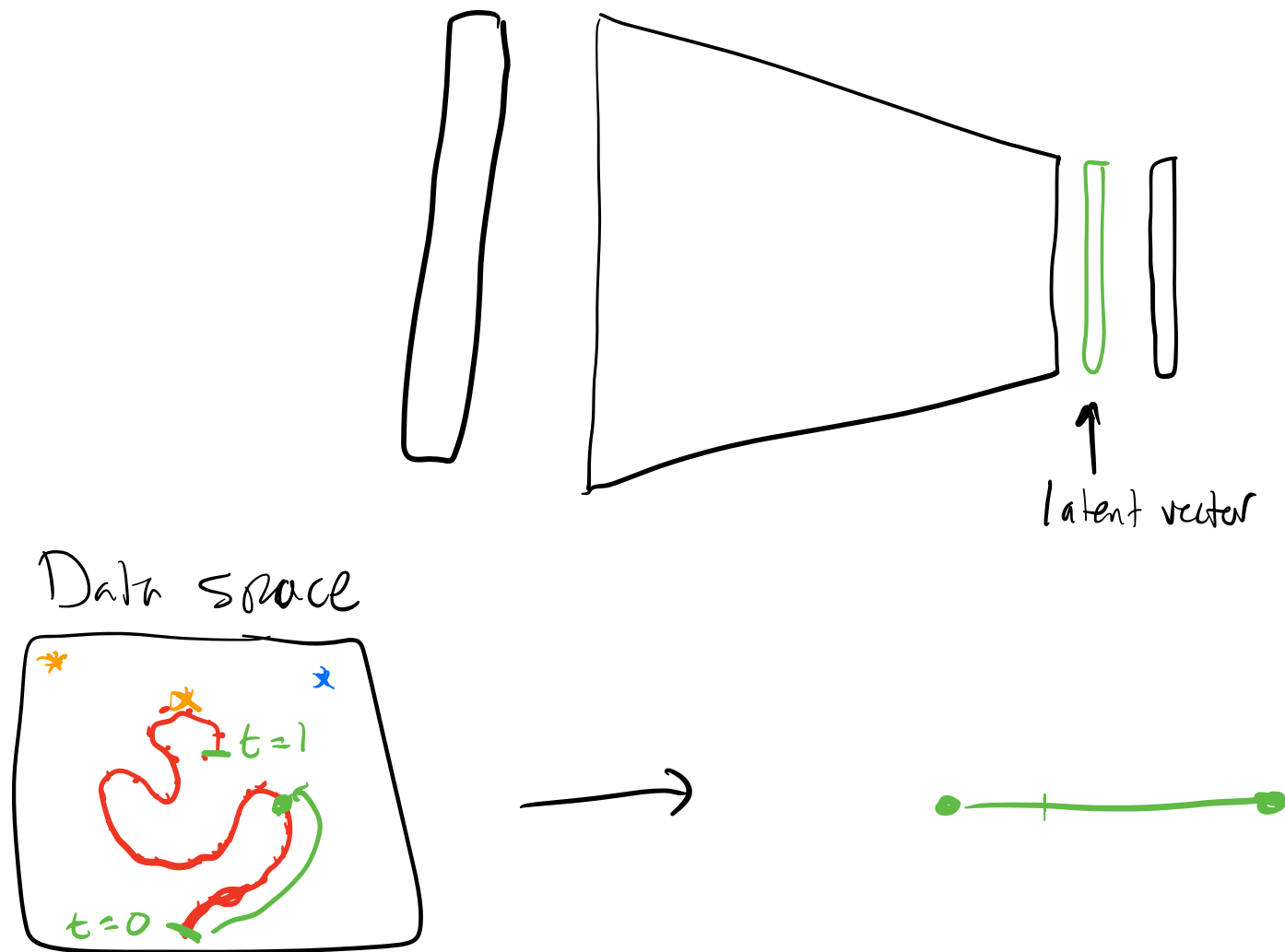


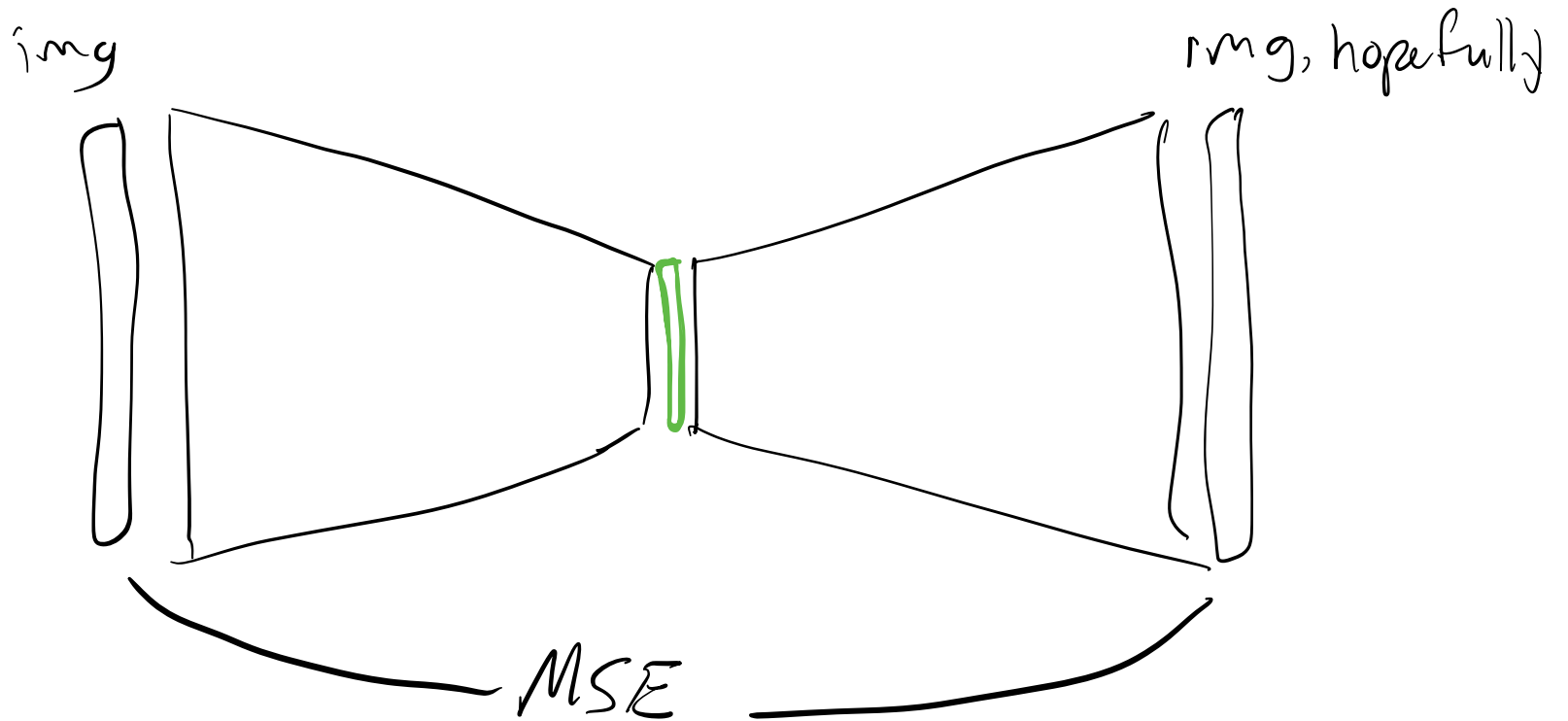


Okay but the data...

