

More Generative Models

CS280

Spring 2025

Angjoo Kanazawa

Mid-term Logistics

- March 19th Wed before spring break
 - Next Wednesday
- Written exam
- One page (8.5 x 11 in/A4) cheatsheet allowed (both sides)

Final project logistics

- Group of 3 is encouraged. Maximum 4, but more people = higher expectations.

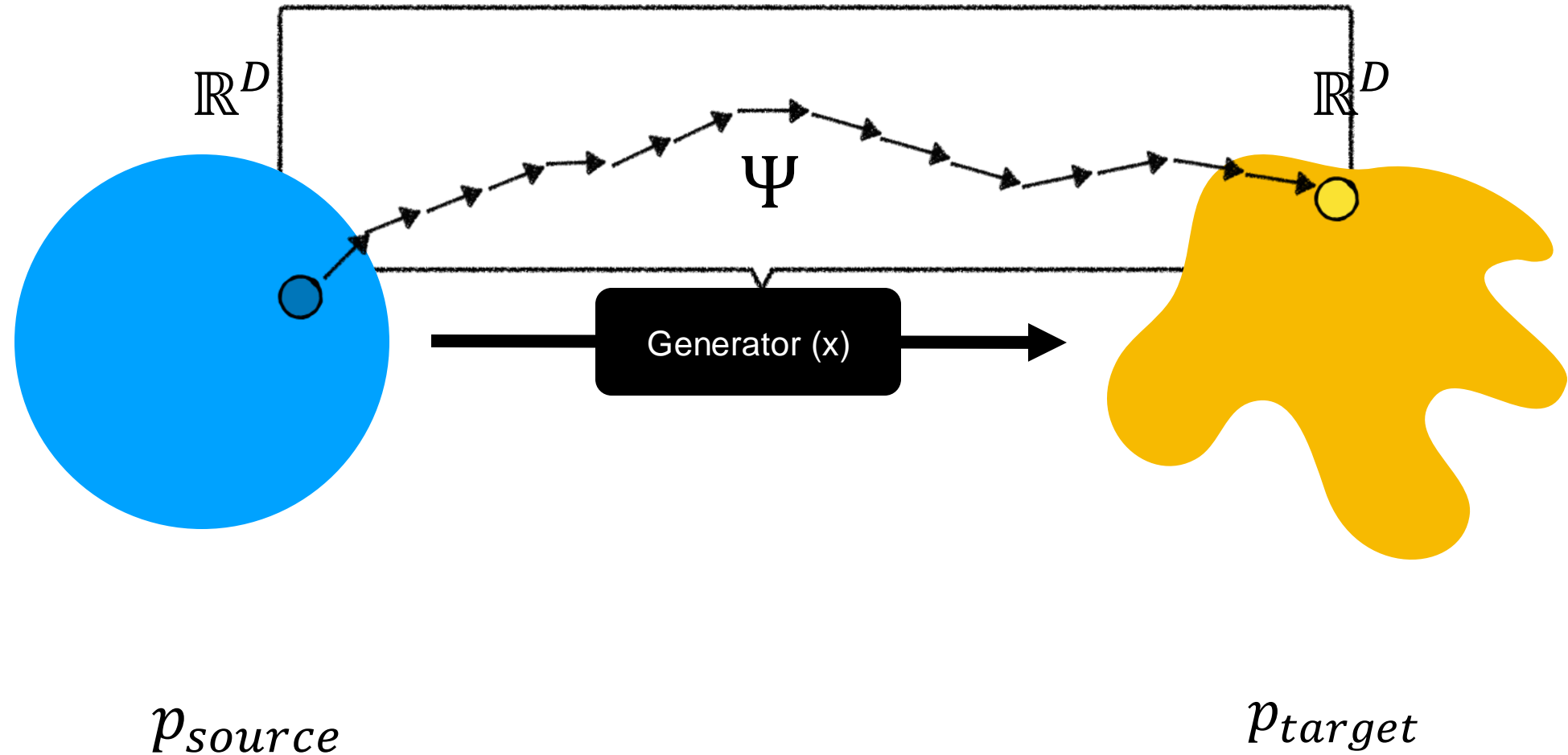
Deliverables:

- 1 page proposal due a week+ after Spring break
- Final project presentation during RRR week over the two dates 5/5 and 5/7
 - Presentation length: 3~5min
- Written report due 5/14

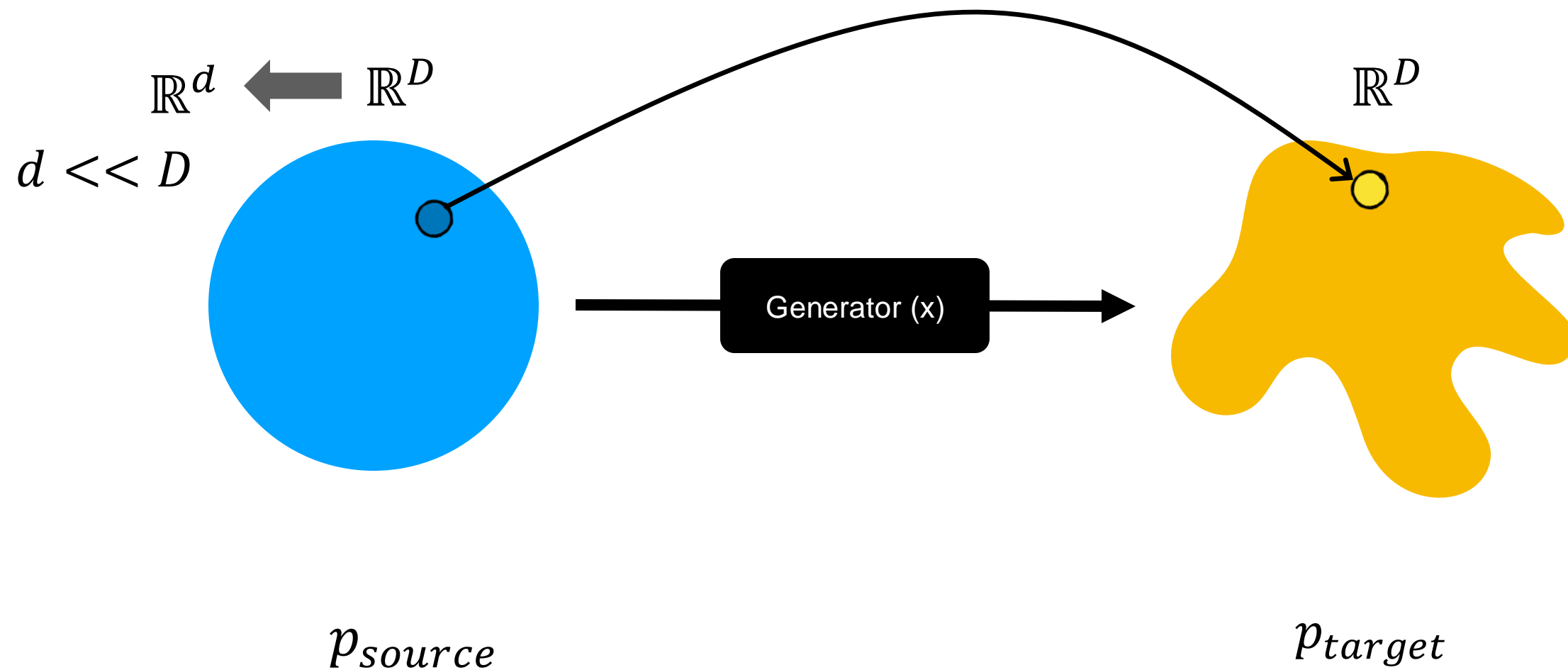
Generative Models

- Autoregressive models
- Flow based generative models
- **Variational autoencoders (VAEs)**
- **Generative adversarial networks (GANs)**

