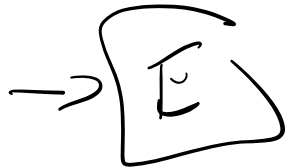# (Sharp?) left turn:
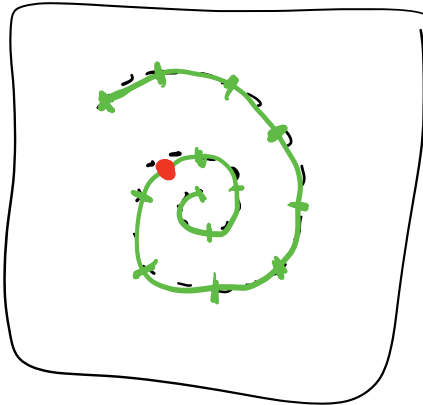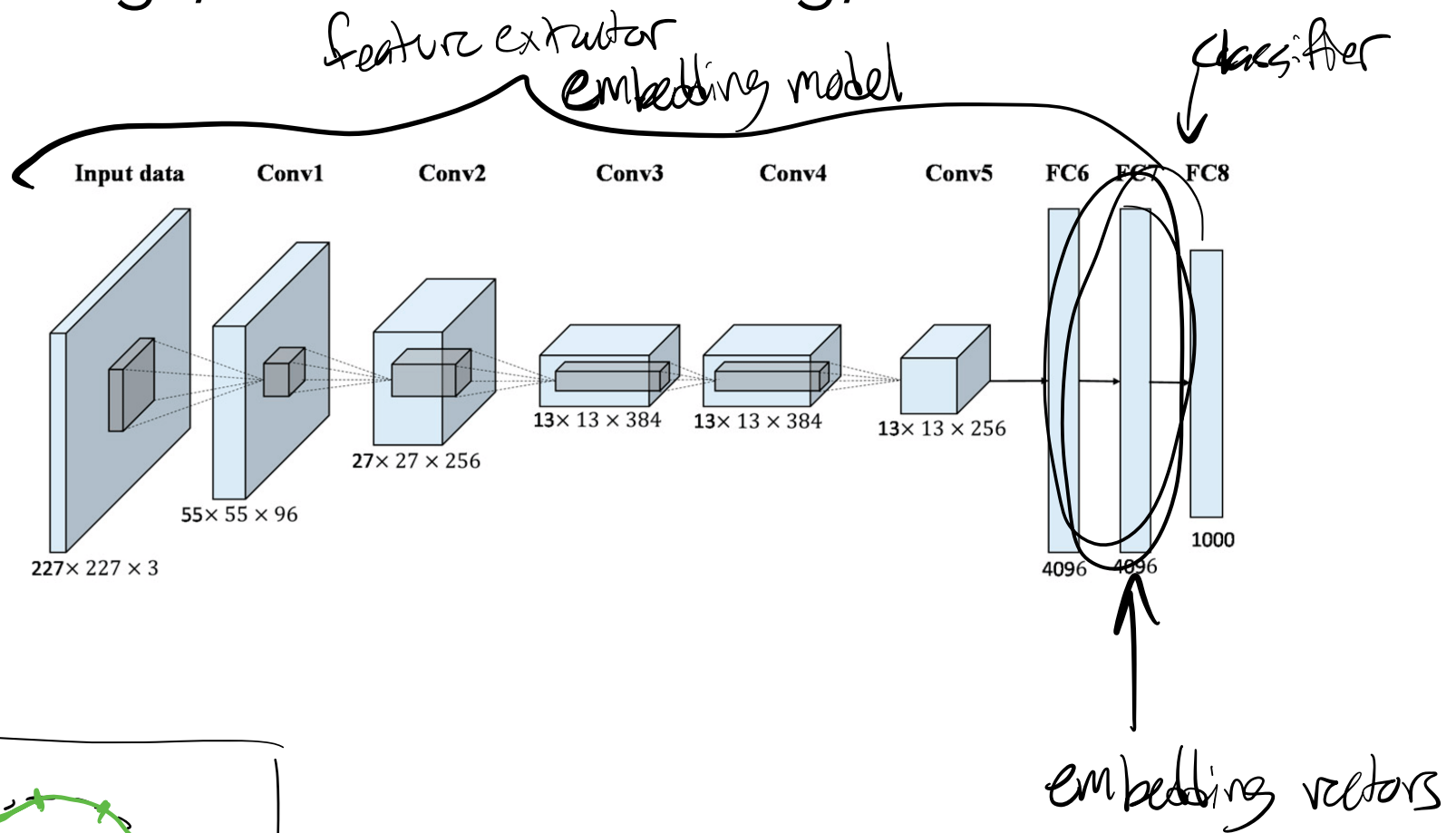# ✓ Embeddings, Manifold Learning, and Autoencoders ✓