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



Auto encoder



Unsupervised / self-supervised learning case study: SimCLR

A Simple Framework for Contrastive Learning of Visual Representations



(a) Original



(b) Crop and resize



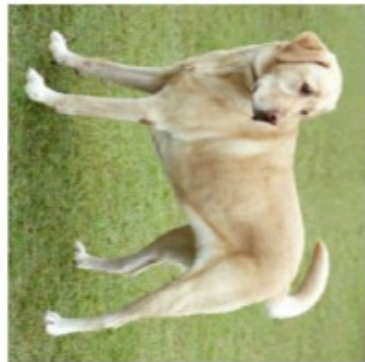
(c) Crop, resize (and flip)



(d) Color distort. (drop)



(e) Color distort. (jitter)



(f) Rotate $\{90^\circ, 180^\circ, 270^\circ\}$



(g) Cutout



(h) Gaussian noise



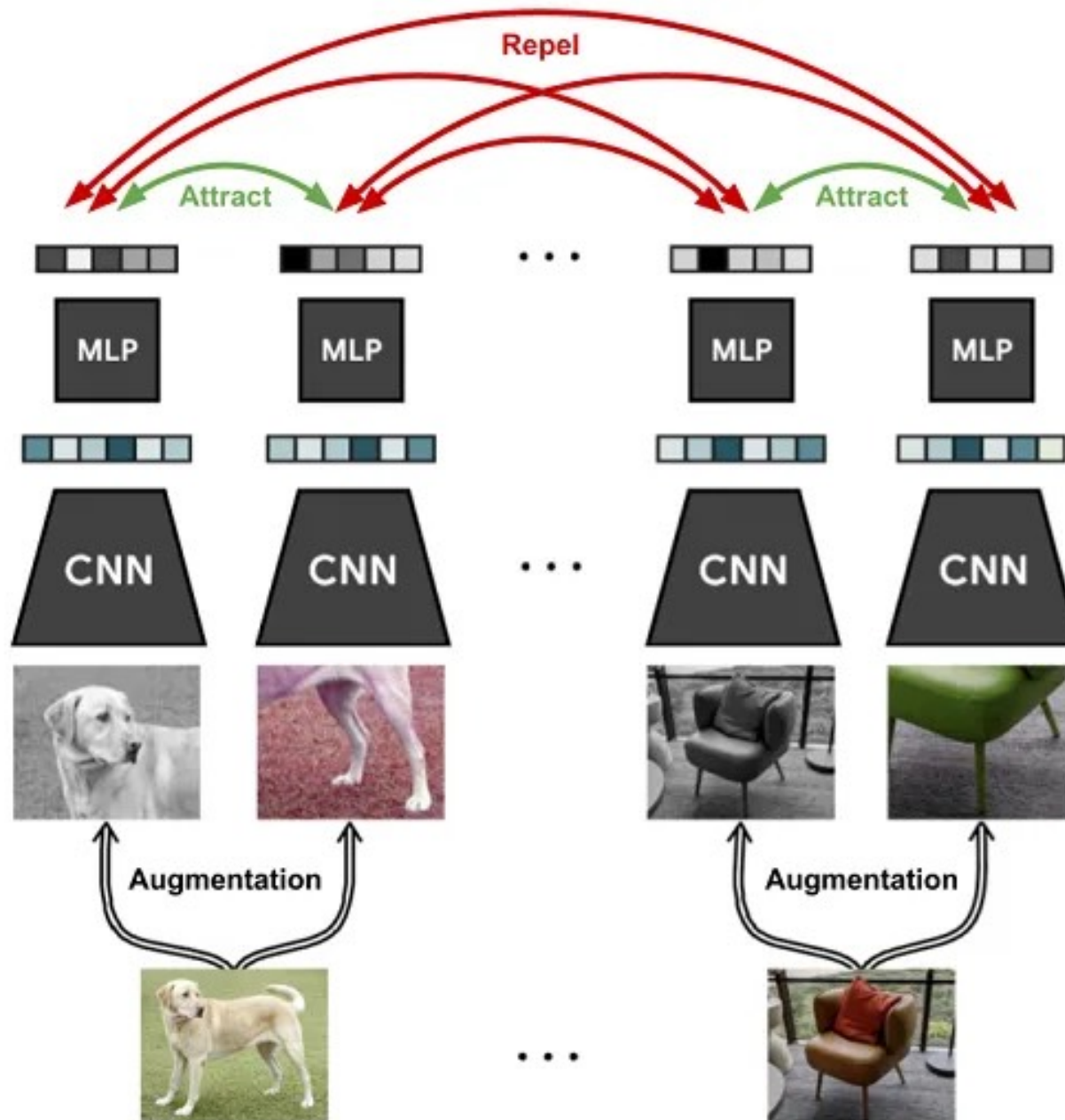
(i) Gaussian blur



