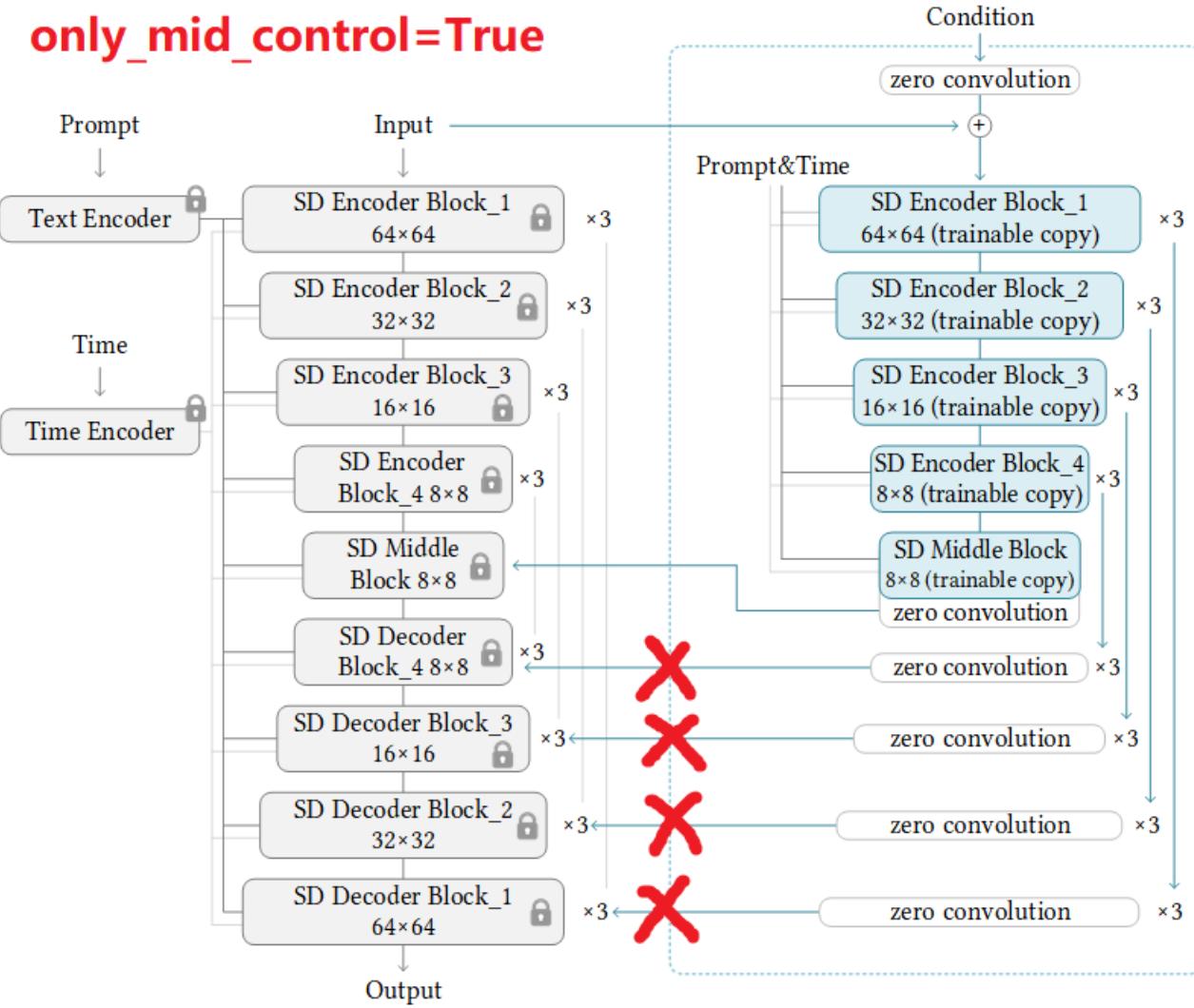
3. Training objective

$$\mathcal{L} = \mathbb{E}_{z_0, t, c_t, c_f, \epsilon \sim \mathcal{N}(0, 1)} [\|\epsilon - \epsilon_\theta(z_t, t, c_t, c_f)\|_2^2]$$

