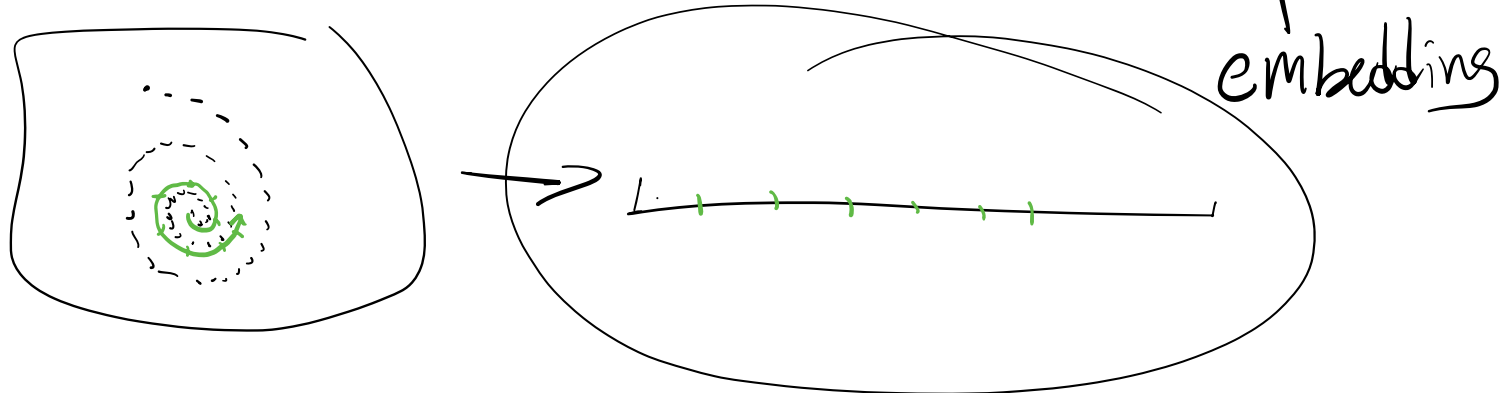# (Sharp?) left turn:
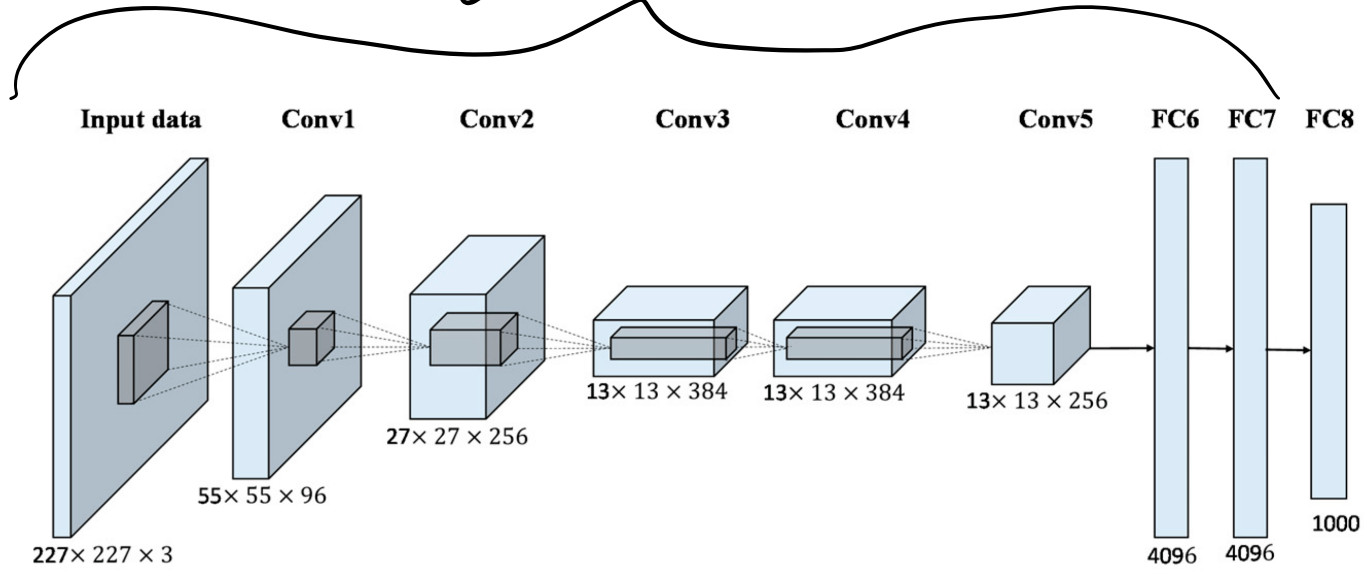# Embeddings, Manifold Learning, and Autoencoders

embedding model



| Input data | Conv1 | Conv2 | Conv3 | Conv4 | Conv5 | FC6 | FC7 | FC8 |

$13 \times 13 \times 384$     $13 \times 13 \times 384$     $13 \times 13 \times 256$