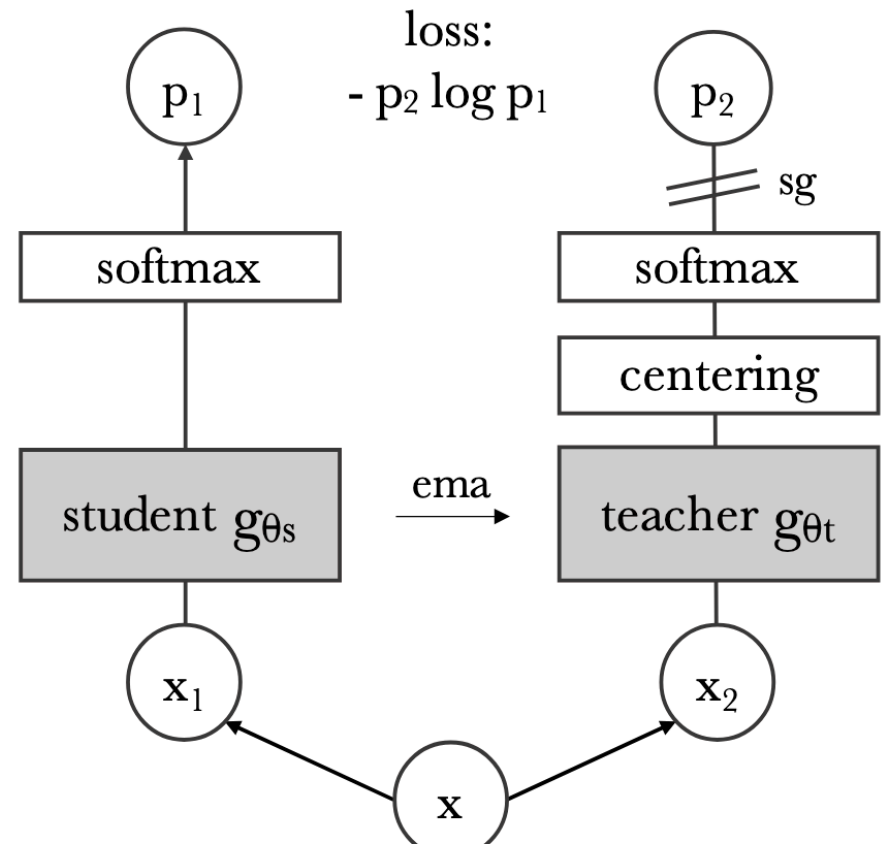
(j) Sobel filtering

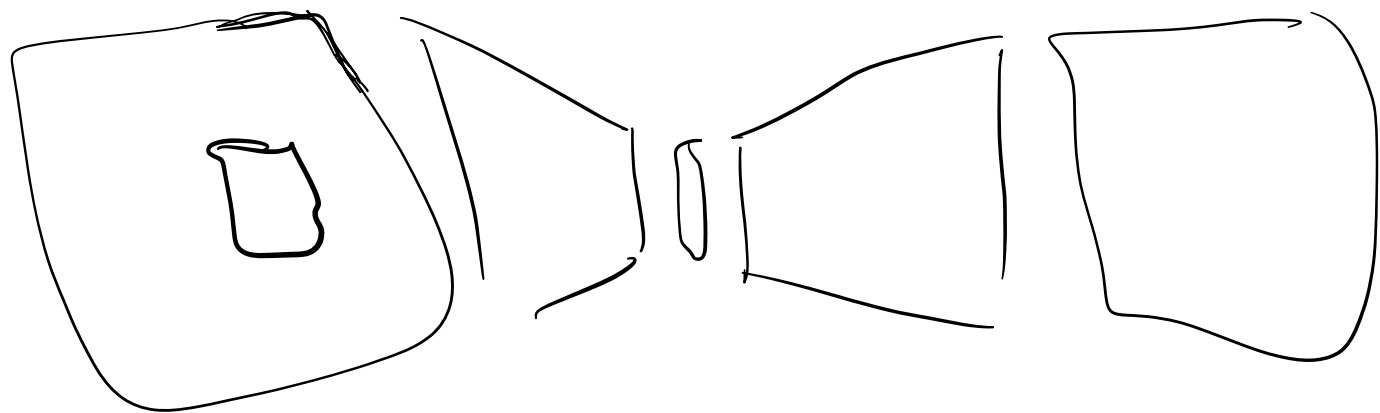
Figure 4. Illustrations of the studied data augmentation operators. Each augmentation can transform data stochastically with some internal parameters (e.g. rotation degree, noise level). Note that we *only* test these operators in ablation, the *augmentation policy* used to train our models only includes *random crop* (with *flip* and *resize*), *color distortion*, and *Gaussian blur*. (Original image cc-by: Von.grzanka)