ControlNet with Limited Computing Resources

13

only_mid_control=True



1. 速度1.6x
2. 收敛后再设成False
finetune一些steps



ControlNet

14

□ More Examples



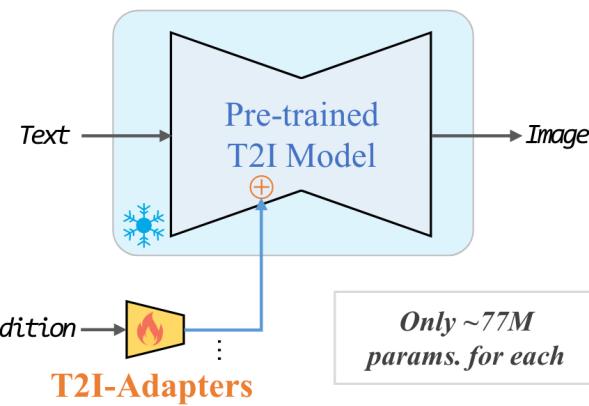
T2I-Adapter

15

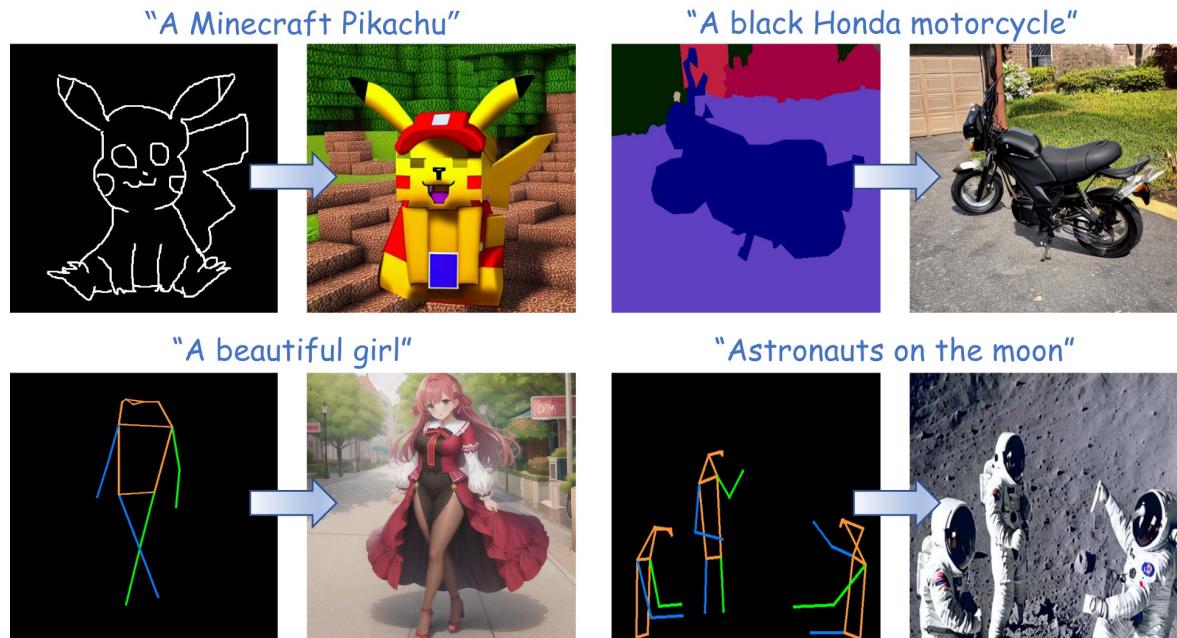
T2I-Adapter: Learning Adapters to Dig out More Controllable Ability for Text-to-Image Diffusion Models

Chong Mou^{*1,2} Xintao Wang^{†2} Liangbin Xie^{*2,3,4} Jian Zhang^{†1}
Zhongang Qi² Ying Shan² Xiaohu Qie²

¹Peking University Shenzhen Graduate School ²ARC Lab, Tencent PCG ³University of Macau ⁴Shenzhen Institute of Advanced Technology