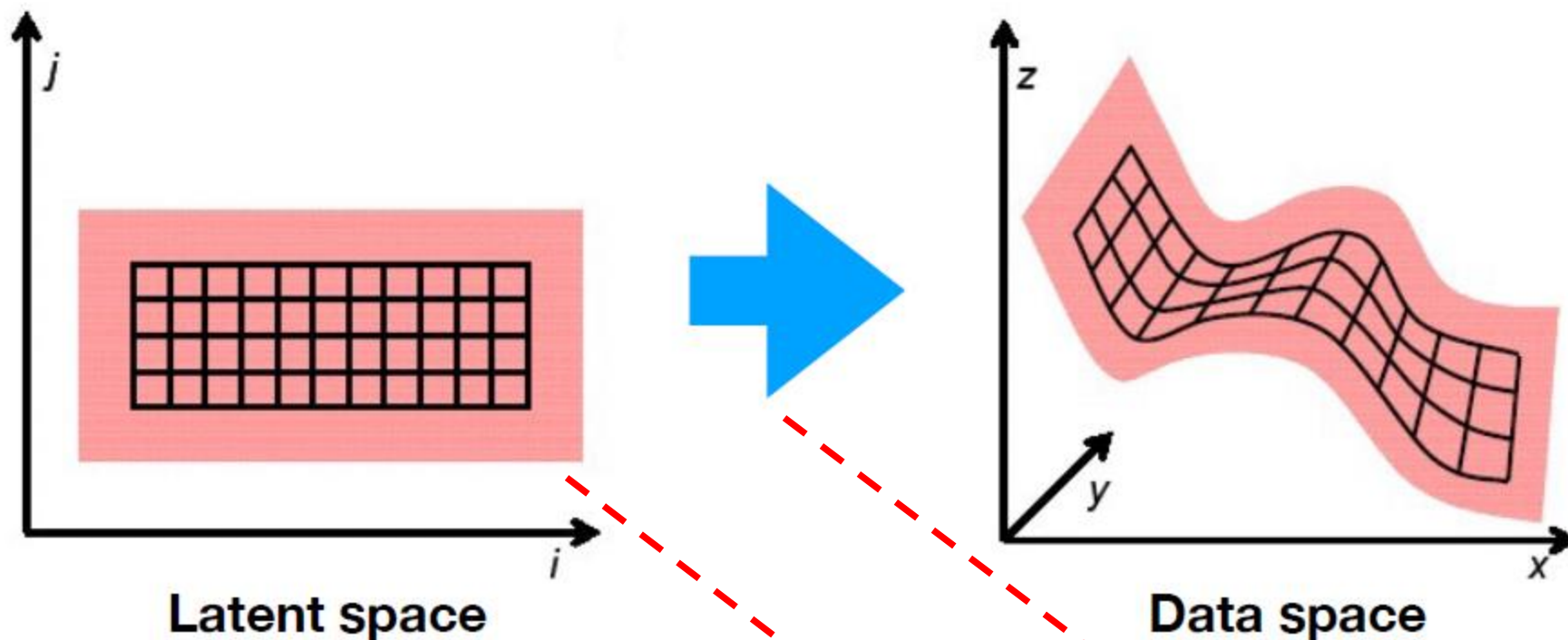
Flow based Generative Models



Latent Space Generative Models



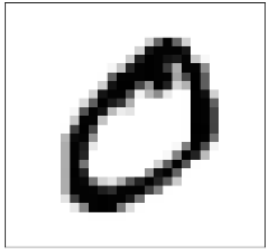
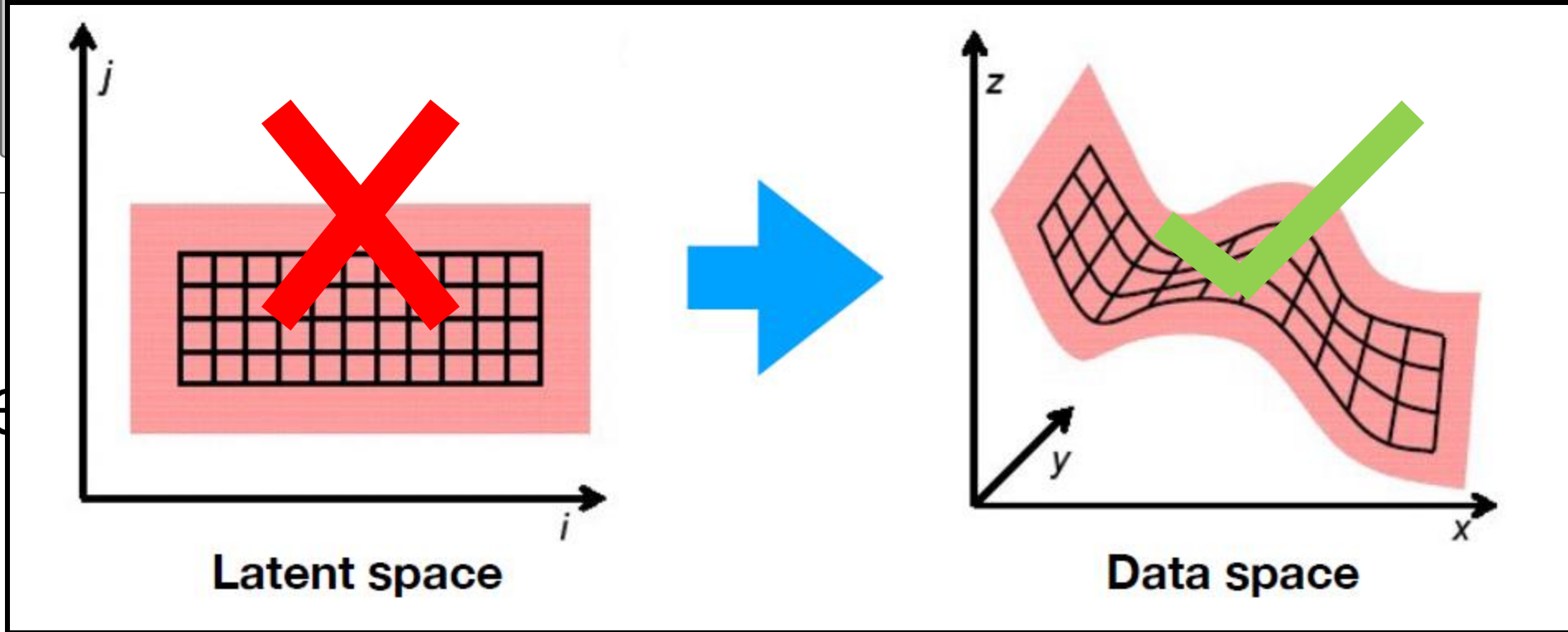
Latent space mapping approach



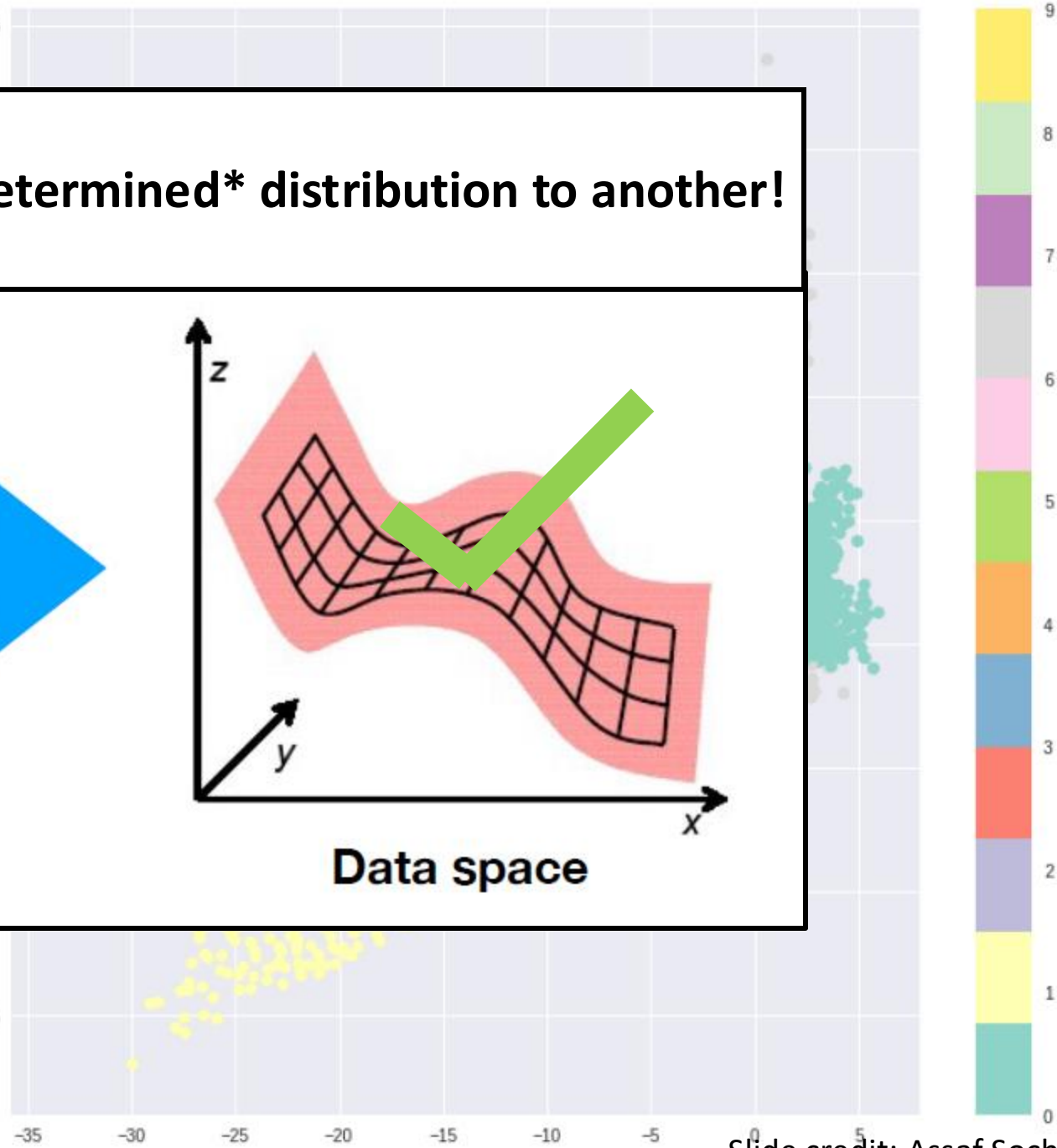
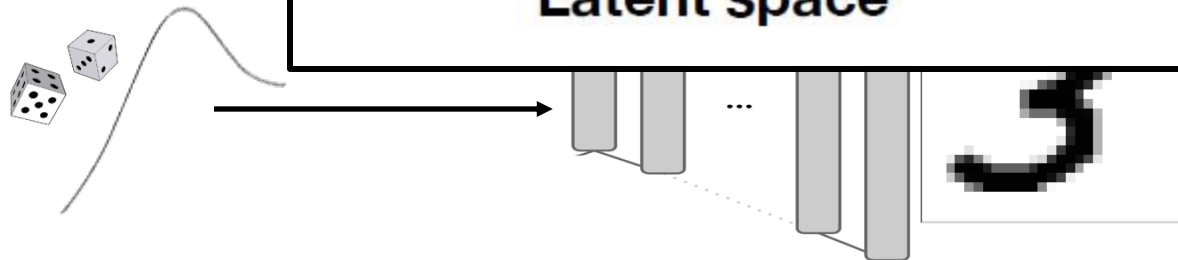
Data likelihood:
$$p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$$

Autoencoders

AE does not transform one ***pre-determined*** distribution to another!



- Gene



Variational Autoencoders

$$f(x) = z$$

- The Encoder maps x to z .
- But need to be able to sample from z ! like $p(z)$
- So we need the encoder to learn $p(z|x)$
- But how to get $p(z|x)$?
- What's wrong?
- So learn $q(z|x)$ and make it close to $p(z)$
- Through ELBO (will do later):

$$p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

$$DL(q(z|x)||p(z))$$

$$\max \log p(x) \geq \mathbb{E}_{q(z|x)} \log p(x|z) - DL(q(z|x)||p(z))$$

- i.e. max $p(x|z)$ while regularize $q(z|x)$ to be close to $p(z)$

KL divergence on $q(z|x)$ with $p(x)$

- Has a nice closed form if we use gaussians for q and p

$$q(z|x) = \mathcal{N}(\mu, \sigma^2)$$

$$p(z) = \mathcal{N}(0, I)$$

$$D_{\text{KL}}(q(z|x) \parallel p(z)) = \frac{1}{2} \sum_{i=1}^n [\sigma_i^2 + \mu_i^2 - 1 - \log \sigma_i^2]$$

$$\text{Loss: } \mathbb{E}_{q(z|x)} \log p(x|z) - \frac{1}{2} \sum_{i=1}^n [\sigma_i^2 + \mu_i^2 - 1 - \log \sigma_i^2]$$

$$\text{Objective: } \mathbb{E}_{z \sim q(z|x)} \log p(x_i|z) - DL(q(z|x_i) || p(z))$$