Unsupervised / self-supervised learning case study: SimCLR



DINO - Self-Supervised Learning Image Features

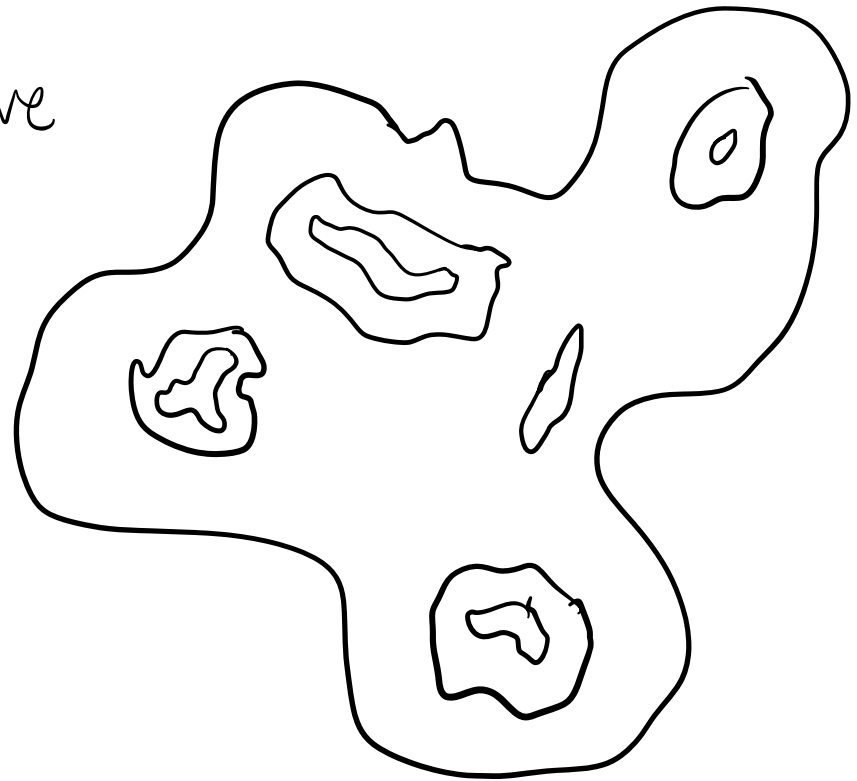




Generative Modeling