feature extractor

embedding model

classifier

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

$27 \times 27 \times 256$

$13 \times 13 \times 384$

$13 \times 13 \times 384$

$13 \times 13 \times 256$

$55 \times 55 \times 96$

$227 \times 227 \times 3$

4096   4096   1000

embedding vectors

embedding vectors

encoder

latent / embedding vector

decoder

E

D

4096

Sample

# Generative Modeling

Dis: $p(y \mid x)$

Gen: $p(x, y)$

# Generative Adversarial Networks

# Diffusion Models

$\sigma = 0.1$      $\sigma = 0.5$      $\sigma = 1$
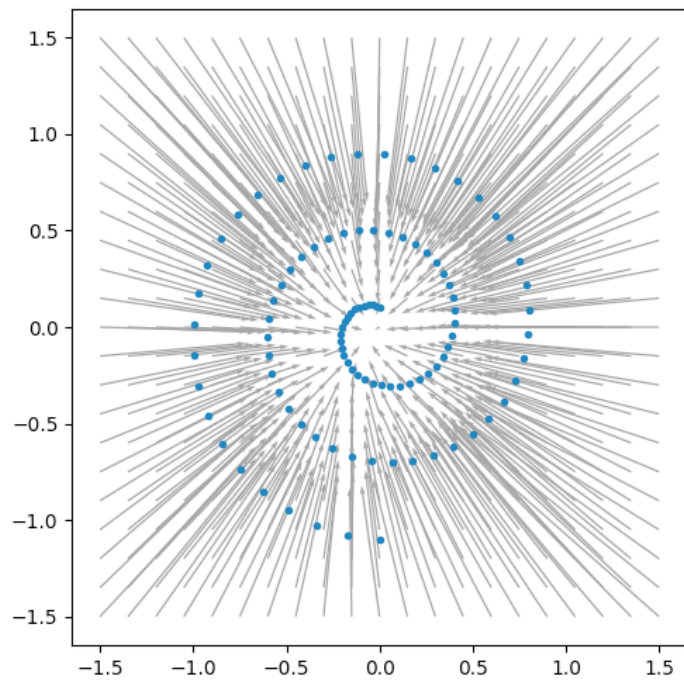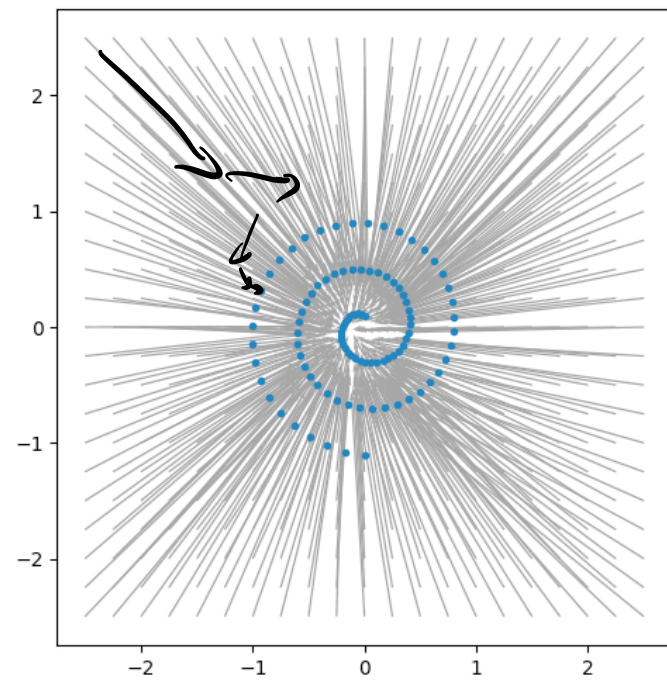
# UNet - a more detailed picture