$27 \times 27 \times 256$

$55 \times 55 \times 96$

$227 \times 227 \times 3$

4096     4096     1000

embedding

embedding
latent space

E

D

L

# Generative Modeling

Disc: $p(y|x)$

Gen: $p(x,y)$

Sample

embedding
latent space

E

D

rand
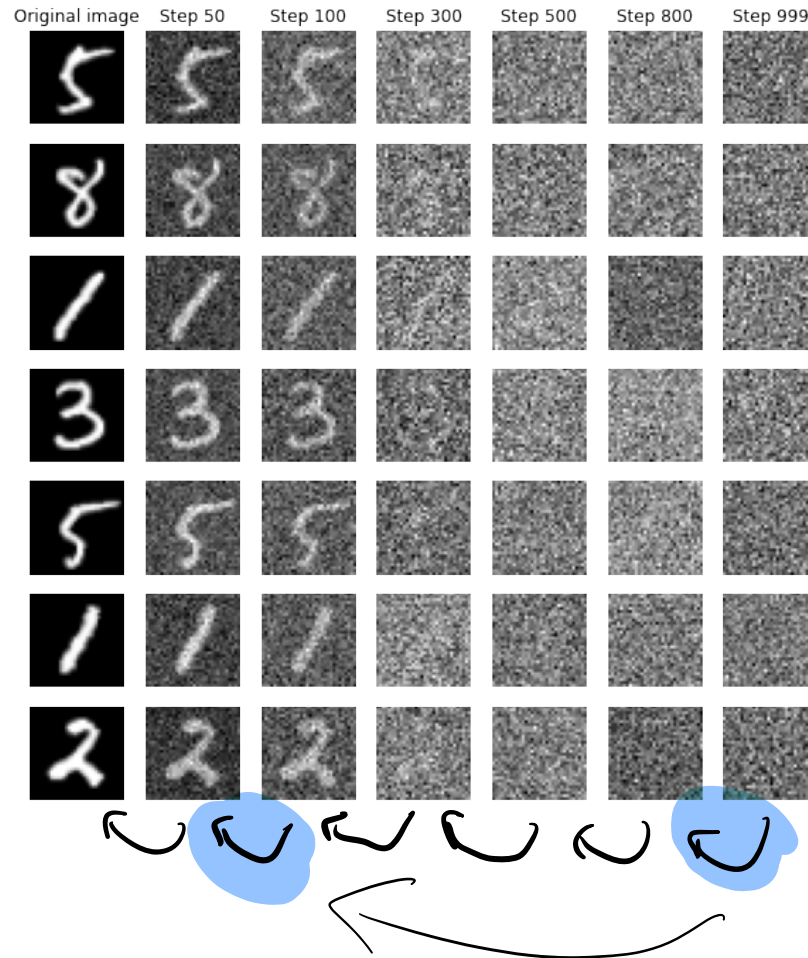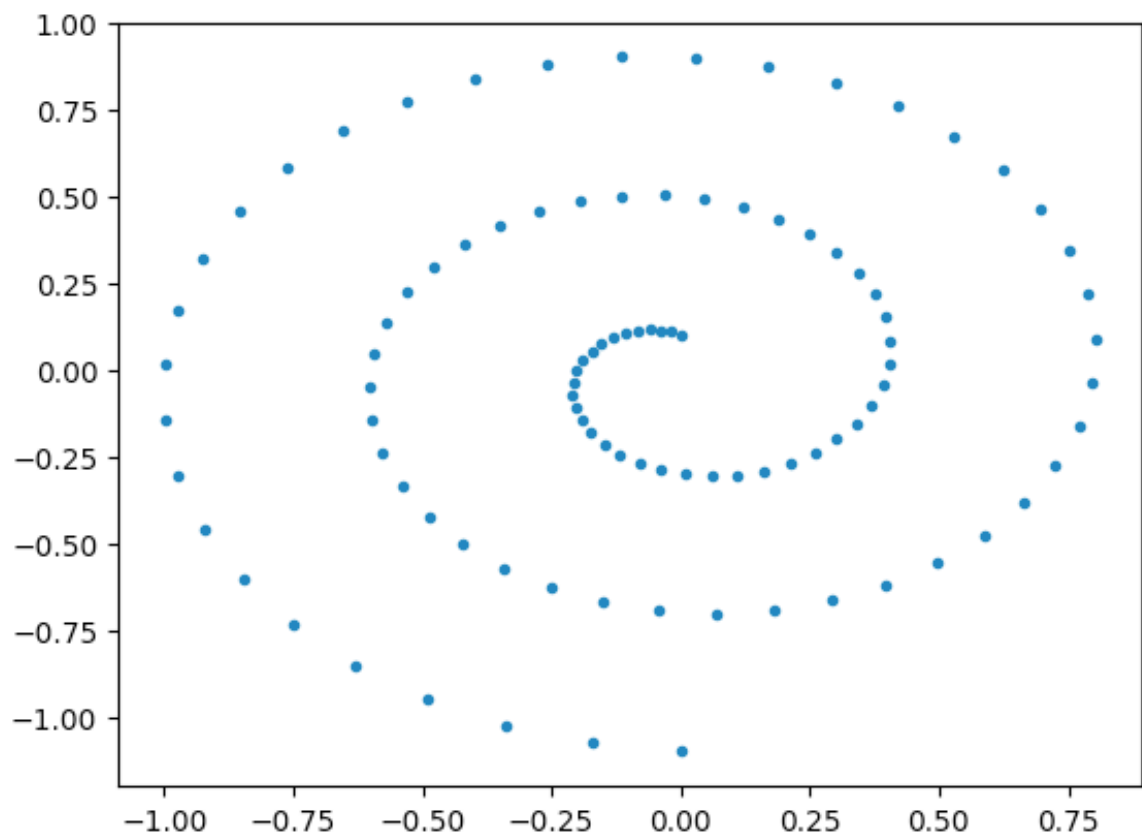
# Generative Adversarial Networks