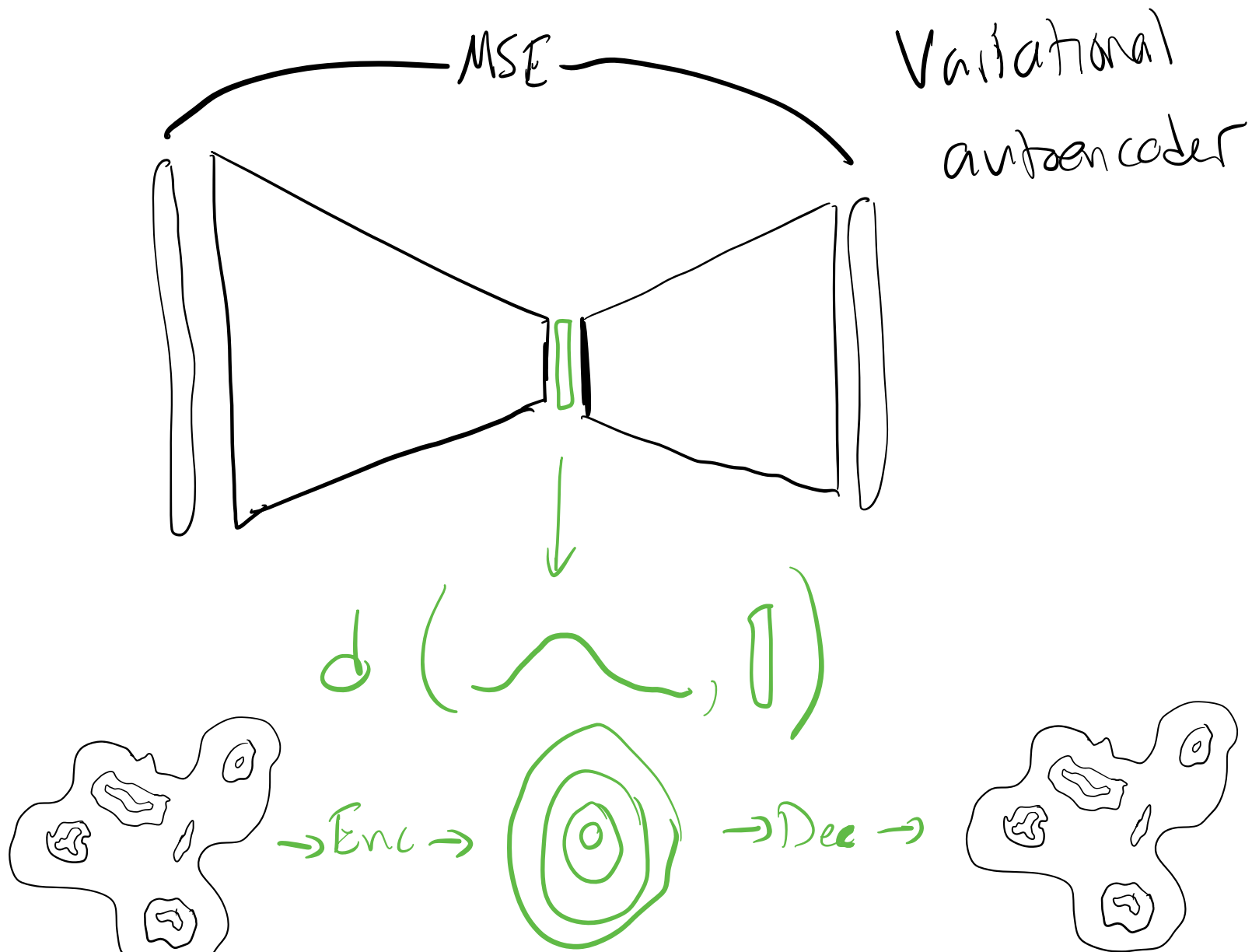
$P(y|x)$ discriminative

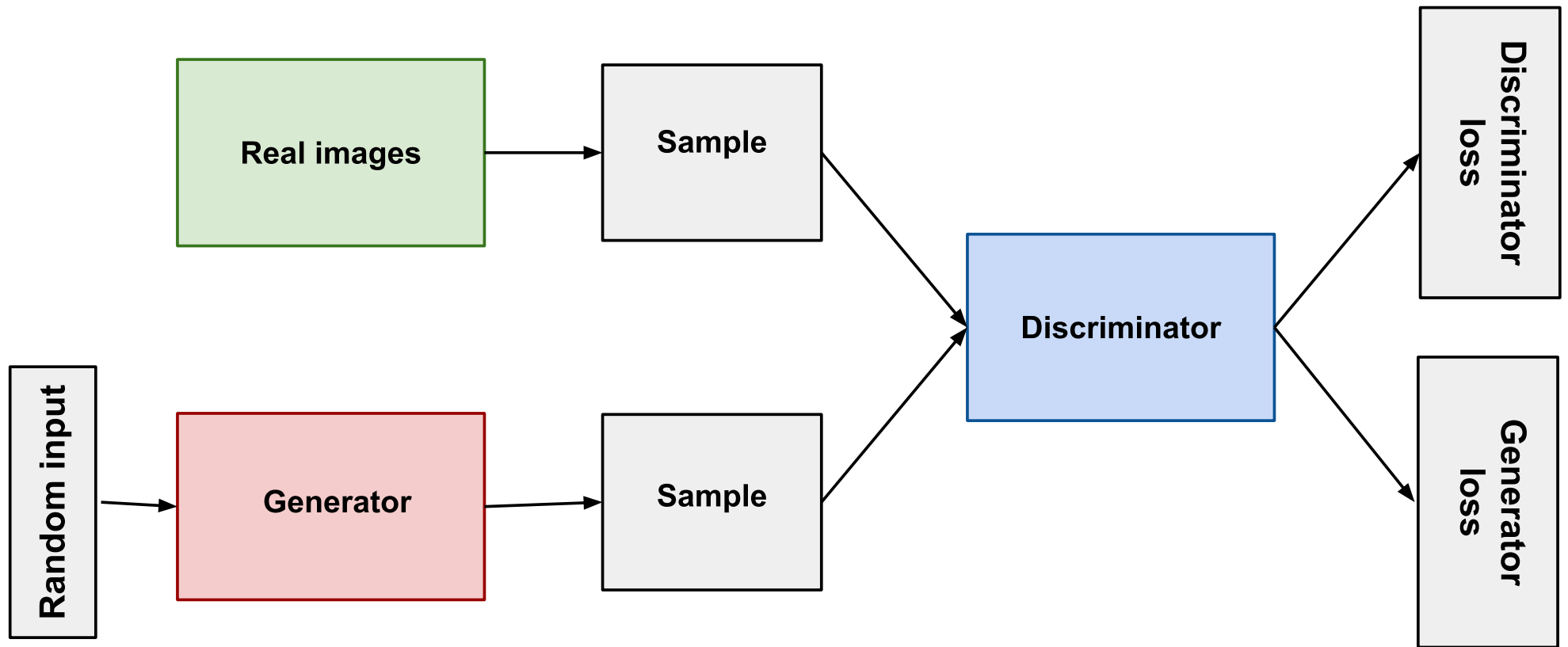
$P(x,y)$ generative

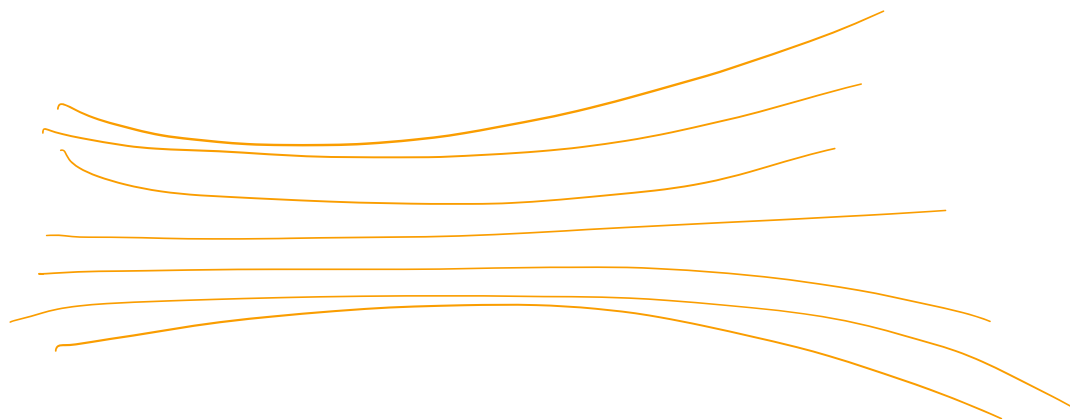
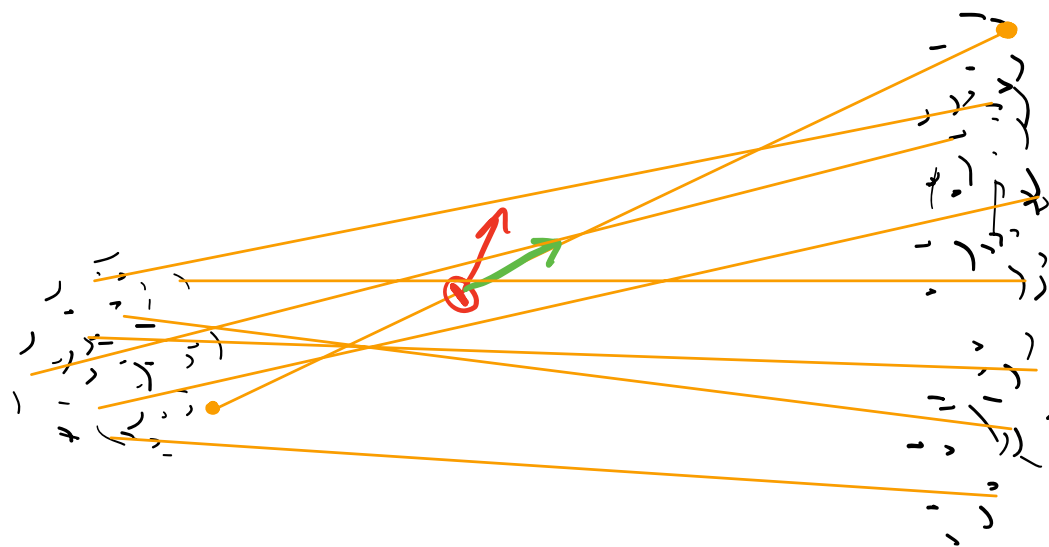