Another perspective

$$D_{\text{KL}}(P || Q) = \sum_{x \in \mathcal{X}} P(x) \log \left(\frac{P(x)}{Q(x)} \right).$$

$$D_{\text{KL}}(q_\phi(z|x_i) || p(z)) = \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[\log \frac{q_\phi(z|x_i)}{p(z)} \right]$$

$$= \underbrace{\mathbb{E}_{z \sim q_\phi(z|x_i)} [\log q_\phi(z|x_i)]}_{\text{Entropy!}} - \mathbb{E}_{z \sim q_\phi(z|x_i)} [\log p(z)]$$

$$- \mathcal{H}(q_\phi(z|x_i))$$

$$\mathcal{H}(p) = -\mathbb{E}_{x \sim p(x)} [\log p(x)] = -\int_x p(x) \log p(x) dx$$

$$- D_{\text{KL}}(q_\phi(z|x_i) || p(z)) = \mathbb{E}_{z \sim q_\phi(z|x_i)} [\log p(z)] + \mathcal{H}(q_\phi(z|x_i))$$

$$\text{Re-written Objective: } \mathbb{E}_{z \sim q(z|x)} \log p(x_i|z) + \log p(z) + \mathcal{H}(q(z|x_i))$$

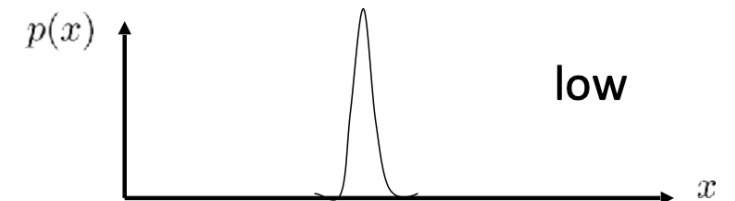
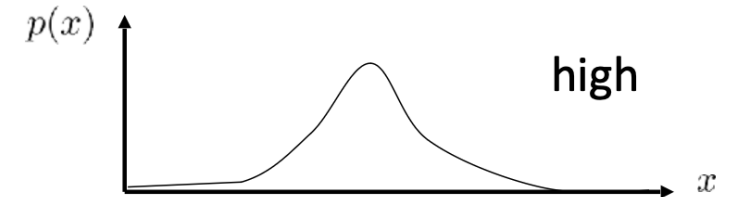
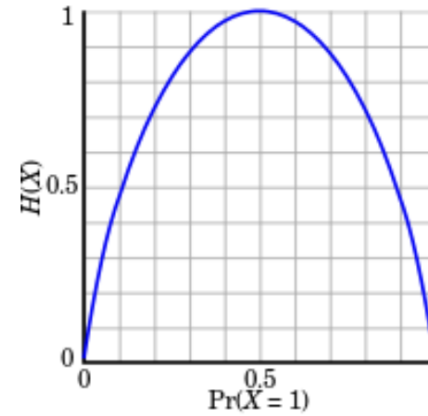
A brief aside...

Entropy:

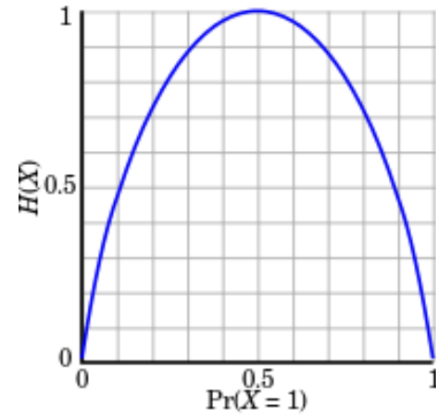
$$\mathcal{H}(p) = -E_{x \sim p(x)}[\log p(x)] = -\int_x p(x) \log p(x) dx$$

Intuition 1: how *random* is the random variable?

Intuition 2: how large is the log probability in expectation *under itself*



A brief aside...



Entropy:

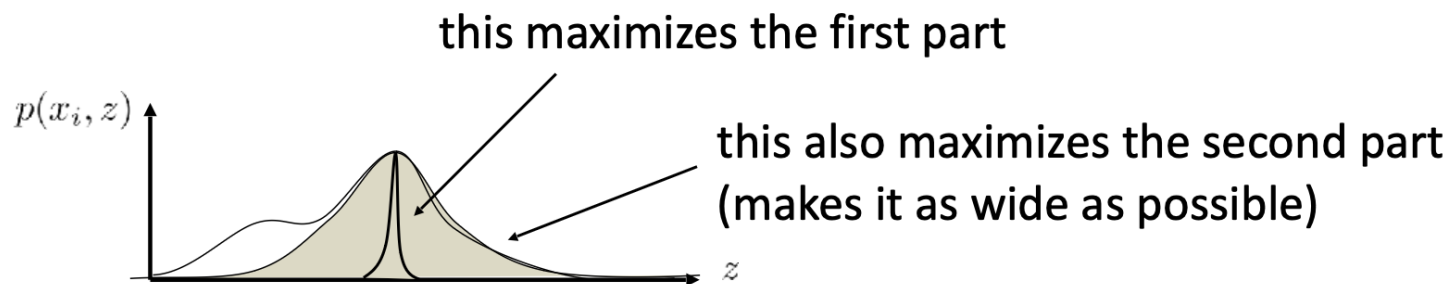
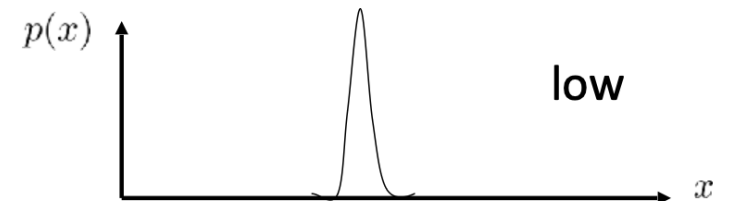
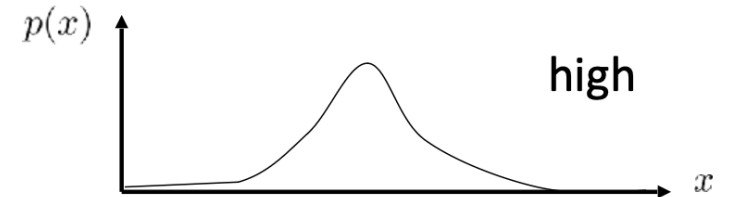
$$\mathcal{H}(p) = -E_{x \sim p(x)}[\log p(x)] = -\int_x p(x) \log p(x) dx$$

Intuition 1: how *random* is the random variable?

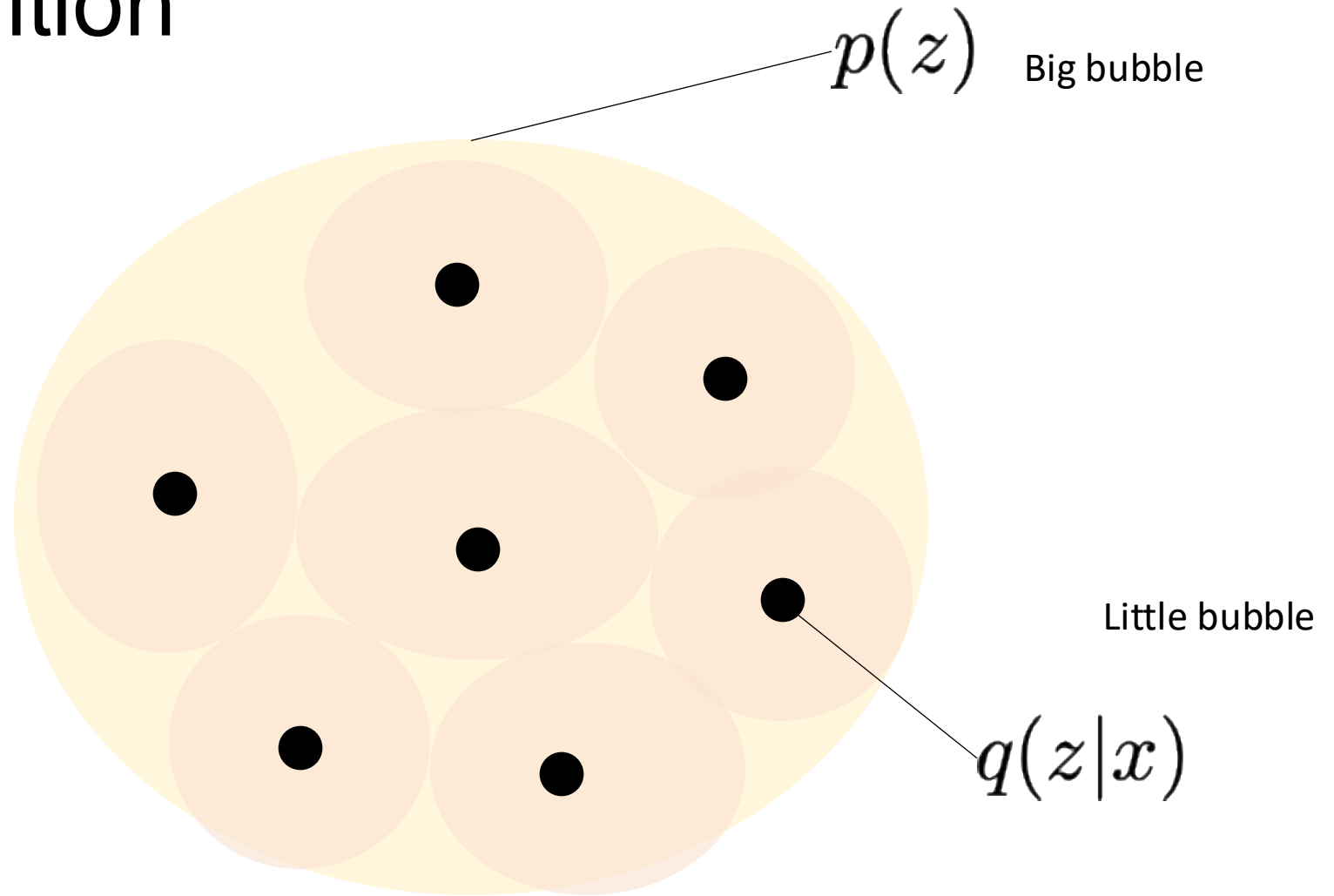
Intuition 2: how large is the log probability in expectation *under itself*

what do we expect this to do?