# Diffusion Models



Some other good visuals: https://www.chenyang.co/diffusion.html
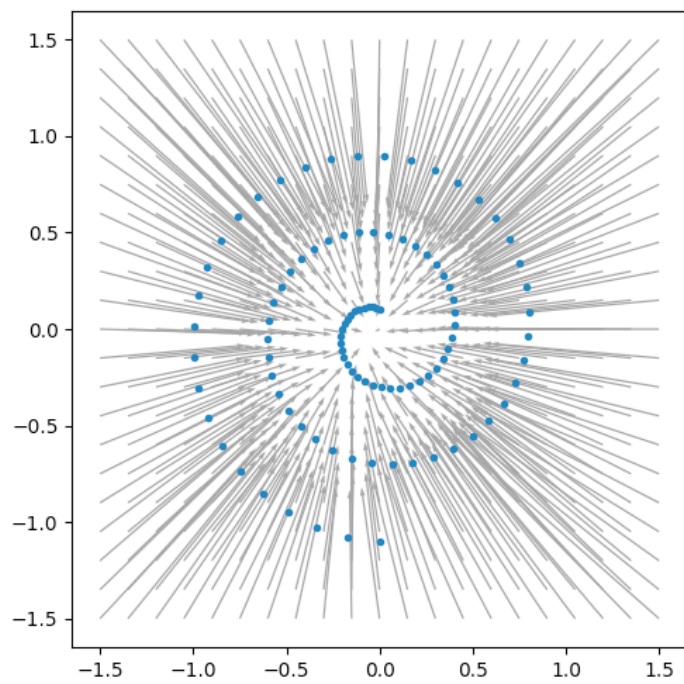
$\sigma = 0.1$          $\sigma = 0.5$          $\sigma = 1$

# UNet - a more detailed picture

# Stable Diffusion
## (without the conditioning)



Sample

denoised latent

**Pixel Space**

$x$   $\mathcal{E}$   $z$

**Latent Space**

Diffusion Process

$z_T$

Denoising U-Net $\epsilon_\theta$

$\times (T-1)$

$z$   $z_{T-1}$

$Q$ $KV$   $Q$ $KV$   $Q$ $KV$   $Q$ $KV$

$z_T$

$\tilde{x}$   $\mathcal{D}$

$\tau_\theta$

**Conditioning**

Semantic Map

Text

Repres entations

Images

denoising step    $Q$ $KV$ crossattention    switch    skip connection    concat

# Vision and Language

# Case study: CLIP



(1) Contrastive pre-training

Pepper the aussie pup → Text Encoder → $T_1$, $T_2$, $T_3$, ..., $T_N$

Image Encoder → $I_1$, $I_2$, $I_3$, ⋮, $I_N$

| | $T_1$ | $T_2$ | $T_3$ | ... | $T_N$ |
|---|---|---|---|---|---|
| $I_1$ | $I_1 \cdot T_1$ | $I_1 \cdot T_2$ | $I_1 \cdot T_3$ | ... | $I_1 \cdot T_N$ |
| $I_2$ | $I_2 \cdot T_1$ | $I_2 \cdot T_2$ | $I_2 \cdot T_3$ | ... | $I_2 \cdot T_N$ |
| $I_3$ | $I_3 \cdot T_1$ | $I_3 \cdot T_2$ | $I_3 \cdot T_3$ | ... | $I_3 \cdot T_N$ |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋱ | ⋮ |
| $I_N$ | $I_N \cdot T_1$ | $I_N \cdot T_2$ | $I_N \cdot T_3$ | ... | $I_N \cdot T_N$ |