$P(\mathbf{x}, y)$

...with autoencoders

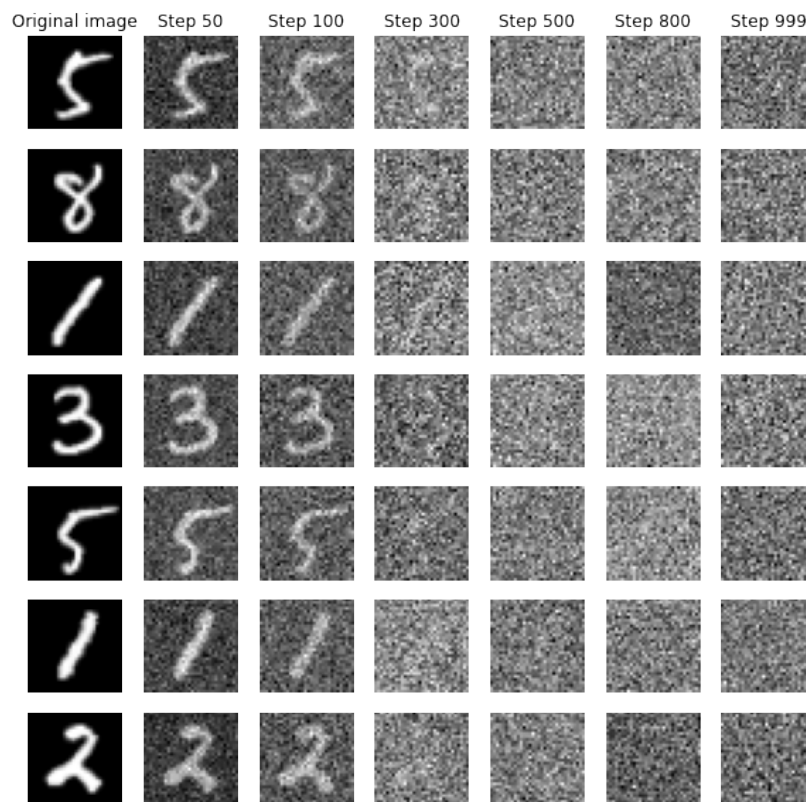




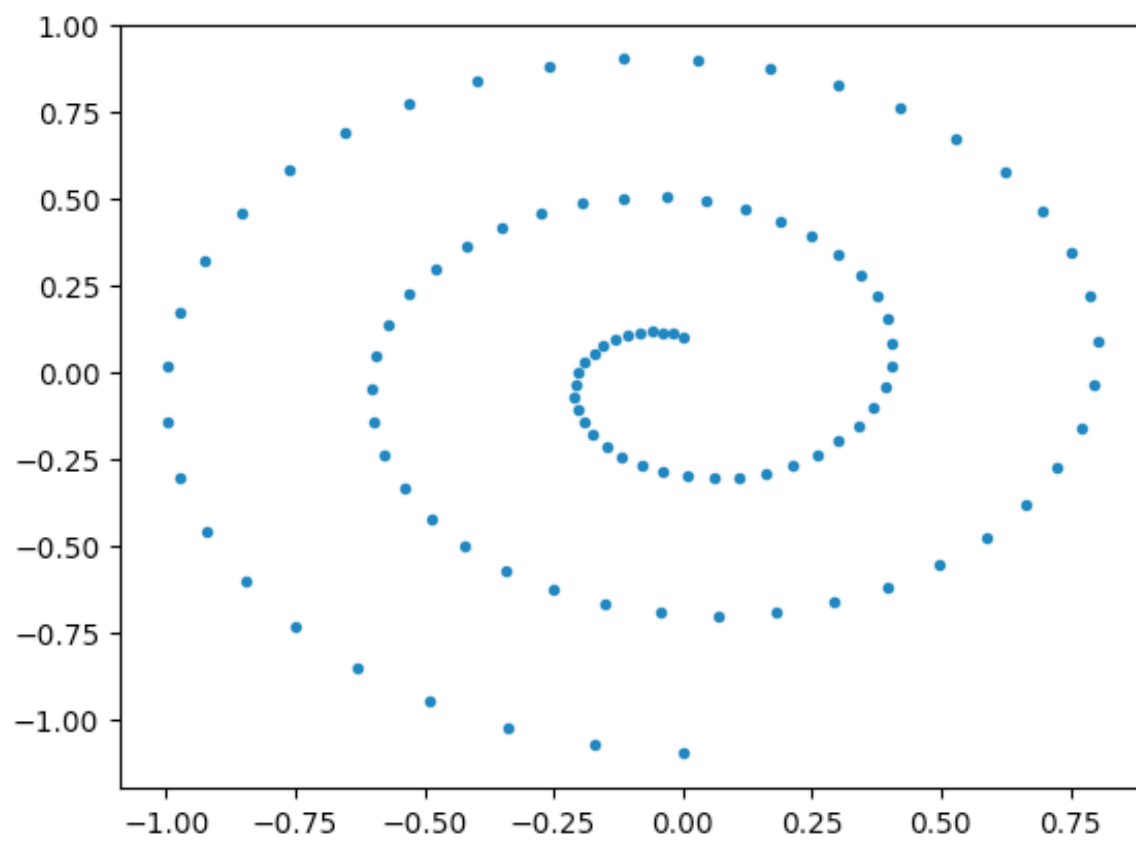




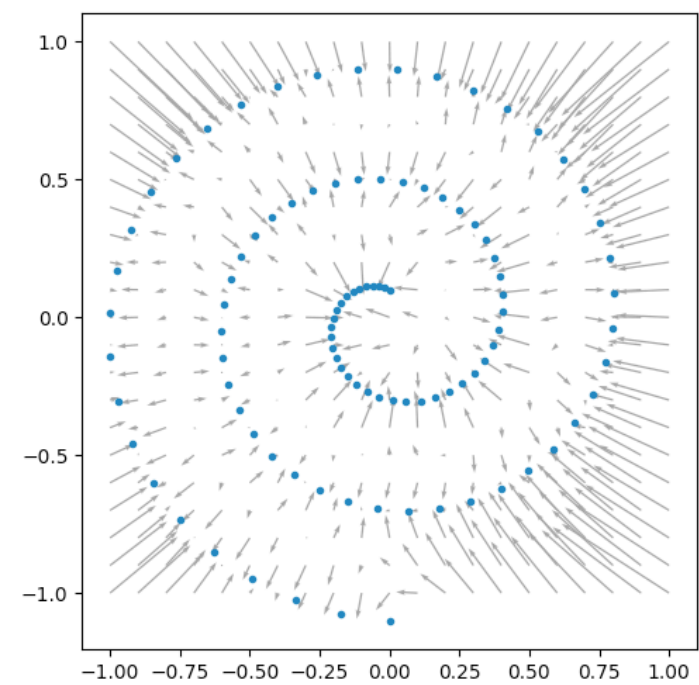
Diffusion Models



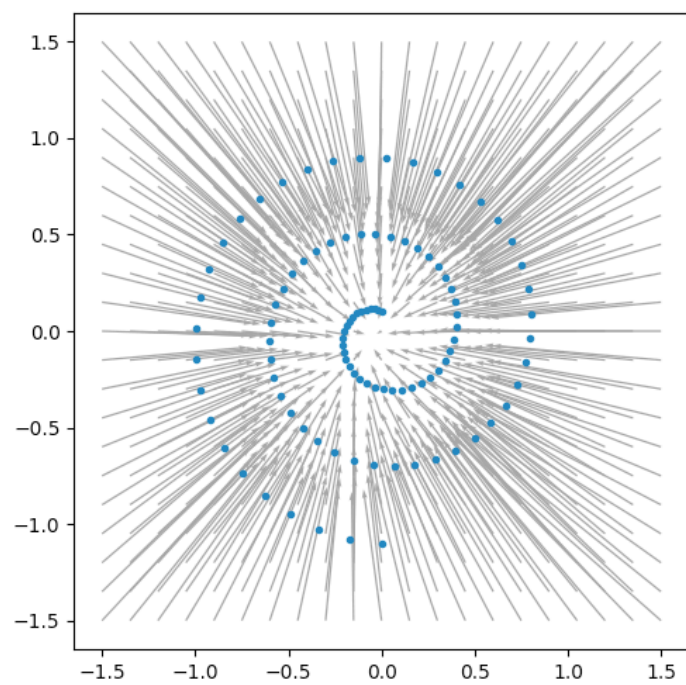
Some other good visuals: <https://www.chenyang.co/diffusion.html>