$$E_{z \sim q_i(z)}[\log p(x_i|z) + \log p(z)] + \mathcal{H}(q_i)$$

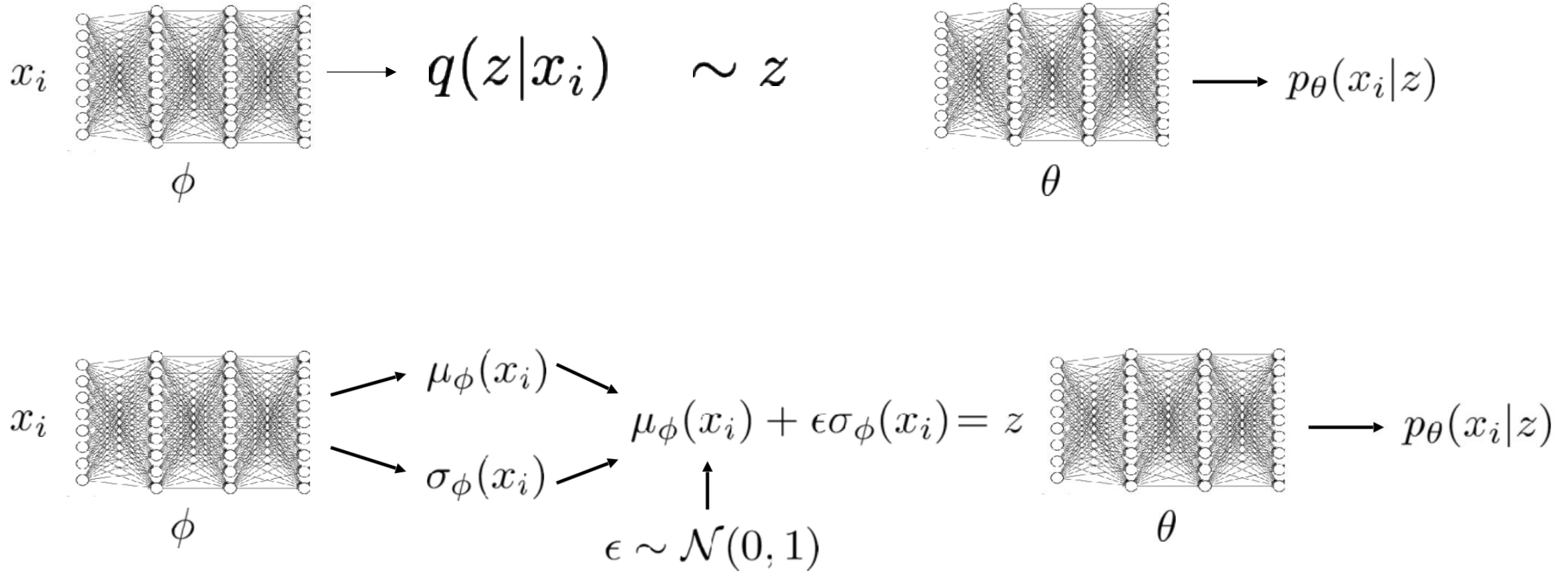


Another intuition

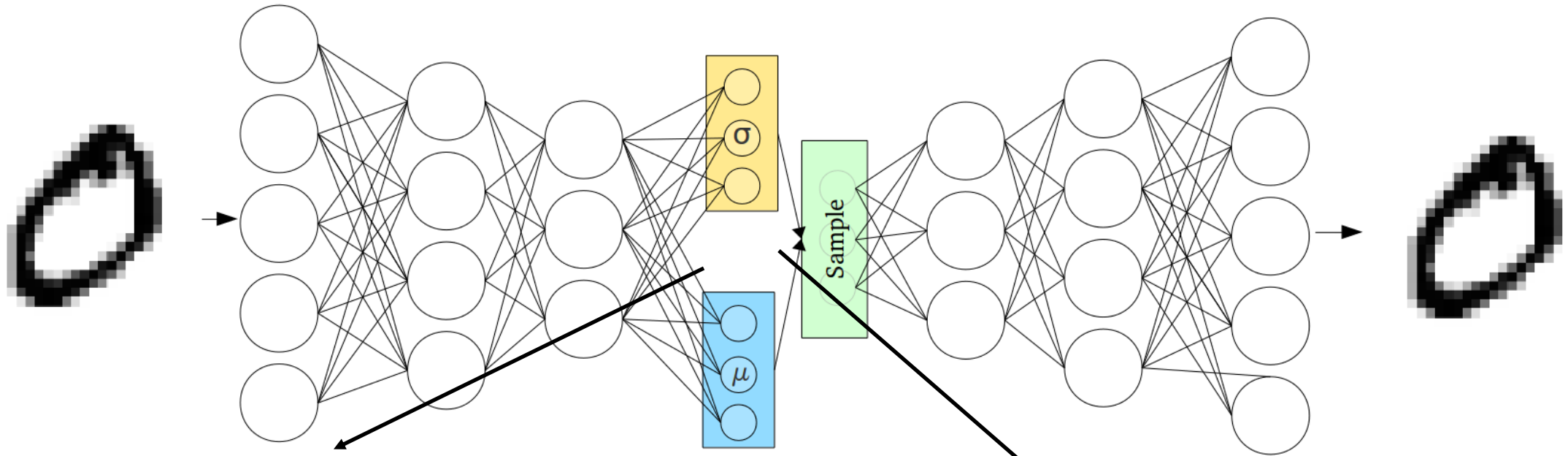


Nice conversation with Yan LeCun

How do you train through sampling?



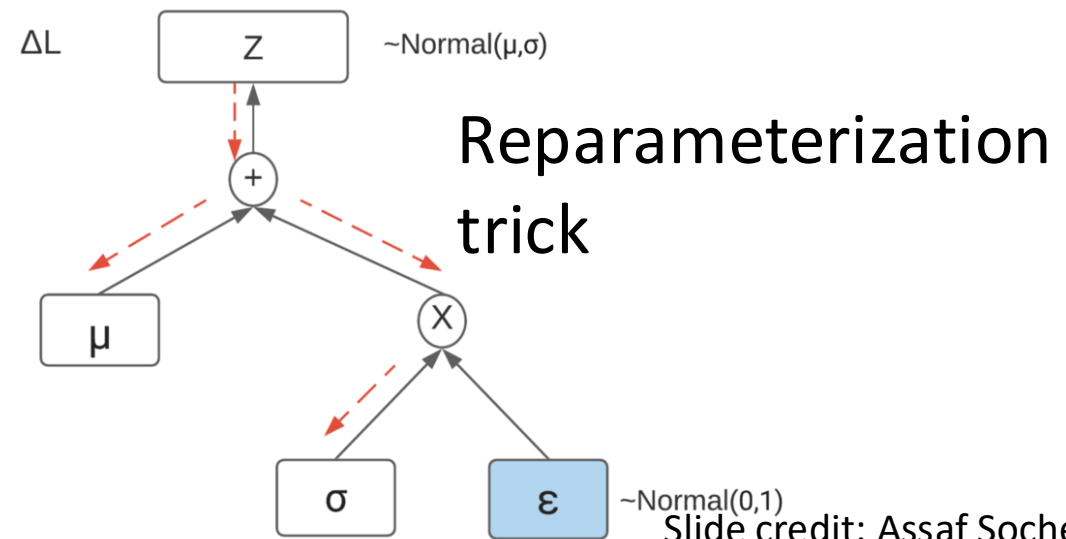
Variational Autoencoders (Kingma&Welling 2014)



Regularization:
encourage $p(z) \sim N(0,1)$

by KL divergence:

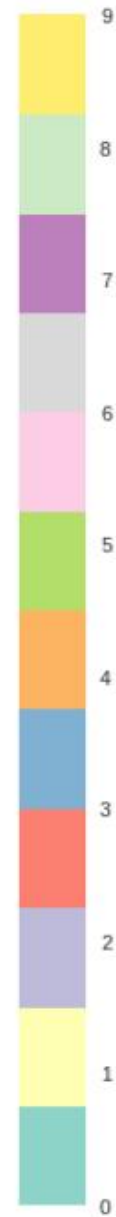
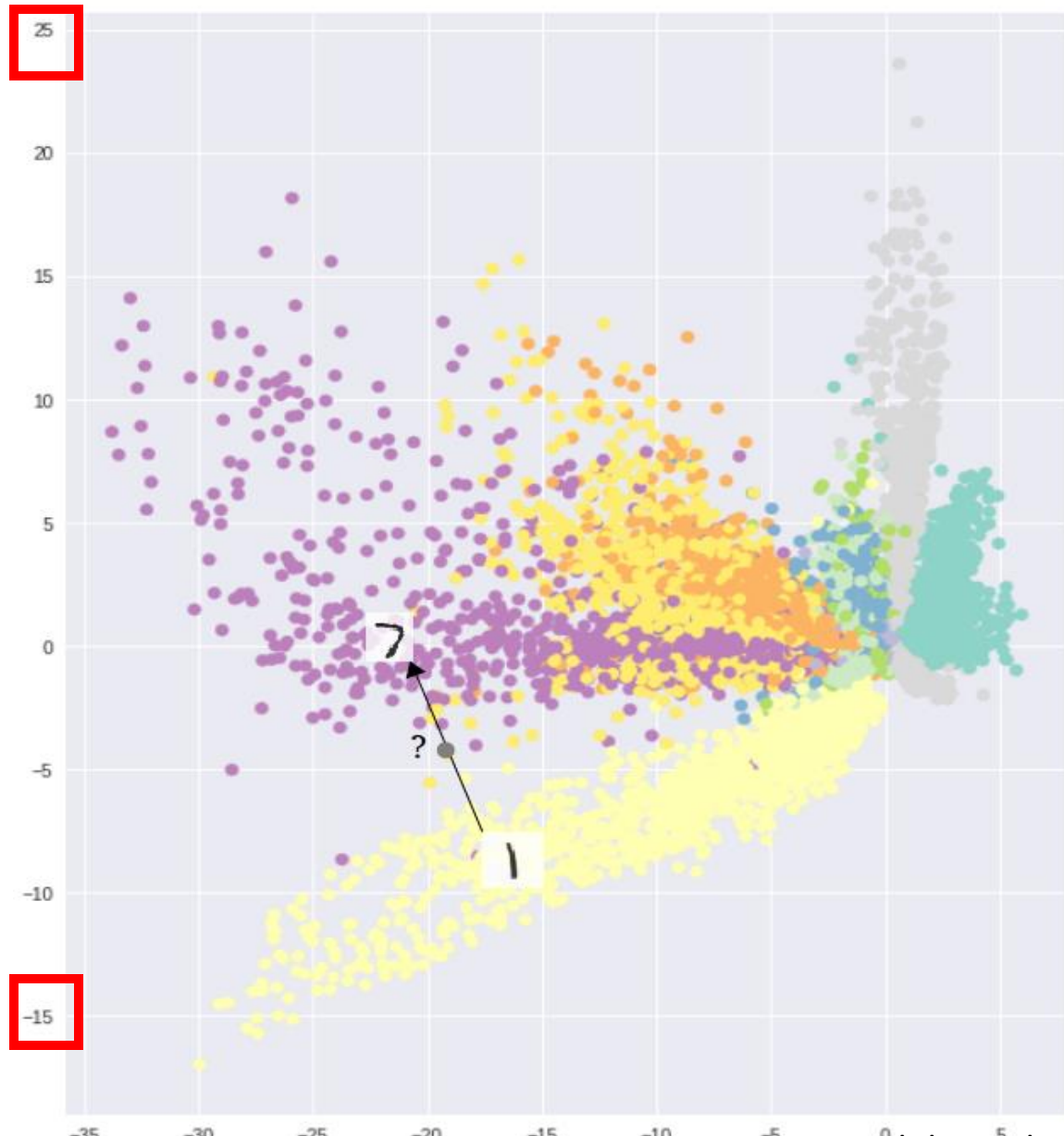
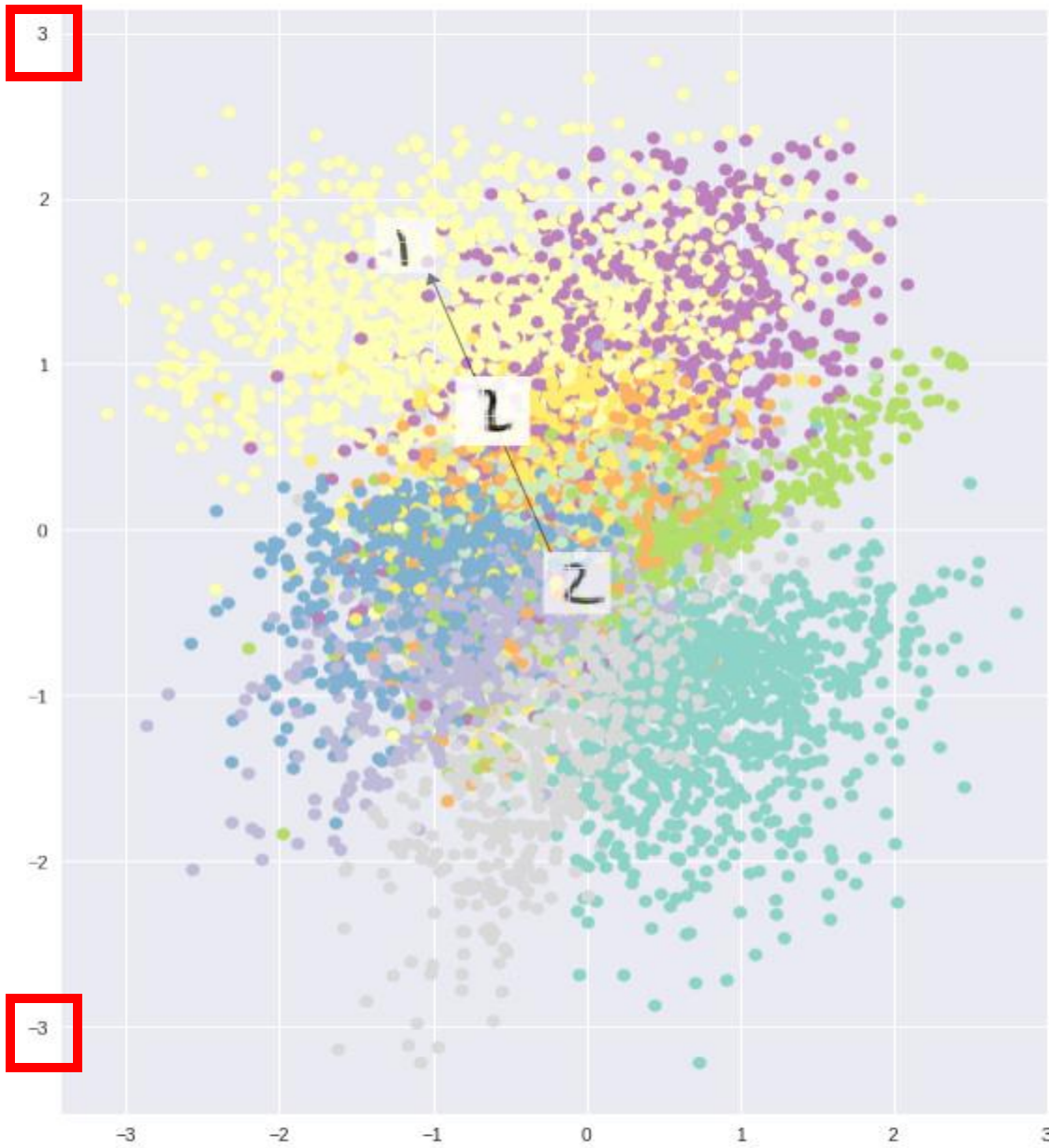
$$\sum_{i=1}^n \sigma_i^2 + \mu_i^2 - \log(\sigma_i) - 1$$



VAE

Also check out the scale!

AE



Probabilistic interpretation

Data likelihood: $p_{\theta}(x) = \int p_{\theta}(z) \overset{\mathcal{N}(0,1)}{p_{\theta}(x|z)} dz$

Goal: make $\log p_{\theta}(x^{(i)})$ as high as possible

Sample z from $z|x \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$



Encoder network
 $q_{\phi}(z|x)$
(parameters ϕ)

Sample $x|z$ from $x|z \sim \mathcal{N}(\mu_{x|z}, \Sigma_{x|z})$



Decoder network
 $p_{\theta}(x|z)$
(parameters θ)

$$\log p_{\theta}(x^{(i)}) = :$$

Slide credit: Stanford cs231n

$$\begin{aligned}
\log p_\theta(x^{(i)}) &= \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[\log p_\theta(x^{(i)}) \right] && (p_\theta(x^{(i)}) \text{ Does not depend on } z) \\
&= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z) p_\theta(z)}{p_\theta(z | x^{(i)})} \right] && (\text{Bayes' Rule}) \\
&= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z) p_\theta(z) q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)}) q_\phi(z | x^{(i)})} \right] && (\text{Multiply by constant}) \\
&= \mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \right] && (\text{Logarithms}) \\
&= \mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z)) + D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z | x^{(i)}))
\end{aligned}$$



Decoder network gives $p_\theta(x|z)$, can compute estimate of this term through sampling. (Sampling differentiable through reparam. trick, see paper.)



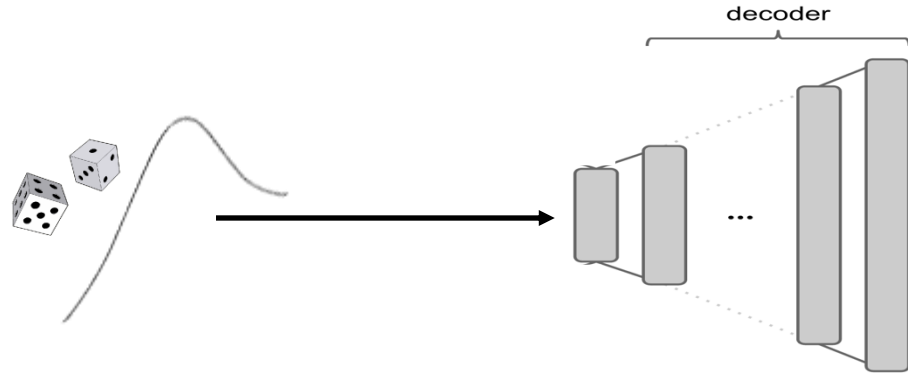
This KL term (between Gaussians for encoder and z prior) has nice closed-form solution!



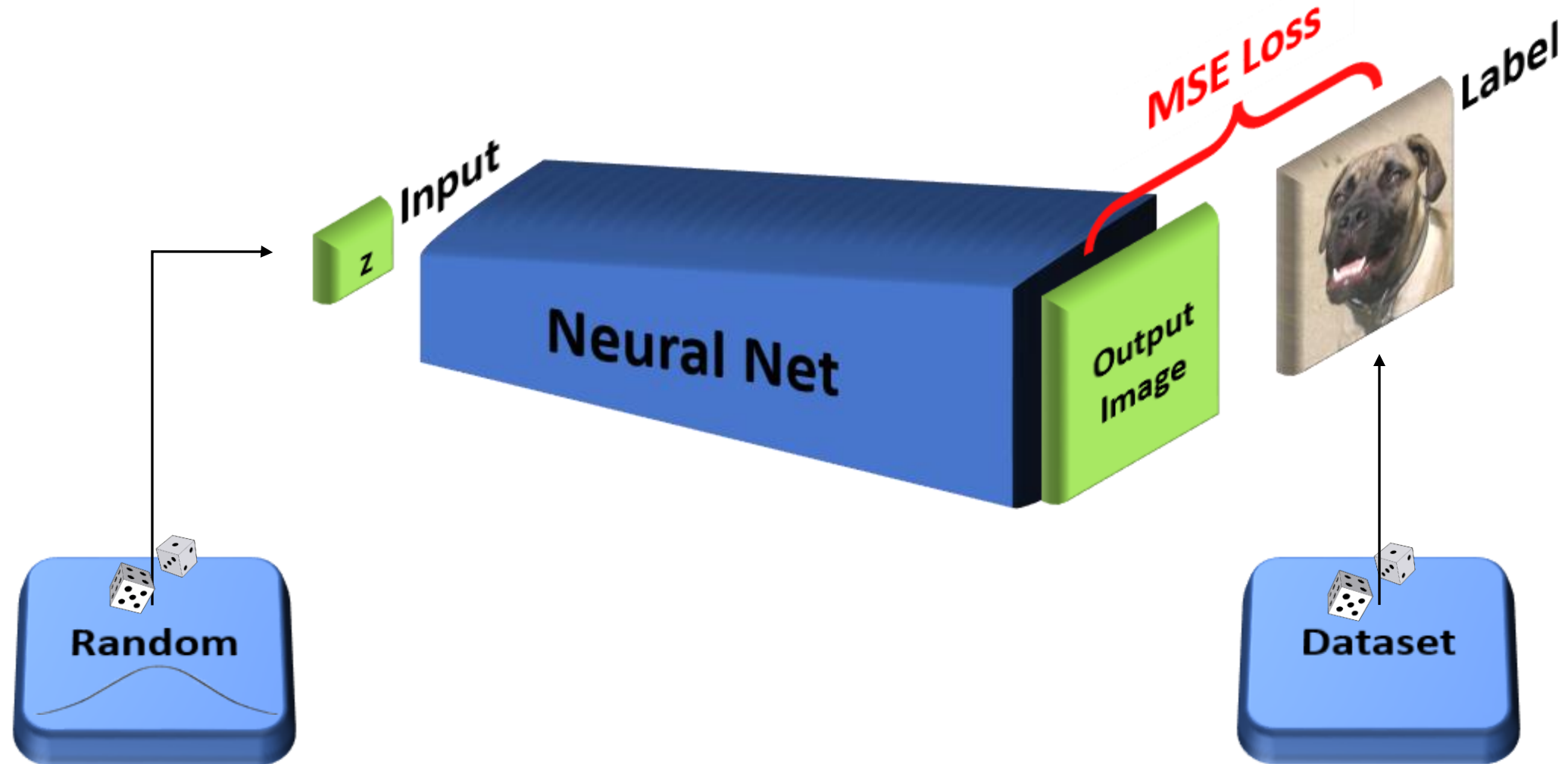
$p_\theta(z|x)$ intractable (saw earlier), can't compute this KL term :(But we know KL divergence always ≥ 0 .