# Stable Diffusion
## (without the conditioning)

**Pixel Space**

$x$    $\mathcal{E}$

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

**Latent Space**

Diffusion Process

$z$                    $z_T$

Denoising U-Net $\epsilon_\theta$

$z_{T-1}$

$z$     $z_{T-1}$                 $z_T$

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

**Conditioning**

Semantic Map

Text

Repres entations

Images

$\tau_\theta$

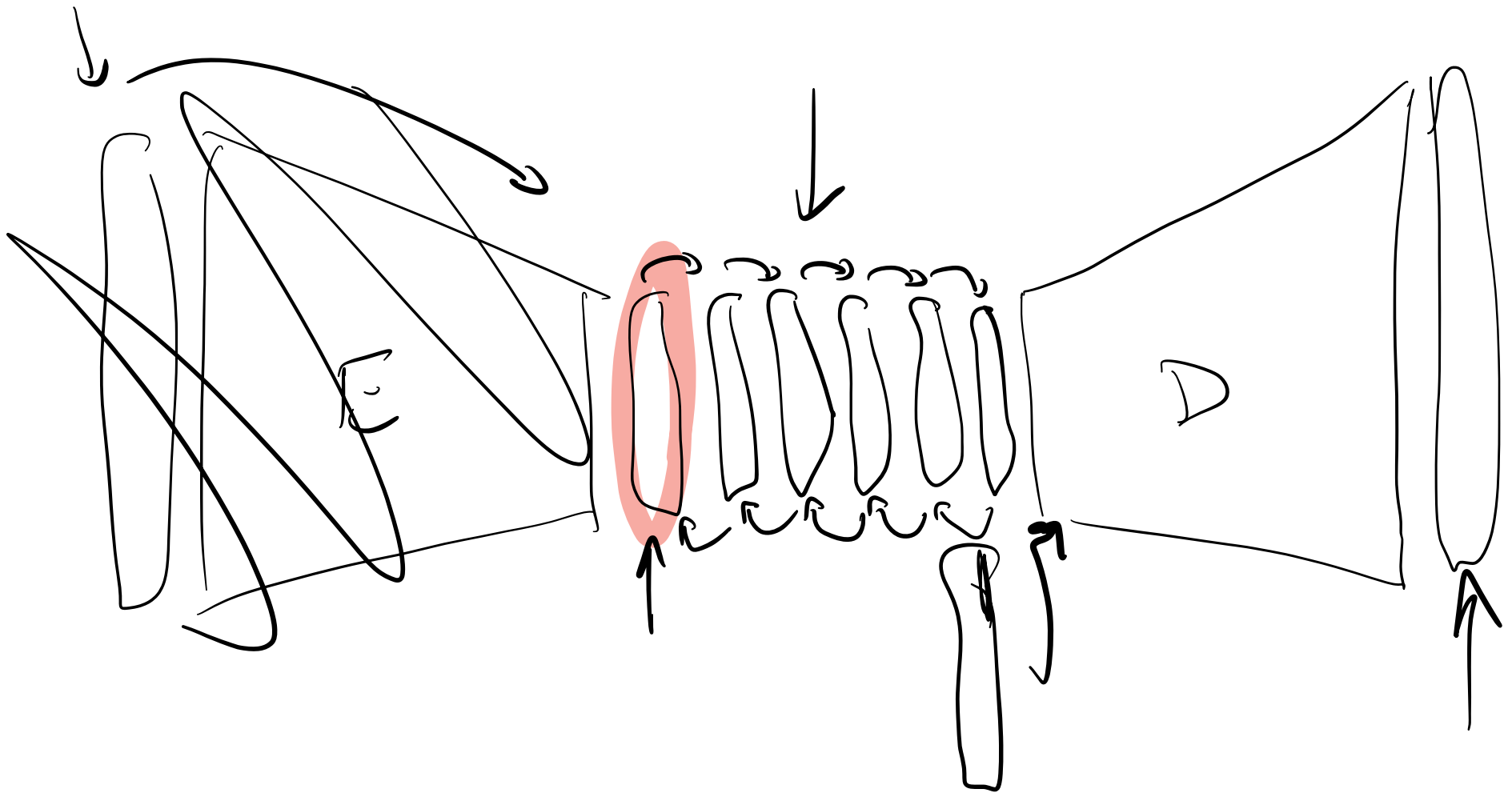denoising step    crossattention    switch    skip connection    concat