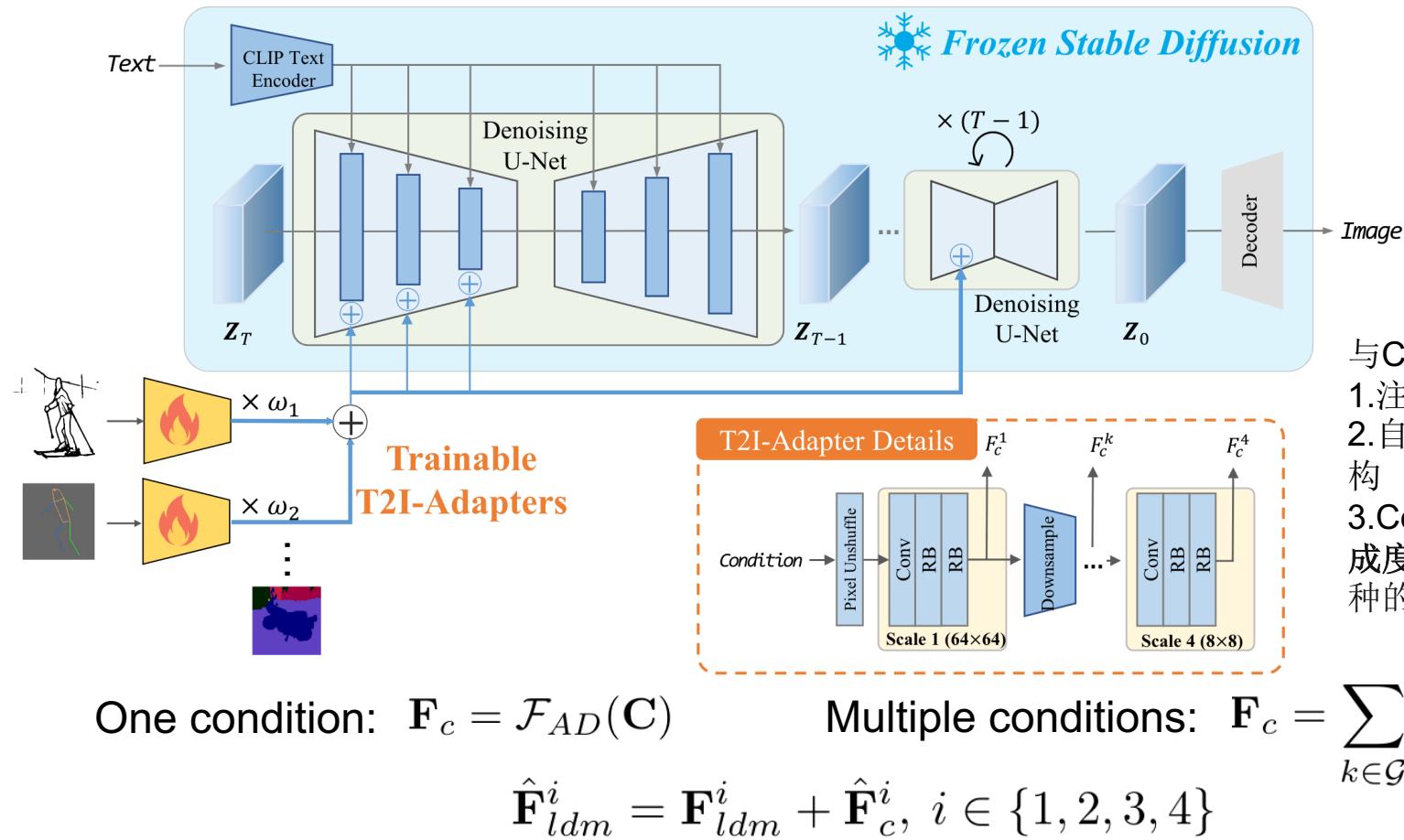


- ✓ **Plug-and-play.** Not affect original network topology and generation ability
- ✓ **Simple and small.** ~77M parameters and ~300M storage
- ✓ **Flexible.** Various adapters for different control conditions
- ✓ **Composable.** Several adapters to achieve multi-condition control
- ✓ **Generalizable.** Can be directly used on customized models



T2I-Adapter

16





T2I-Adapter

17

□ More Examples



Conclusion

18

- 效果为先，需要有令人眼前一亮的效果或**demo**展示
- **Adapter**架构仍是当前研究的热点方向，且其算力开销可以接受
- 直接注入参考图是待研究方向