

DALL·E 2

Hierarchical Text-Conditional Image Generation with CLIP Latents (2022) 452

GAN:

$z \rightarrow \text{Generator} \rightarrow x \searrow$
 $x' \rightarrow \text{discriminator} \rightarrow 0/1$

Pro: 真实性高, 需要的数据量不大

Con: 训练不太稳定, 训练两个网络有平衡的问题, 创造性不好

AE (auto-encoder):

$x \rightarrow \text{Encoder} \rightarrow z \text{ (bottleneck)} \rightarrow \text{Decoder} \rightarrow x' \quad \text{compare } x \text{ and } x'$

DAE (Denoising auto-encoder):

$x \rightarrow x_c \text{ (corrupted } x) \rightarrow \text{Encoder} \rightarrow z \rightarrow \text{Decoder} \rightarrow x' \quad \text{compare } x \text{ and } x'$

更加稳健, 不容易过拟合; 图像冗余性较高, 即使污染后还能抓住本质

VAE (Variational auto-encoder):

$x \rightarrow \text{Encoder} \rightarrow \mu, \sigma \text{ (distribution)} \rightarrow z \rightarrow \text{Decoder} \rightarrow x'$

$$q(z|x) \quad z = \mu + \sigma \cdot \epsilon \text{ prior} \quad p(x|z)$$

改为学习一个分布, 假定为高斯分布则可以用 μ, σ 表示, 从中sample得到 z

VQVAE (Vector Quantised VAE):

$x \rightarrow \text{Encoder} \rightarrow f \rightarrow \text{codebook } K \times D \rightarrow f_n \rightarrow \text{Decoder} \rightarrow x'$

codebook有 K 个长度为 D 的聚类中心, 当图片经过encoder得到特征图 f , 对比特征图与codebook中的向量, 找出最接近的聚类中心, 使用这个聚类中心对应的特征 f_n

DALLE

Text $\xrightarrow{\text{BPE}}$ $f_t \searrow$
concatenated vector \rightarrow GPT

Image $\xrightarrow{\text{VQVAE}}$ $f_q \nearrow$

Diffusion models

$x_0 \ x_1 \ \dots \ x_{t-1} \ x_t \ \dots \ x_T$

forward diffusion \rightarrow add noise

reverse diffusion \leftarrow U-Net, attention... time embedding in U-Net 提醒模型学到第几步了, 从粗略特征学到细微特征

DDPM \rightarrow improved DDPM \rightarrow Diffusion beats GAN \rightarrow GLIDE \rightarrow DALL-E2

DDPM贡献:

1. $x_t \rightarrow x_{t-1}$

$x_t \rightarrow \epsilon = x_{t-1} - x_t$ U-Net only predicts the residual (similar to ResNet)

2. 对于一个分布, 固定方差学习均值, 效果就很好了

Improved DDPM改动:

1. 学了方差, 效果更好

2. 添加噪声的schedule, 从线性改为余弦

3. 实验发现scale效果好

Diffusion beats GAN:

1. 更大更复杂的模型

2. 使用classifier guidance引导模型的采样和生成: 在训练模型的同时, 训练一个 classifier (可以在ImageNet上加了噪声训练), reverse diffusion $x_t \rightarrow x_{t-1}$ 过程中, 将 x_t 扔给classifier看分类得对不对, 算出交叉熵目标函数得到梯度 ∇g , 用梯度帮助模型的采样和生成

牺牲了一部分多样性换取了真实性, 在IS, FID上超过了GAN

一些其他guidance:

$p(x_{t-1}|x_t) = \|\epsilon - f_\theta(x_t, t, y)\|$ — y 是引导

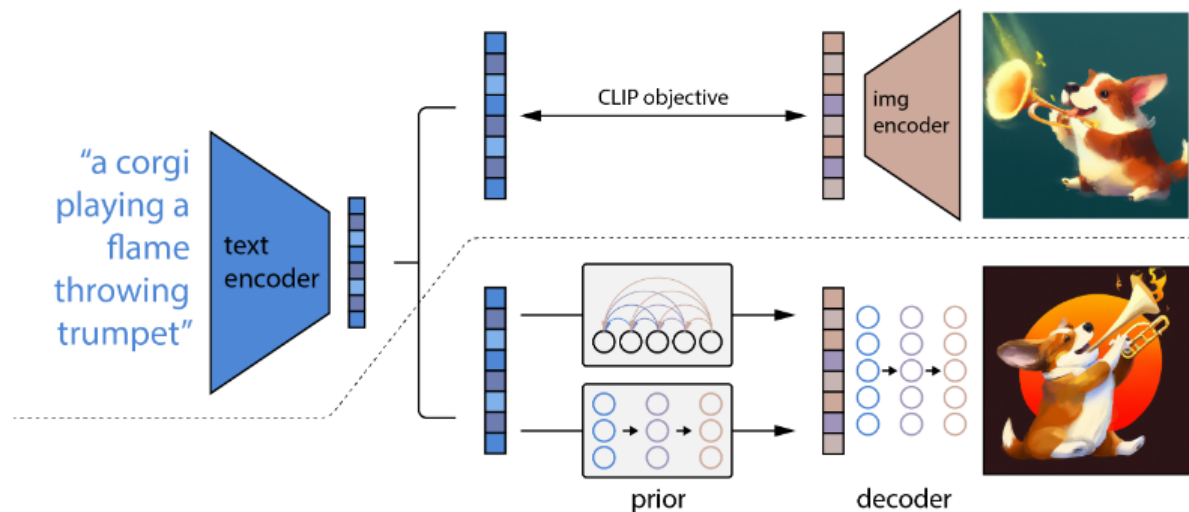
简单的classifier换成CLIP模型, 可以把文本和图像联系起来, 不仅用梯度, 也用文本控制模型的采样和生成; 图像除了像素层面的引导, 可以用特征层面的引导, 风格方面的引导; 文本可以用已经训练得很好的大语言模型引导

GLIDE:

classifier-free guidance: $f_\theta(x_t, t, y) - f_\theta(x_t, t, 0)$ 训练中生成有条件时的输出&无条件时的输出, 得到一个方向可以从无条件时的输出到有条件时的输出, 这样在无条件

时可以推测出有条件时的结果

DALLE 2



CLIP: learn representations of images

text \rightarrow text encoder \iff image encoder \leftarrow image

Training dataset: pairs (x, y) of images x and their corresponding captions y

$$P(x|y) = P(x, z_i|y) = P(x|z_i, y)P(z_i|y)$$

decoder: 给定 y 和 z_i 生成 x

prior: 给定文本 y 生成图像 embedding 的 z_i

prior (diffusion model): generates a CLIP image embedding

decoder (experiment with autoregressive \times & diffusion model \checkmark): generates an image conditioned on the image embedding

CLIP guidance & classifier-free guidance

transformer in decoder