Slide credit: Stanford cs231n

Generate data



How about this idea for a generative model?

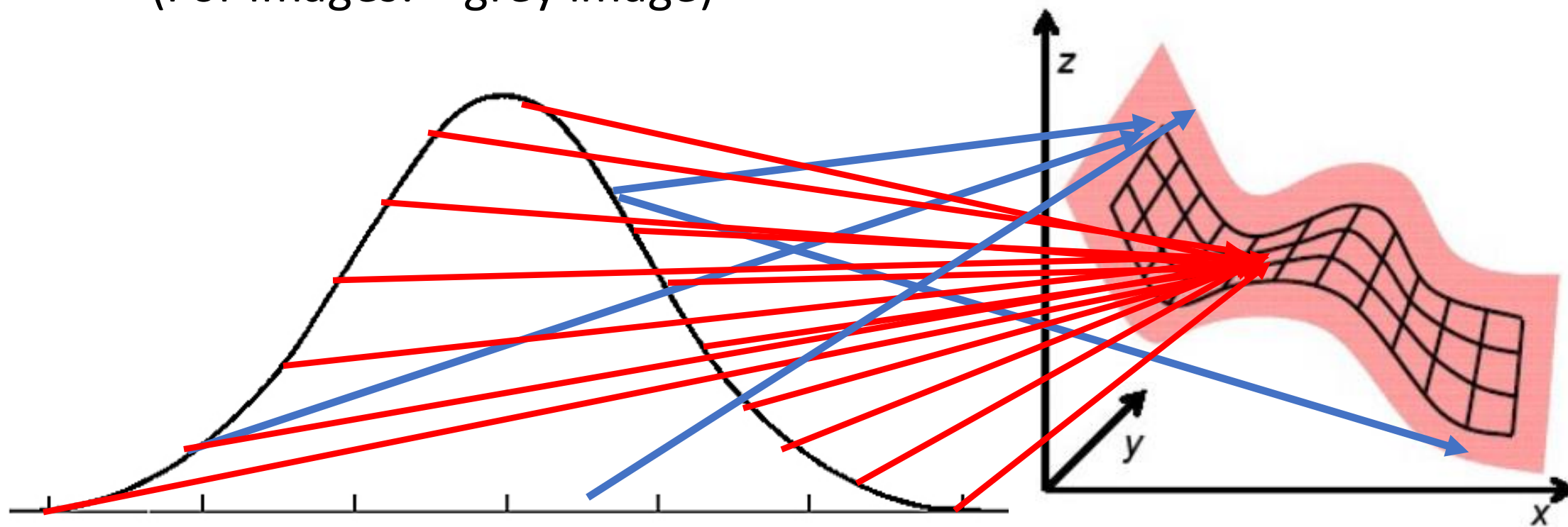


No good!

Multimodality not obtained!

In expectation: every noise is mapped to every instance

Best L2 solution: All noise is mapped to the mean
(For images: \sim grey image)



Generative Adversarial Networks



of GAN related papers per year (Salehi et al.)



Q: What makes a good counterfeiter?

Q: Who do you train first?

A: Alternate training! G,D,G,D....

Minimax game: Make the best cop do the worst mistake. Update weights Don't update weights

Maximize $\min_G \max_D \left\{ \mathbb{E}_{x \sim p_{data}} \log(D(x)) + \mathbb{E}_{z \sim p_z} \log(1 - D(G(z))) \right\}$

FAQ1: Why does it work?

- D learns probability! G trains to sample instance with high probability!
- Objective does not determine mapping directly- arrangement of latent space is learned!
- Theory: minimizes JS divergence between generated and real distributions.

FAQ2: Why alternating?

- Gradients are meaningless when game is unbalanced.
- Pre-train D? Negative examples?
- Pre-train G? What loss?
For G, D is a **learned loss function**



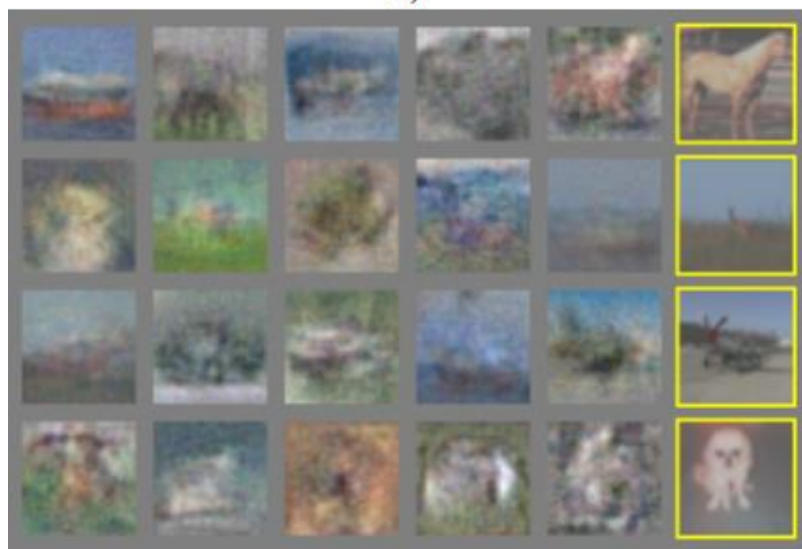
GANs, Goodfellow 2014



a)



b)



c)

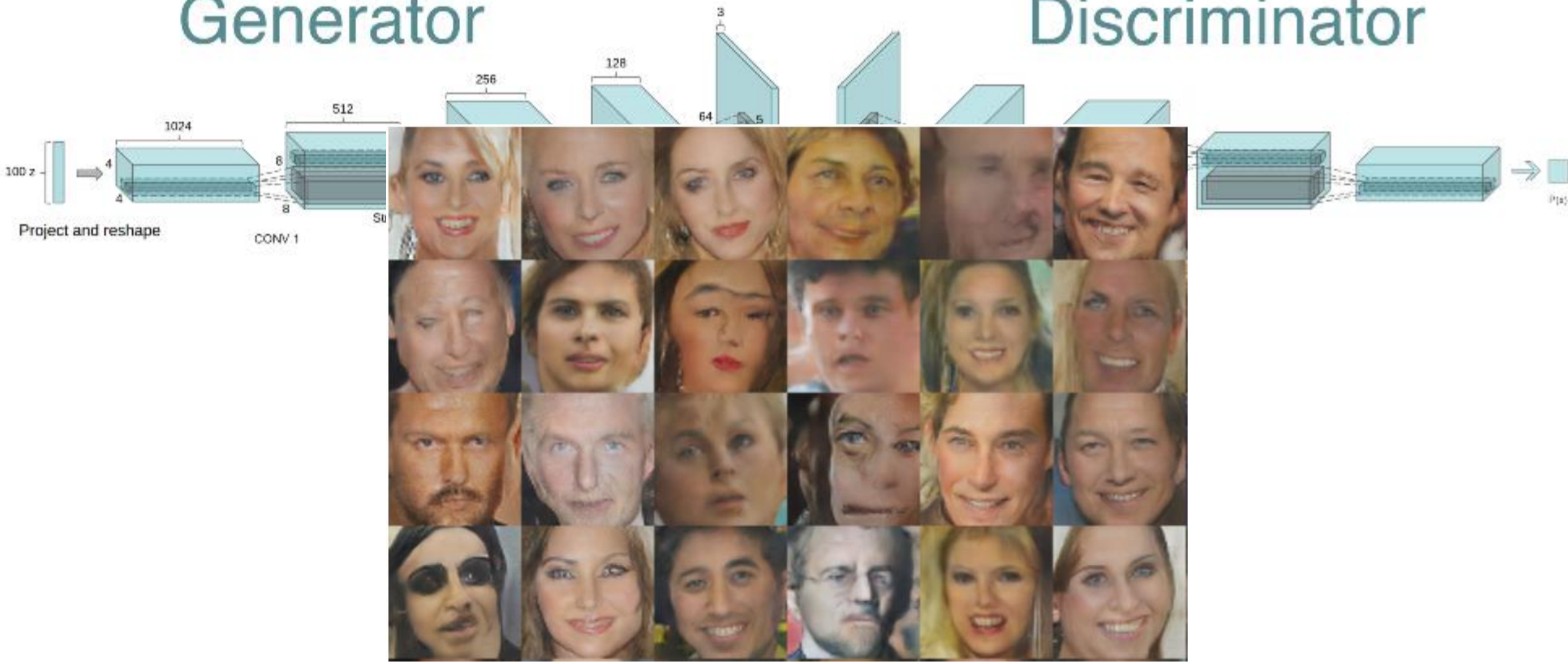


d)

