# Dragondiffusion: Enabling drag-style manipulation on diffusion models

Chong Mou[1]; Xintao Wang[2]; Jiechong Song[1]; Ying Shan[1];Jian Zhang[1]

[1]Peking University Shenzhen Graduate School, Shenzhen, China

[2]ARC Lab, Tencent, PCG

https://villa.jianzhang.tech/

➢ **Stable Diffusion**



$$L_{DM} = \mathbb{E}_{x,\epsilon \sim \mathcal{N}(0,1),t} \left[ \| \epsilon - \epsilon_\theta(x_t, t) \|_2^2 \right]$$

$$L_{LDM} := \mathbb{E}_{\mathcal{E}(x),\epsilon \sim \mathcal{N}(0,1),t} \left[ \| \epsilon - \epsilon_\theta(z_t, t) \|_2^2 \right]$$

Relying on the powerful generation capabilities of SD:
➢ How can we edit existing images?
We propose DragonDiffusion.

PEKING UNIVERSITY VILLA Visual-Information Intelligent Learning LAB

ARC Applied Research Center

ICLR

SDEdit

Perturb with SDE

Reverse SDE

Stroke

Image

Stroke

Image

Input

Output

Prompt2Prompt

synthesized image        "a        furry        bear

"Photo of a cat riding on a bicycle."
car

**InstructPix2Pix:**
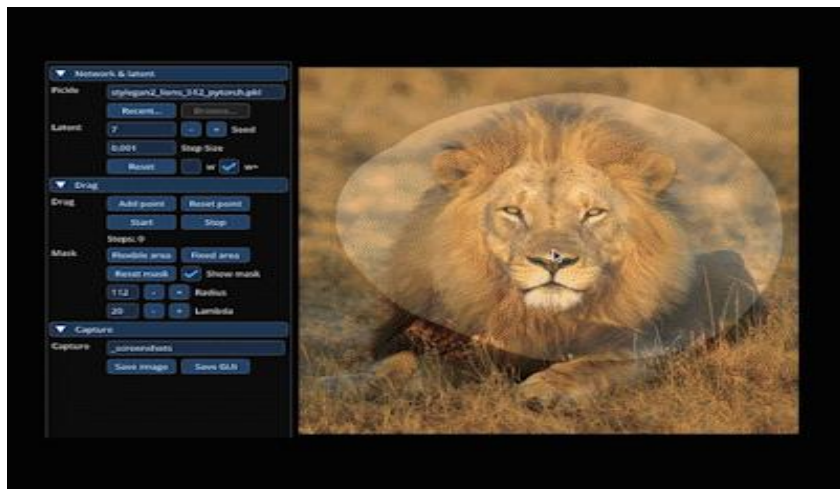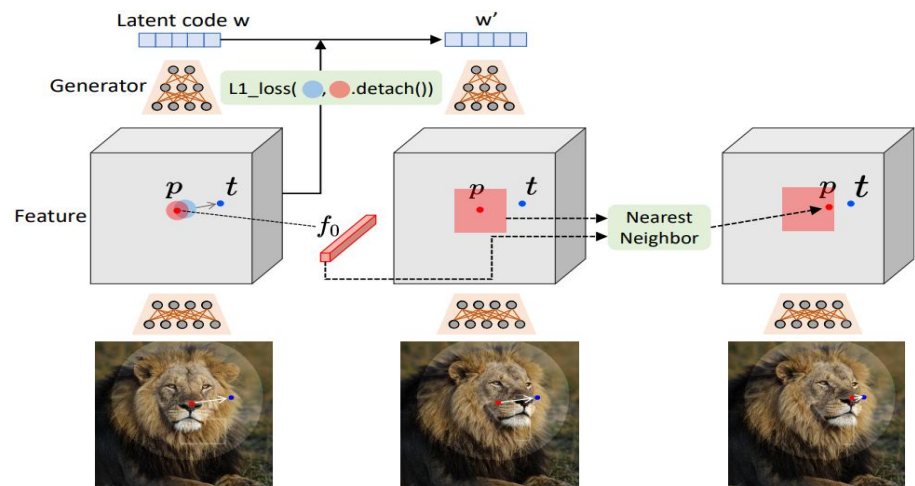


However, the correspondence between text and image features is weak, heavily relying on the design of prompts.

# Background: DragGAN

➢ **DragGAN**
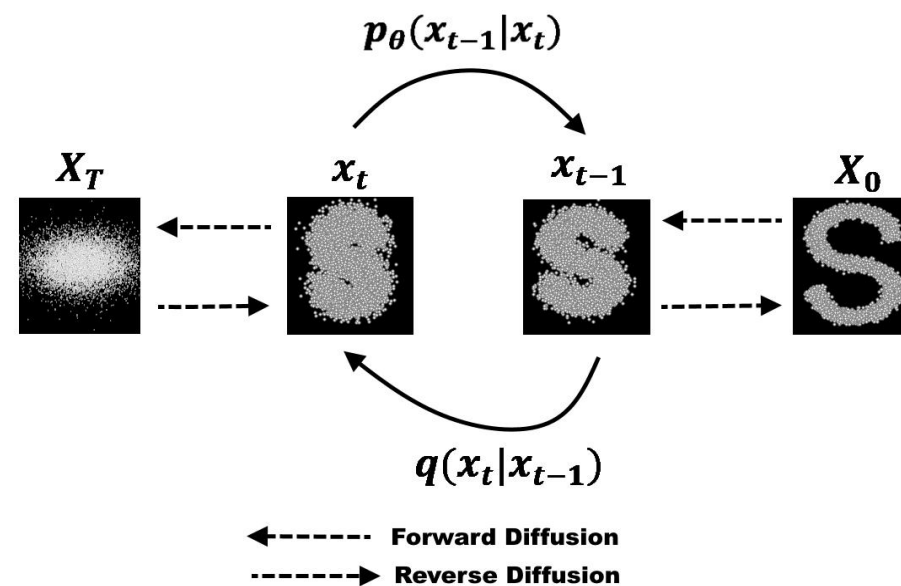


Latent code w → w'

Generator L1_loss( ●, ●.detach())

Feature

$p$ $t$ $f_0$ Nearest Neighbor


**Original Image**


**Edit w/o alignment**


**Edit w alignment**

😔 Due to the limited capability of the GANs.