# "Attention"

$$The \quad t_1 \quad \leftarrow \quad t_i^? \quad = \quad \sum w_i \, f(t_i)$$

dos $t_2$

is $t_3$

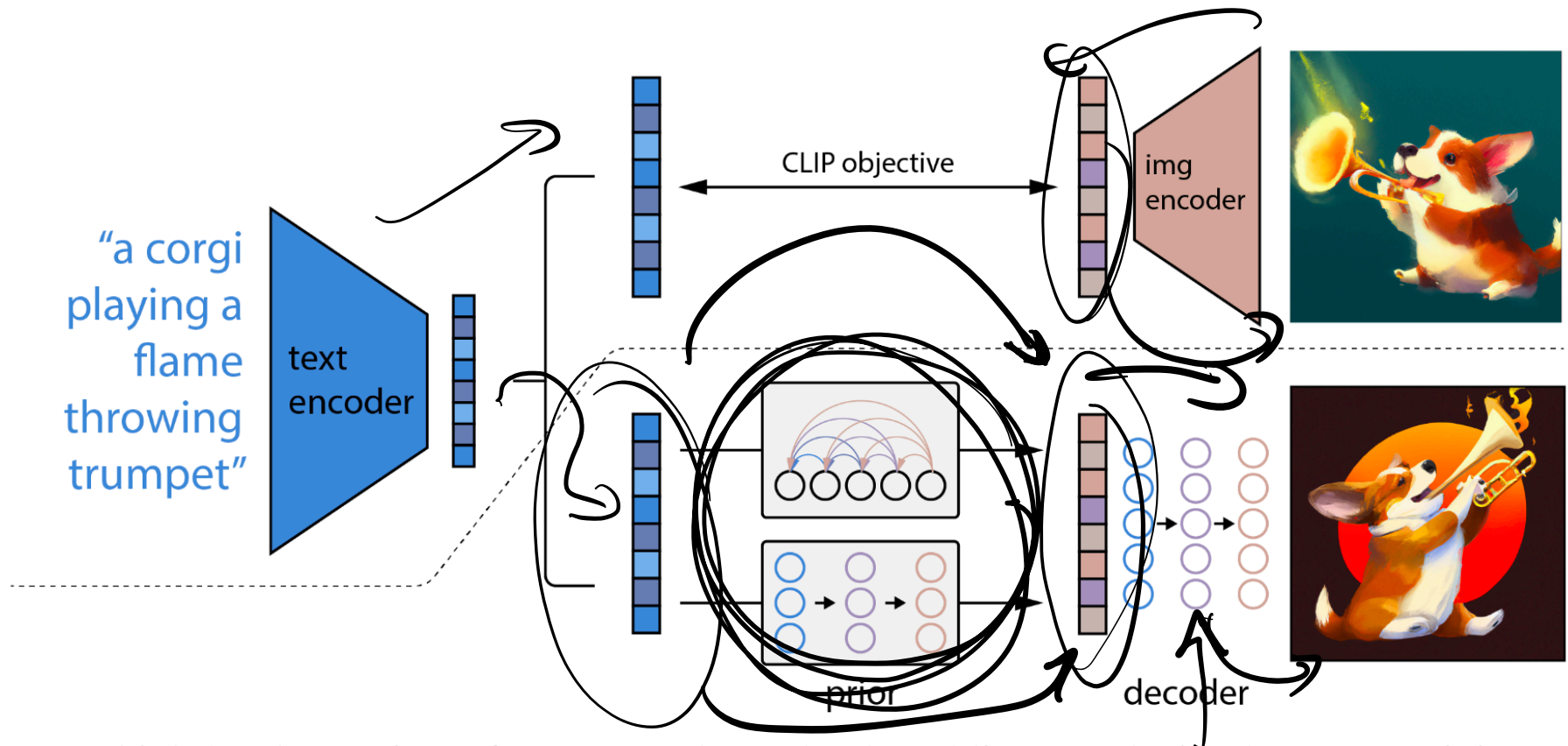cute $t_4$

# unCLIP aka DALL-E 2



Figure 2: A high-level overview of unCLIP. Above the dotted line, we depict the CLIP training process, through which we learn a joint representation space for text and images. Below the dotted line, we depict our text-to-image generation process: a CLIP text embedding is first fed to an autoregressive or diffusion prior to produce an image embedding, and then this embedding is used to condition a diffusion decoder which produces a final image. Note that the CLIP model is frozen during training of the prior and decoder.

# Stable Diffusion
# (with the conditioning)