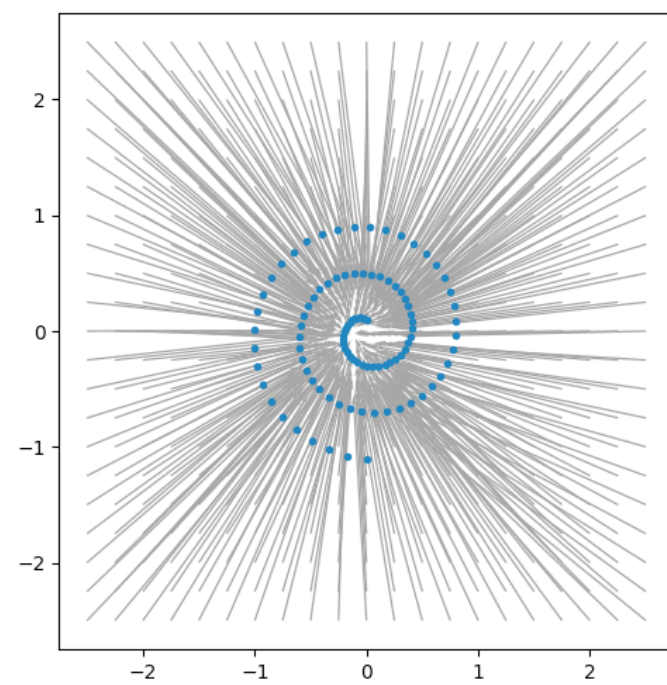
$\sigma = 0.1$



$\sigma = 0.5$



$\sigma = 1$

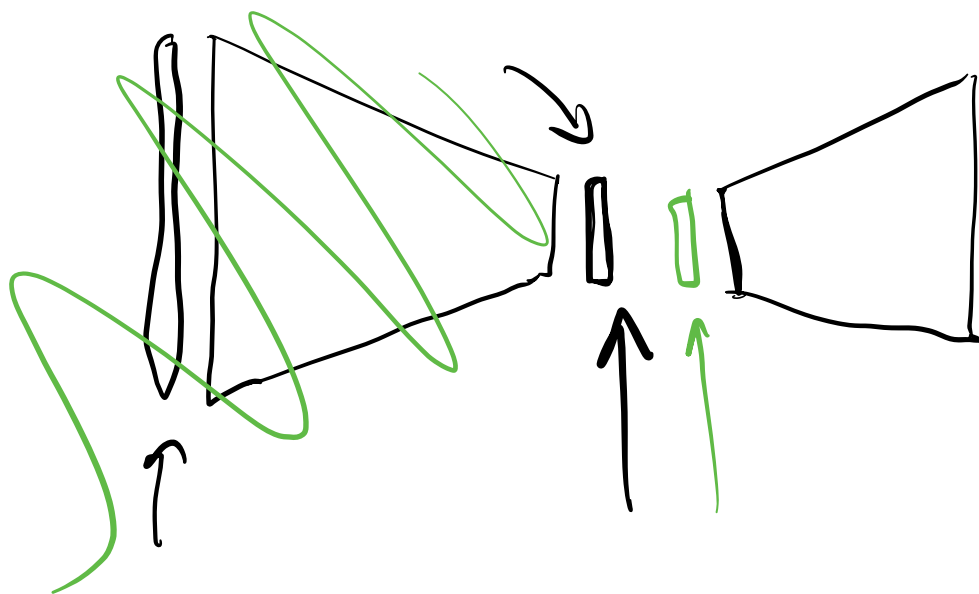


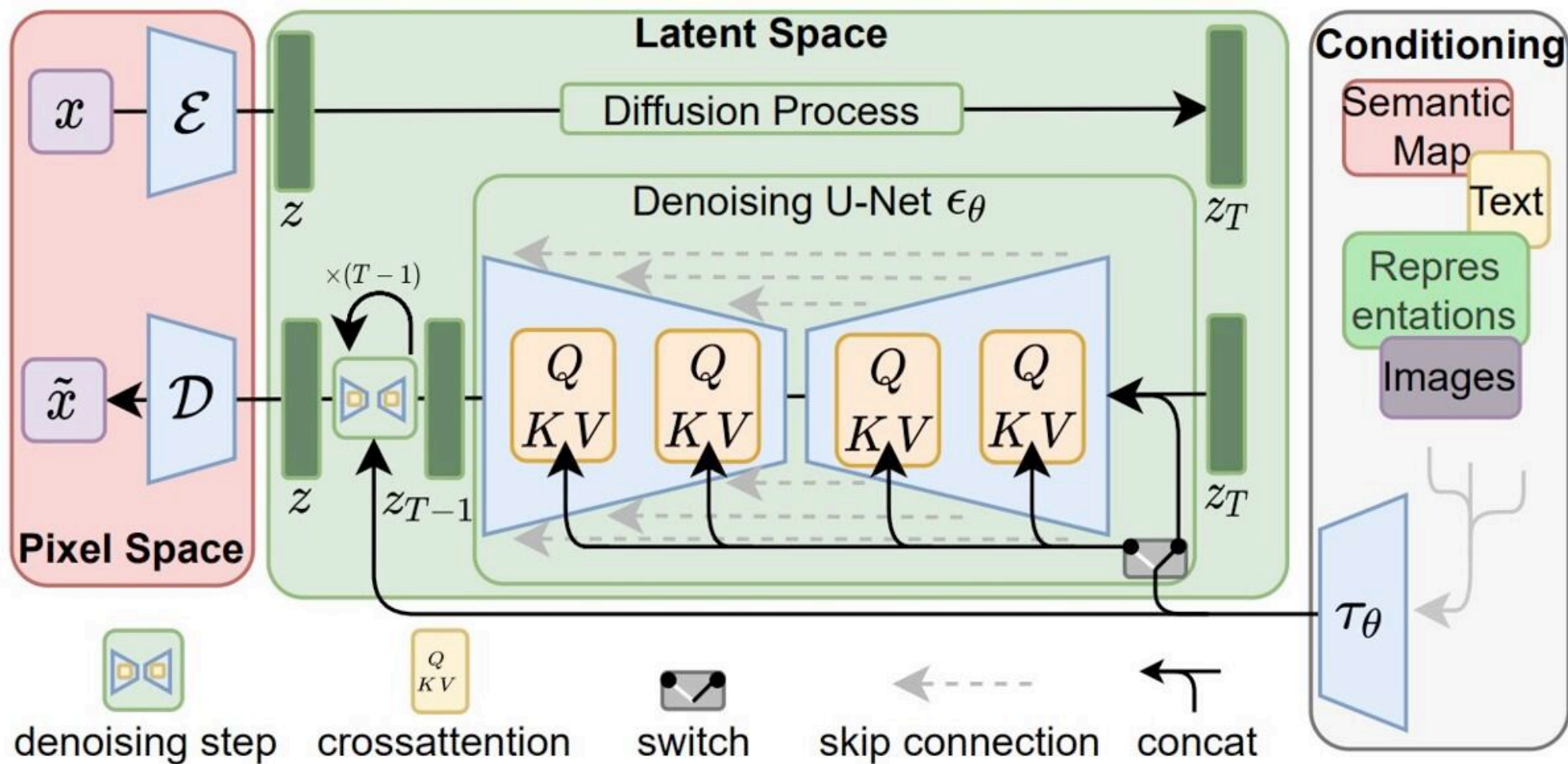
Flow Models



Stable Diffusion

(without the conditioning)