**GAN model: compact and editable latent space**

$$p_\theta(x_{t-1}|x_t)$$

$X_T$  $x_t$  $x_{t-1}$  $X_0$

$$q(x_t|x_{t-1})$$

◄------- **Forward Diffusion**
------► **Reverse Diffusion**
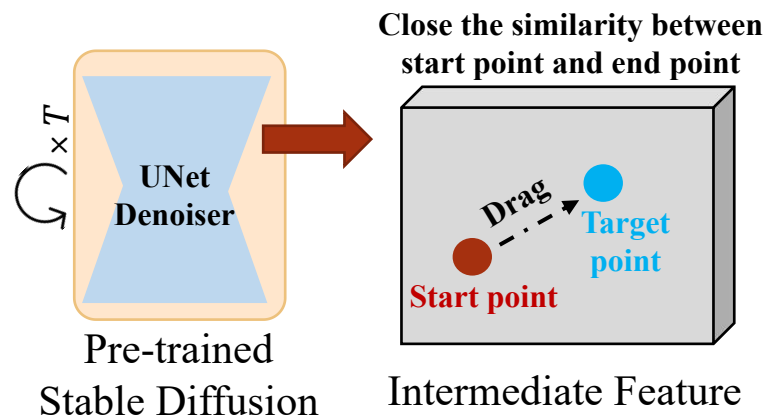
**Diffusion model: Discrete latent space**

Is there a way to perform fine-grained image editing based on SD?

**Emergent Correspondence from Image Diffusion (NeurIPS 2023)**

➤ **Score-based Guidance**

● **Editing modeling：**

Close the similarity between
start point and end point



UNet
Denoiser
×T

Pre-trained
Stable Diffusion

Intermediate Feature
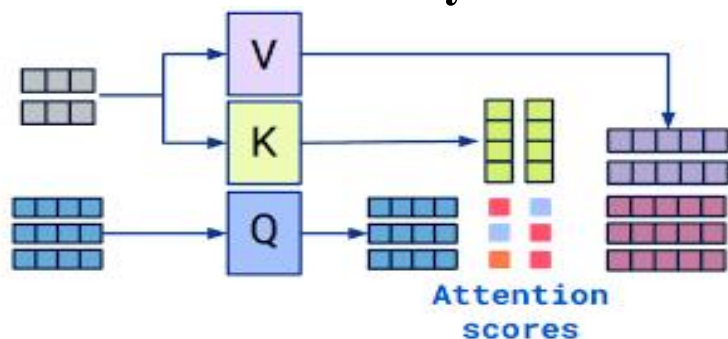
Drag

Target
point

Start point

● **Design energy function for guidance：**

$$\mathcal{E} = \underbrace{w_e \cdot \mathcal{E}_{edit}}_{\text{Editing term}} + \underbrace{w_c \cdot \mathcal{E}_{content}}_{\text{Content consistency term}}$$

$$\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t|\mathbf{y}) = \nabla_{\mathbf{x}_t} \log \left( \frac{q(\mathbf{y}|\mathbf{x}_t)q(\mathbf{x}_t)}{q(\mathbf{y})} \right)$$

$$\propto \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log q(\mathbf{y}|\mathbf{x}_t),$$

$$\tilde{\boldsymbol{\epsilon}}_\theta^t(\mathbf{x}_t) = \boldsymbol{\epsilon}_\theta^t(\mathbf{x}_t) + \eta \cdot \nabla_{\mathbf{x}_t} \mathcal{E}(\mathbf{x}_t, \mathbf{y}),$$
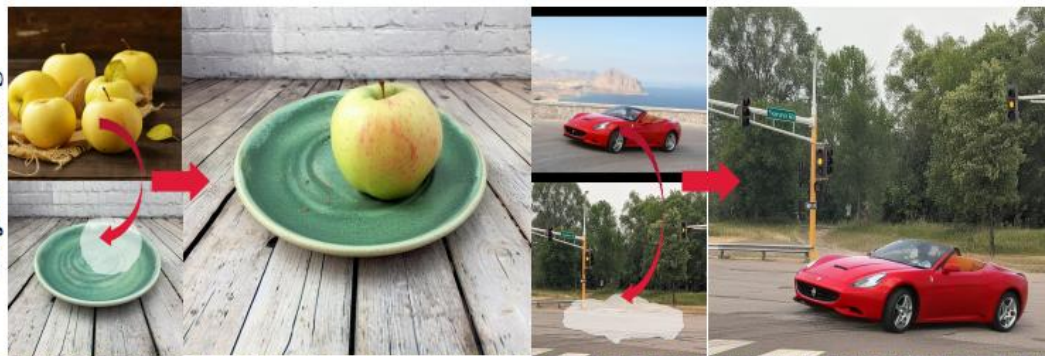
➤ **Content consistency via visual cross-attention**



V

K

Q

Attention
scores

**Key and Value are the diffusion feature from
the reference image.**

# Thanks!

Chong Mou[1]; Xintao Wang[2]; Jiechong Song[1]; Ying Shan[1];Jian Zhang[1]

[1]Peking University Shenzhen Graduate School, Shenzhen, China

[2]ARC Lab, Tencent, PCG