DCGAN Radford 2015

Generator

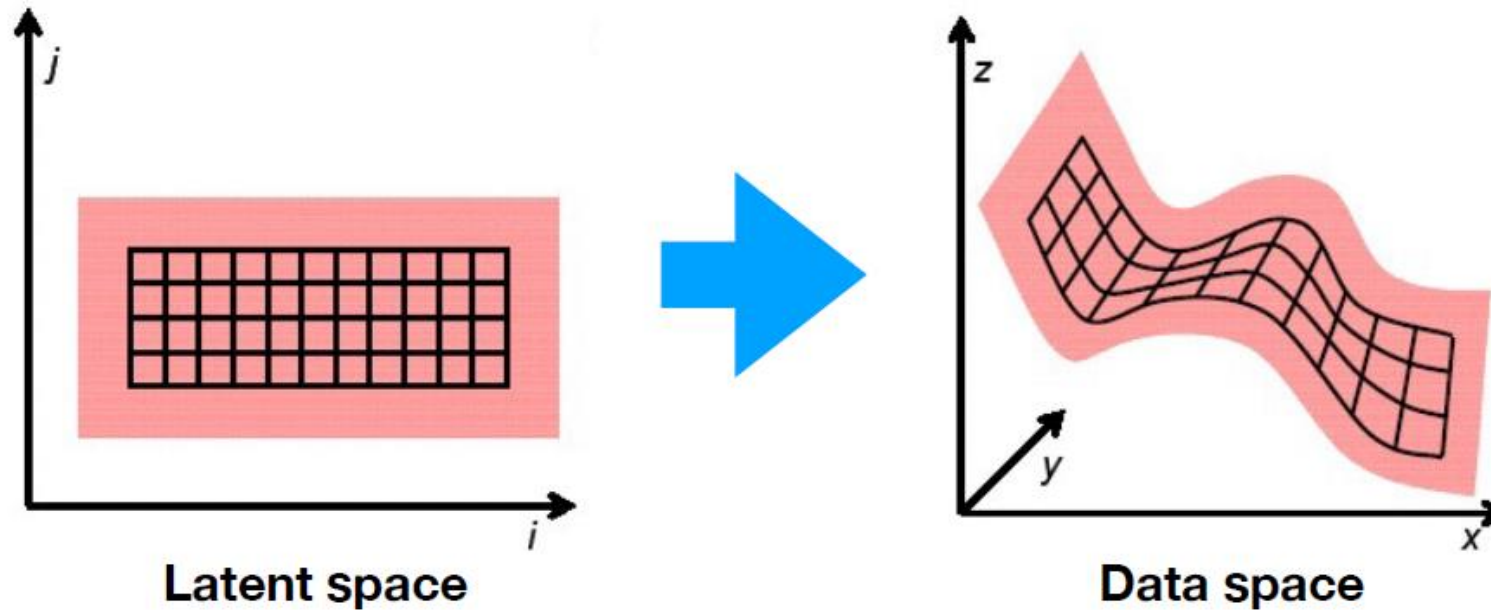
Discriminator



Latent space interpolation

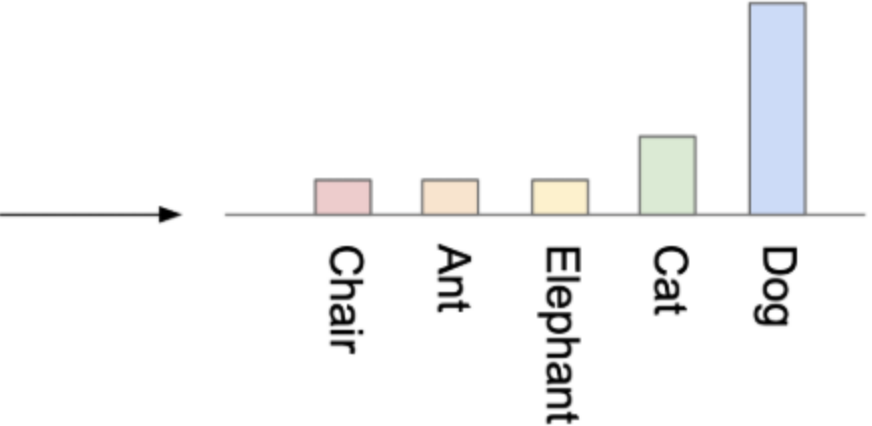


Why does it work?



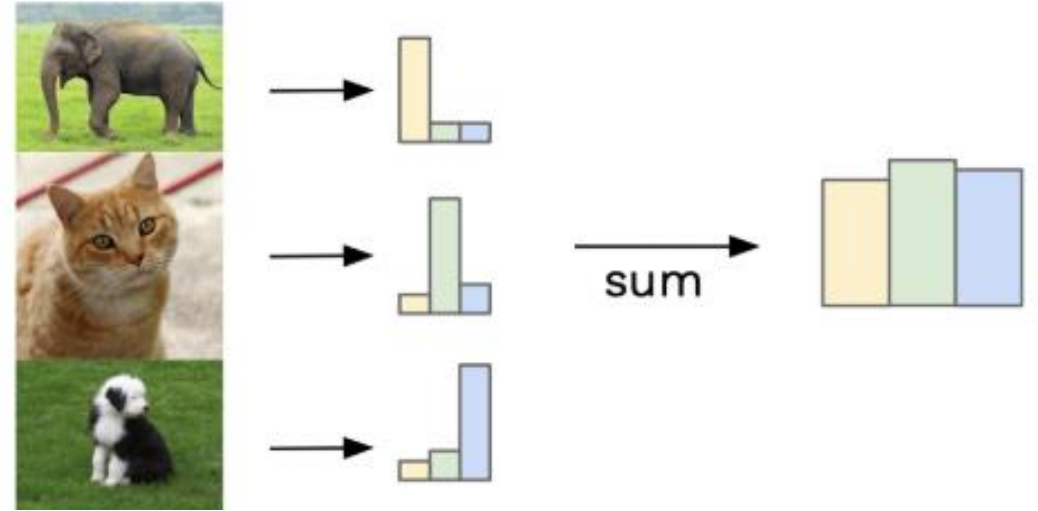
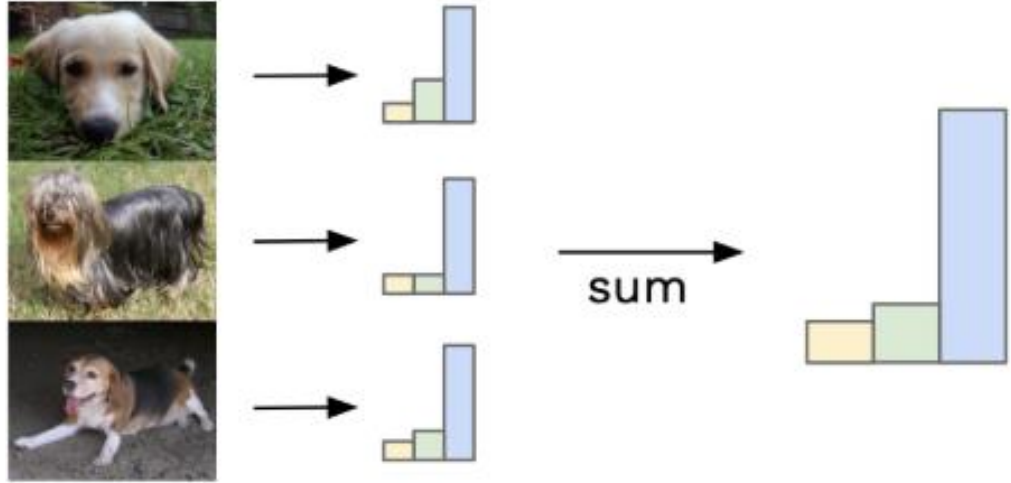
1. Every point is mapped to a valid example.
2. Network is continuous.

Evaluation metrics: Inception score



Similar labels sum to give focussed distribution

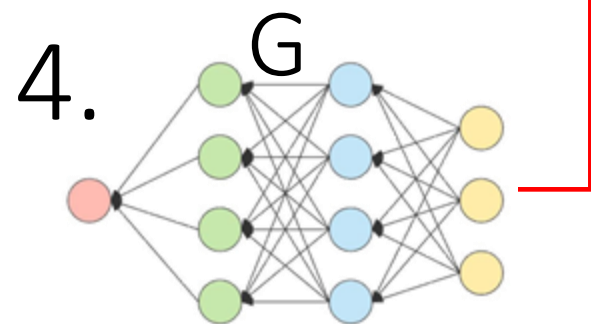
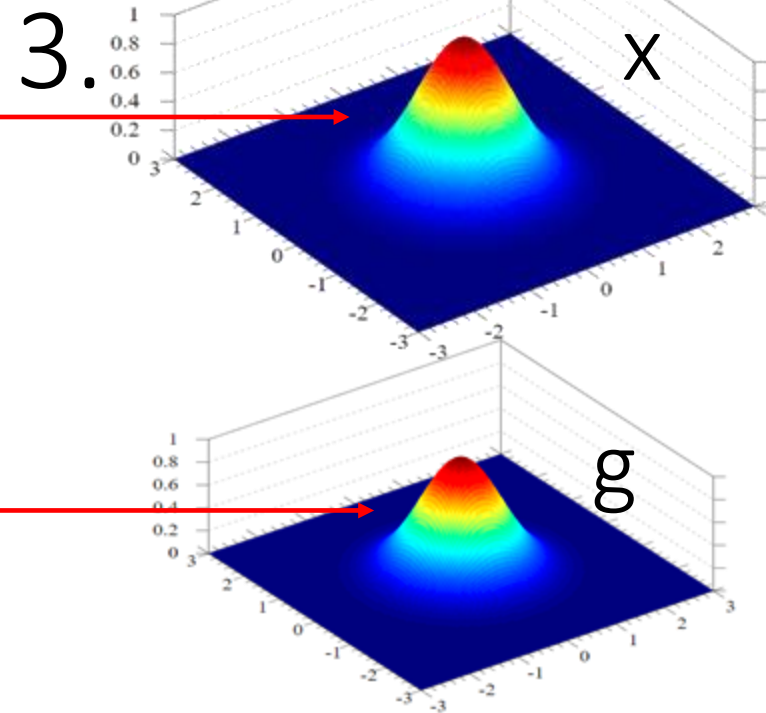
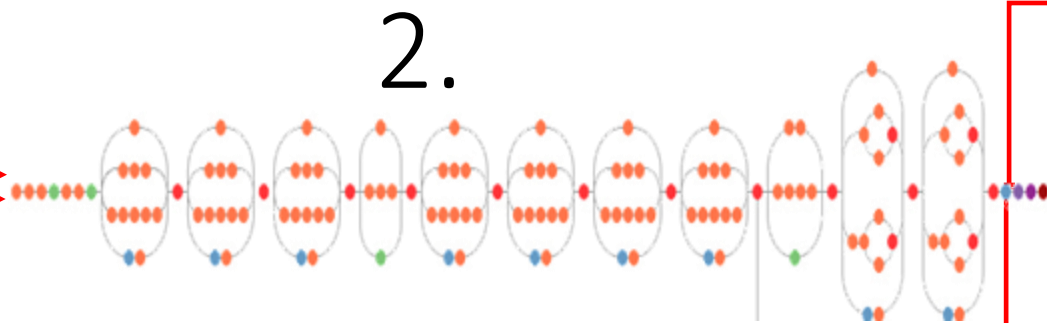
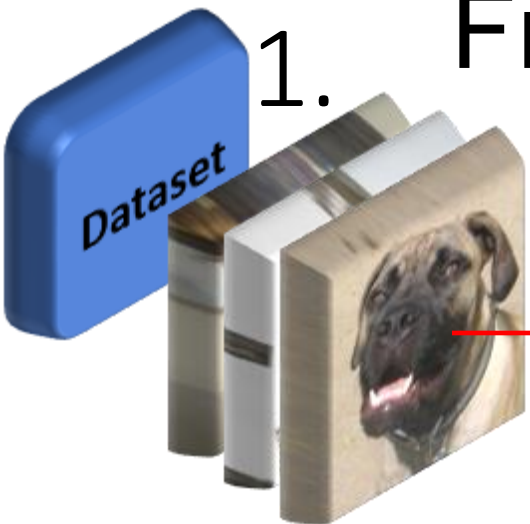
Different labels sum to give uniform distribution



$$IS(G) = \exp \left(\mathbb{E}_{\mathbf{x} \sim p_a} D_{KL} (p(y|\mathbf{x}) \parallel p(y)) \right)$$

Fréchet Inception Distance (FID)

Depends on the number of samples!



