

GLIDE(Guided Language to Image Diffusion for Generation and Editing)

ソース

- [GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models](#)

概要

- テキストから画像の生成の精度を上げた
- classifier-free guidance, clip guidanceという既存の手法に対して大量の学習を行って実験した
- classifier-free guidanceの方が精度良かった

前提知識

diffusion model

- 詳しくはDDPM, DDIM, Improved DDPMなどを見てください。

$$\begin{aligned} q(x_t|x_{t-1}) &:= \mathcal{N}(x_t; \sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I) \\ p_\theta(x_{t-1}|x_t) &:= \mathcal{N}(\mu_\theta(x_t), \Sigma_\theta(x_t)) \\ L_{\text{simple}} &:= E_{t \sim [1, T], x_0 \sim q(x_0), \epsilon \sim \mathcal{N}(0, I)} [\|\epsilon - \epsilon_\theta(x_t, t)\|^2] \end{aligned} \tag{1}$$

guided diffusion model

- 詳しくはADMを見てください

$$\hat{\mu}_\theta(x_t|y) = \mu_\theta(x_t|y) + s \cdot \Sigma_\theta(x_t|y) \nabla_{x_t} \log p_\phi(y|x_t)$$

classifier-free guidance

- 詳細はClassifier-Free Diffusion Guidance(<https://openreview.net/forum?id=qw8AKxfYbl>)
- conditional diffusion modelで条件付けを行う
- GLIDEではcaptionで重み付ける

- 特徴としては、diffusion modelの学習に条件づけを入れているために異なるguidanceをする際に再学習が必要になる

$$\hat{\epsilon}_{\theta}(x_t|y) = \epsilon_{\theta}(x_t|\emptyset) + s \cdot (\epsilon_{\theta}(x_t|y) - \epsilon_{\theta}(x_t|\emptyset))$$

- $\epsilon_{\theta}(x_t|\emptyset)$ は無条件ノイズのこと
- s は1以上の値をとる
- GLIDEではcaption c を用いて以下

$$\hat{\epsilon}_{\theta}(x_t|c) = \epsilon_{\theta}(x_t|\emptyset) + s \cdot (\epsilon_{\theta}(x_t|c) - \epsilon_{\theta}(x_t|\emptyset))$$

clip guidance

- classifier modelの損失勾配で重み付けするというADMのideaを利用
- ただし、CLIPの潜在空間での類似度を利用
 - $\log p_{\phi}(y|x_t)$ ではなく $f(x_t) \cdot g(c)$
- 2つのモデルが必要だが、同時に学習する必要はない
- CLIPが、きれいな画像だけではなくノイズ付きの画像に対しても学習されていないといけない

$$\hat{\mu}_{\theta}(x_t|y) = \mu_{\theta}(x_t|y) + s \cdot \Sigma_{\theta}(x_t|y) \nabla_{x_t} (f(x_t) \cdot g(c))$$

- f はimage encoder
- g はtext encoder
- c はcaption

学習条件

classifier-free guidance

- 35億のパラメータから成るテキスト条件付け拡散モデル(text conditional diffusion model)に対して64×64の解像度の画像を学習

clip guidance

- 64×64の解像度でVit-L CLIPで学習