$Q$
$KV$

# Vision and Language

# Case study: CLIP



**(1) Contrastive pre-training**

# "Attention"



This      dog      is      cute

$Q, K, V$

$W_1$      $W_2$      $W_3$      $W_4$

Linear      Linear
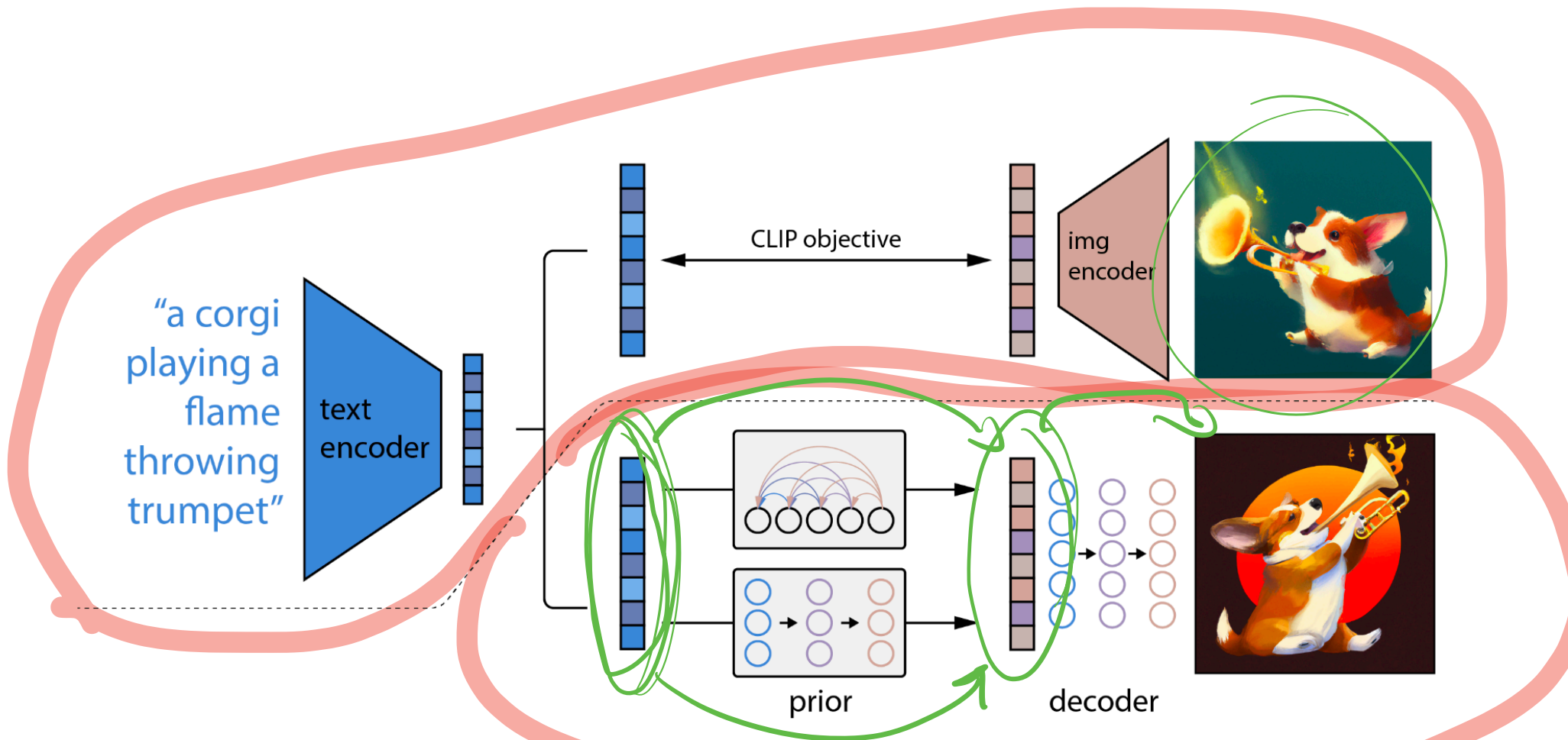
# unCLIP aka DALL-E 2



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)