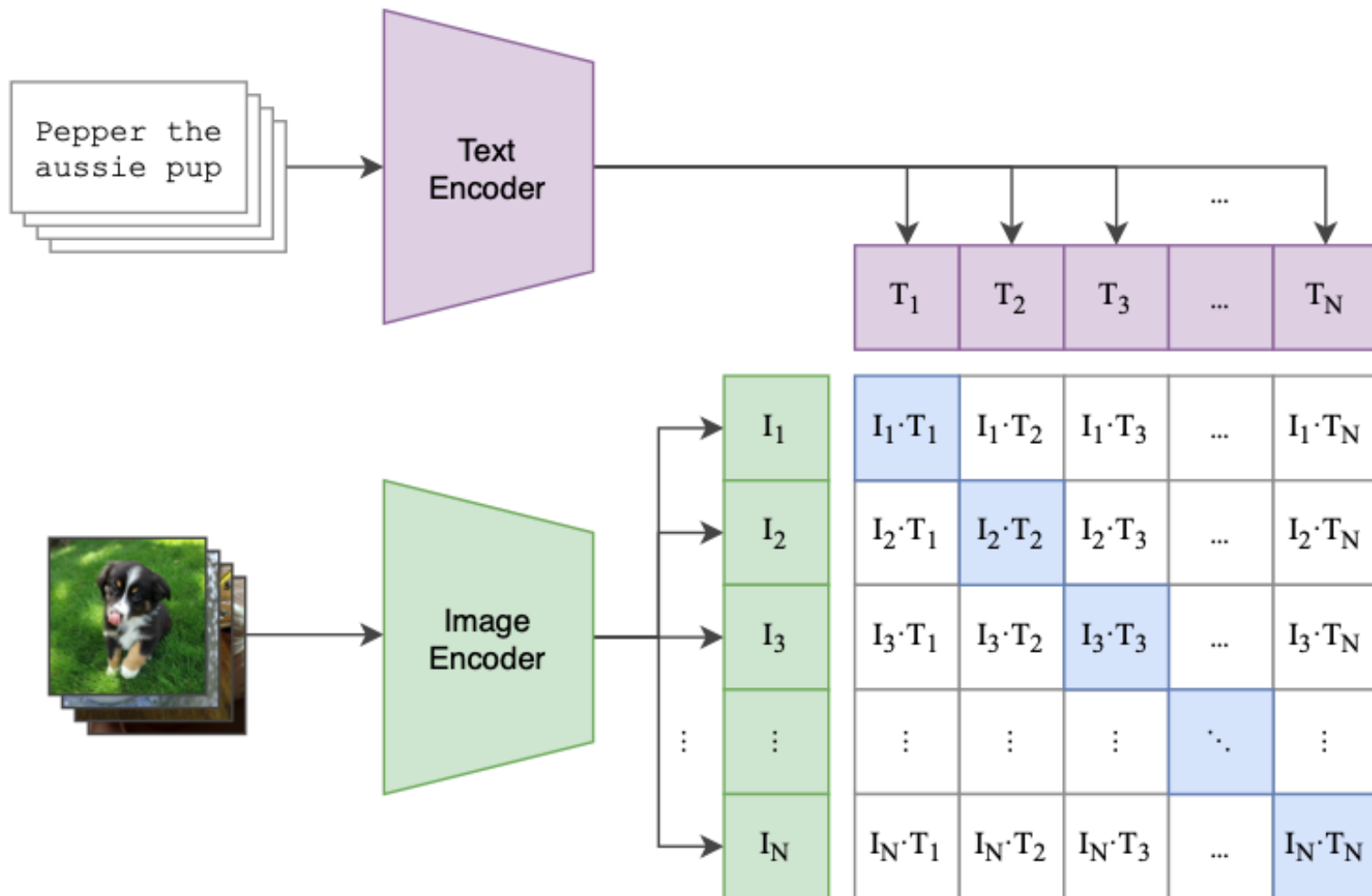




Vision and Language

Case study: CLIP

(1) Contrastive pre-training



unCLIP aka DALL-E 2

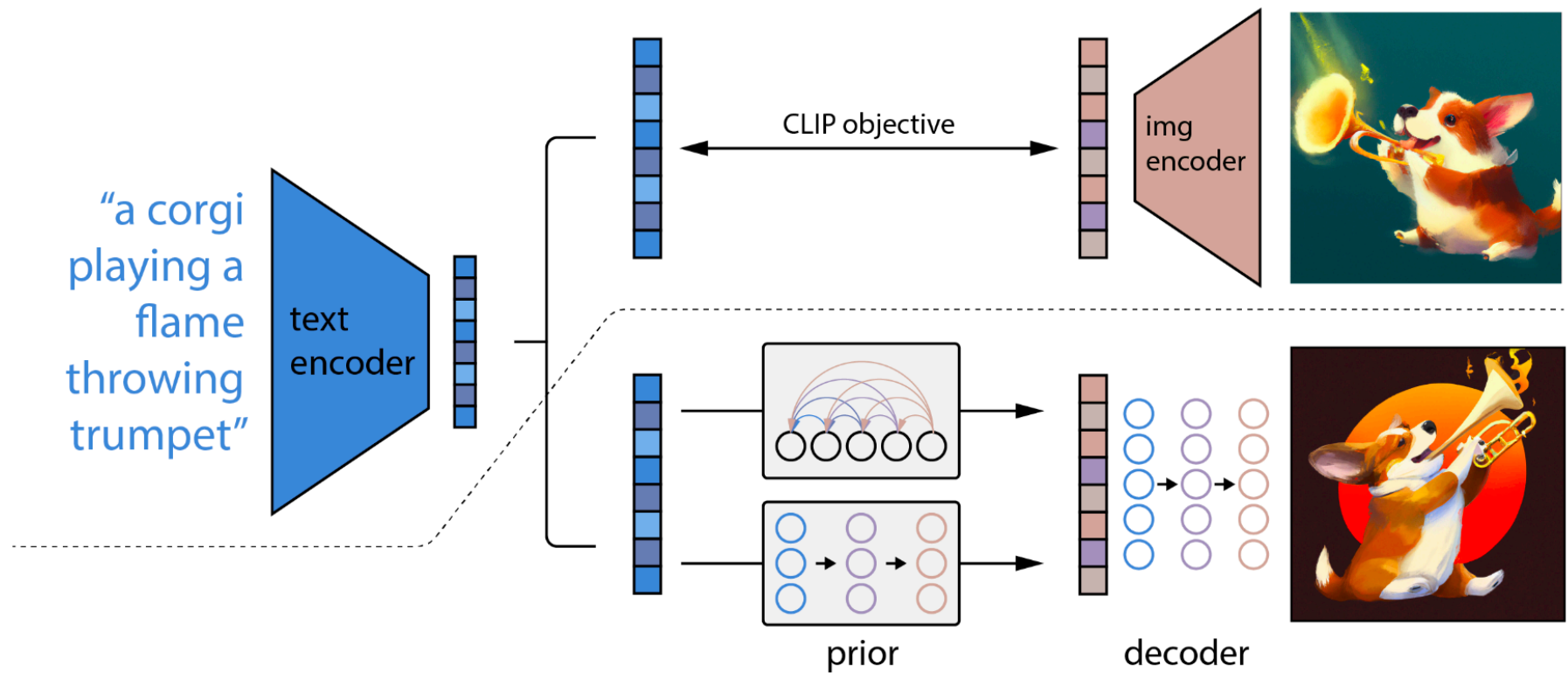


Figure 2: A high-level overview of unCLIP. Above the dotted line, we depict the CLIP training process, through which we learn a joint representation space for text and images. Below the dotted line, we depict our text-to-image generation process: a CLIP text embedding is first fed to an autoregressive 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)