

## ControlNet 论文带读

### 一. Introduction

1. 拟解决问题: 数据少, 算力大

### 二. Related Works

#### 1. Diffusion

DDPM  $\rightarrow$  DDIM  $\rightarrow$  Latent Diffusion Model.  
score-based  $\downarrow$  解决算力

#### 2. Fine-tuning

(1) Hyper Network: 在预训练模型后再加几层神经网络.

(2) zero convolution

#### 3. Text-to-Image Diffusion

(1) CLIP: 基于对比学习的文本和图像的多模态模型.  
(对图像打标签、对文本编码)

(2) Disco Diffusion

#### 4. Control of Pretrained Diffusion Model

(1) img2img (Stable Diffusion): color-level detail variations

(2) inpainting: 对图中某个区域进行修改.

#### 5. Img2Img

(1) Training Transformer.

(像素对齐关系)

### 三. Method.

#### 1. 设计哲学

遗忘 $\times$

很难做到强  
空间控制

以前控制扩散模型用微调/HyperNetwork, 所以采用保护与扩展的思想. (保护Stable Diffusion的权重).



## 2. ControlNet 网络搭建

(1) 令 SD 的 U-Net 里的一个神经网络层记作  $F$ , 参数为  $\theta$ .

① 输入是  $x$ , 输出是  $y$ .

$$y = F(x; \theta)$$

边缘图  $c$

② 一个分支是将  $\theta$  锁定  $\rightarrow \theta_{lock}$

另一分支复制相同结构  $\rightarrow \theta_{trainable}$ . 为了可以接收额外输入

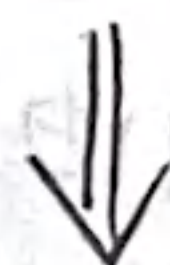
入, 该分支的输入为  $x + Z(c)$ .  $Z$  是简单的特征提取器.

③ Zero Convolution.

( $\Rightarrow$ ) 不能直接将 trainable copy 的输出加回去.

引入零卷积层  $\rightarrow Z(\cdot; \theta_{zero})$

★其中  $W, B$  初始化为 0.



$$y_c = F(x; \theta_{lock}) + Z(F(x + Z(c); \theta_{trainable}); \theta_{zero})$$

(2) Why Zero Convolution?

在训练的第一步, 零卷积的输出是零  $\Rightarrow y_c = F(x; \theta_{lock}) + 0 = y$

which means, 在训练开始时, ControlNet 对模型的影响为 0, 等价

于原本的 SD  $\Rightarrow$  保证了训练的稳定性.

$\uparrow$  逐步调整  $\theta_{zero}$  和  $\theta_{trainable}$ , 逐步注入控制信息.

(3) 复制了什么?

• 只复制了 Encoder 和 Middle Block (共 13 个模块), 并没有复制 Decoder.  $\leftarrow$  为了省参数

• 注入点: 加到 SD 的 decoder 的每一层上.

• ControlNet 的每一层输出, 会作为残差加到相应的 SD 层的输出上, 然后一起送入下一层的 SD-Decoder 中.



### 3. 训练策略.

(1) 数据构建: (原图  $x_0$ , Prompt, Condition Map ( $c$ ))  
↑  
自动生成

(2) Loss

$$L = E_{x_0, t, c, \varepsilon} [\|\varepsilon - \varepsilon_\theta(x_t, t, c)\|^2]$$

• SD 权重锁死, 梯度只会回传到 copy 层和零卷积层.

(3) 突然收敛/顿悟 (Grokking) 现象.

(4) 空文本训练 (Classifier-Free Guidance Support)

• 为了让 ControlNet 在推理时支持调节 Prompt 对生成结果产生影响, 训练时必须采用 Dropout 策略: 50% 将 Text Prompt 替换为 "" 字符串.

→ 迫使不仅仅依赖文本, 学会从 Condition Map 寻找线索.

(5) Zero convolution 中的 backpropagation.

① 前向:  $y = wx + b$   $\xrightarrow{\text{first step}} 0 \cdot x + 0 = 0 \Rightarrow$  传给 SD 的是 0.

② 反向: 需要更新权重  $w$ . (根据  $L$ )

Chain Law:  $\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial w} = \frac{\partial L}{\partial y} \cdot x$

$\Rightarrow$  只要输入不是 0, 并且模型有误差,  $w$  就会获得  $\nabla$ , 从而更新 ControlNet 的

(4) ControlNet 的头部有一个卷积网络, 用来做降维和匹配通道数与 U-Net 输入一致.



# 卷积与全连接的联系

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} & w_{15} \\ w_{21} & & & & \\ w_{31} & & & & \\ w_{41} & & & & \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix} \rightarrow \text{计算 } y_1 \text{ 需用到 } x_1 \sim x_5 \text{ 的所有输入信息.}$$

•  $[a, b]$  的卷积核:

$$\begin{aligned} y_1 &= ax_1 + bx_2 \\ y_2 &= ax_2 + bx_3 \\ &\vdots \\ y_4 &= ax_4 + bx_5 \end{aligned} \rightarrow \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} = \begin{bmatrix} a & b & 0 & 0 & 0 \\ 0 & a & b & 0 & 0 \\ 0 & 0 & a & b & 0 \\ 0 & 0 & 0 & a & b \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix}$$

⇒ 区别 { 稀疏连接: 局部感受野 ← Attention 改进  
参数共享: 平移不变性 (对特征位置不敏感)

Conclusion: 卷积是一种被限制了自由度, 但换来了极高的参数少

效率和泛化能力的全连接层.