5.
$$\text{FID}(x, g) = \|\mu_x - \mu_g\|_2^2 + \text{Tr}(\Sigma_x + \Sigma_g - 2(\Sigma_x \Sigma_g)^{\frac{1}{2}})$$

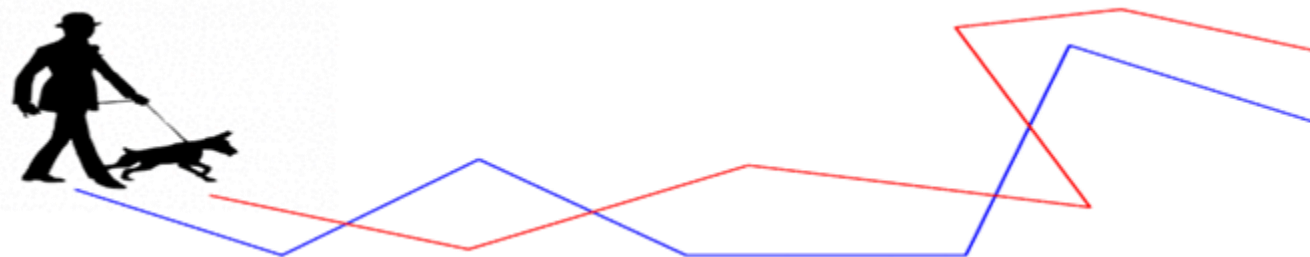
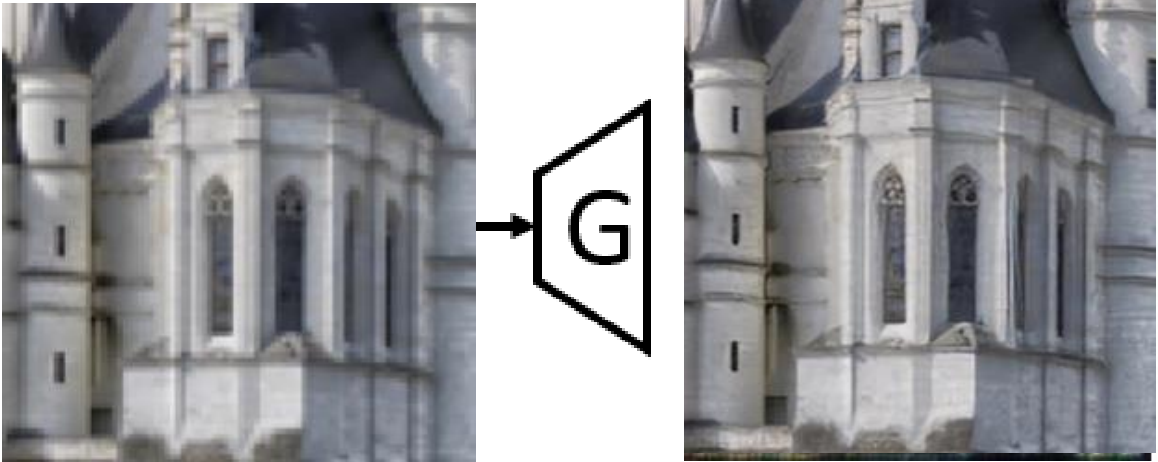
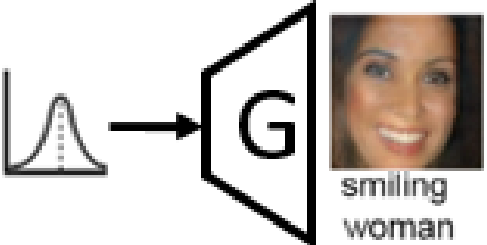
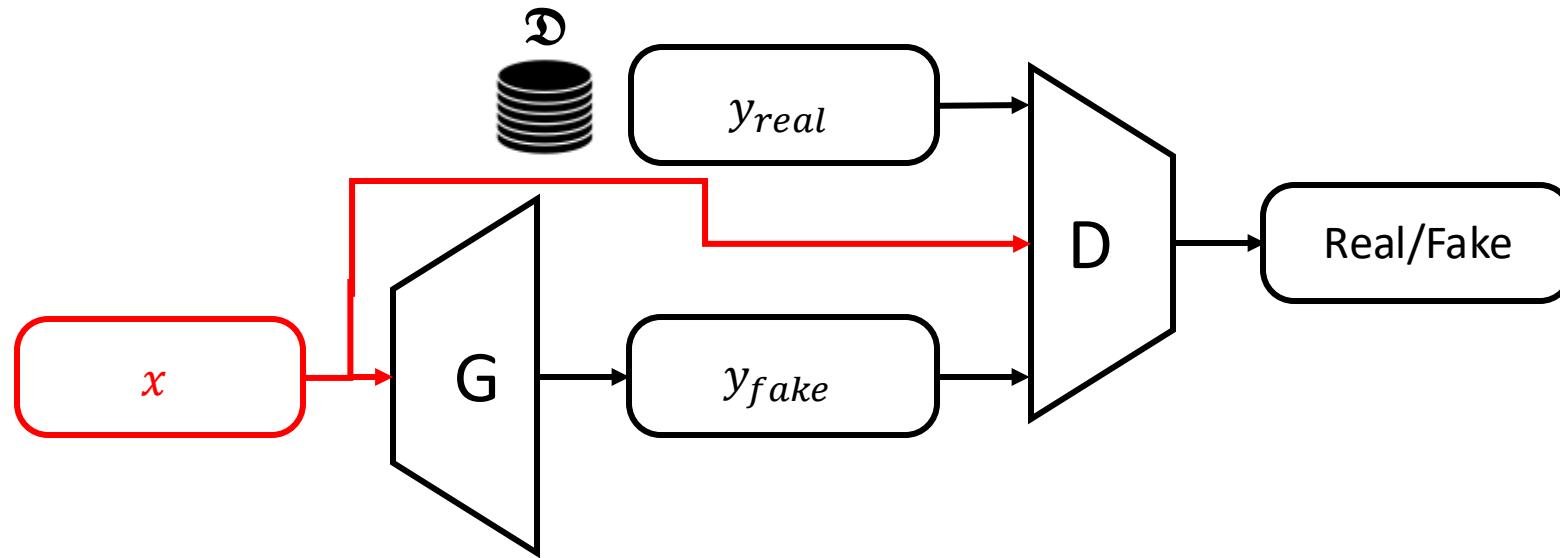


Image to Image translation

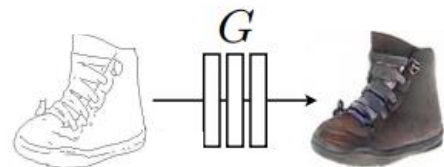


Conditional GAN



$$\mathcal{L}_{C-GAN} = \min_G \max_D \mathbb{E}[\log D(y, x)] + \mathbb{E}[\log(1 - D(G(x), x))]$$

Pix2Pix



$$\mathcal{L}_{C-GAN} = \min_G \max_D \mathbb{E}[\log D(y, \mathbf{x})] + \mathbb{E}[\log(1 - D(G(\mathbf{x}), \mathbf{x}))]$$

$$\mathcal{L}_{L1} = \|y - G(x, z)\|_1$$

$$\text{Objective} = \mathcal{L}_{C-GAN} + \lambda \cdot \mathcal{L}_{L1}$$

Training GANs is hard

- Stability



- Mode collapse

Training GANs
Be like:



Step 15k

Step 20k

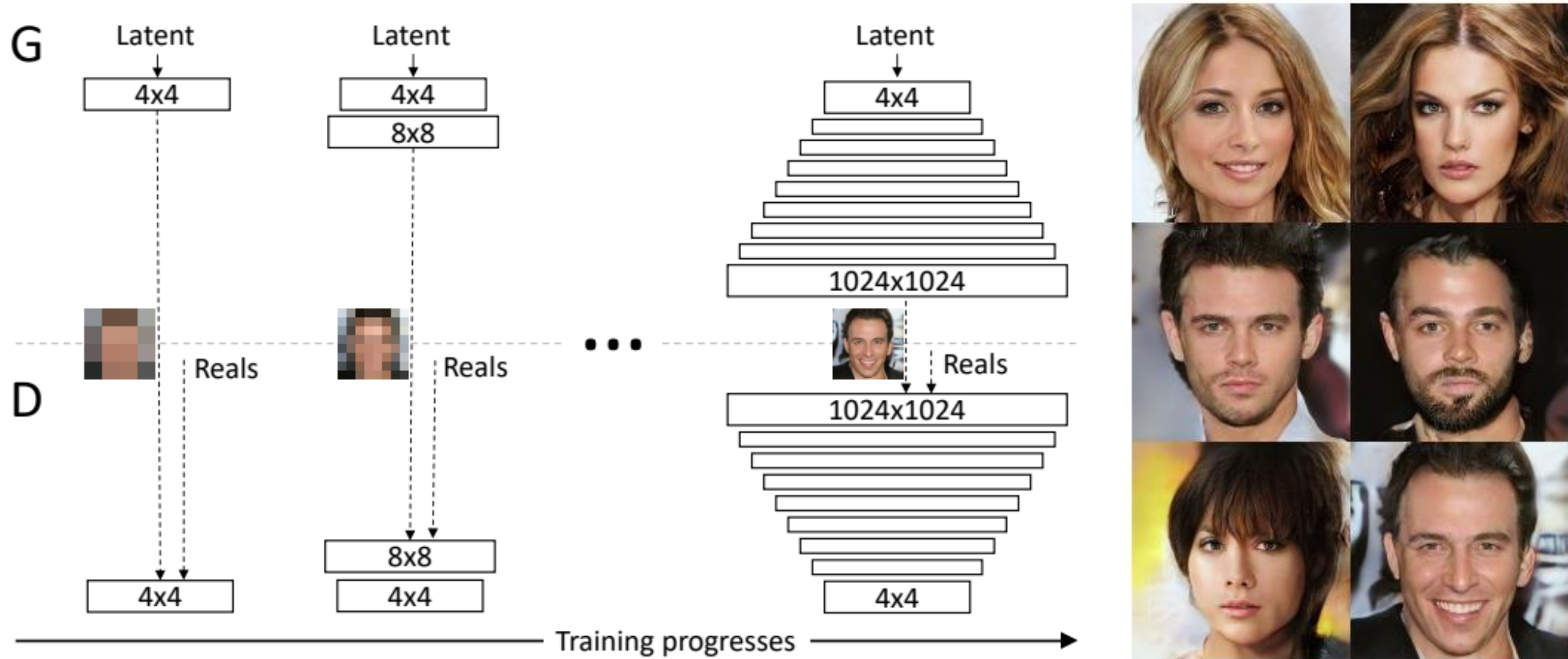
Step 25k

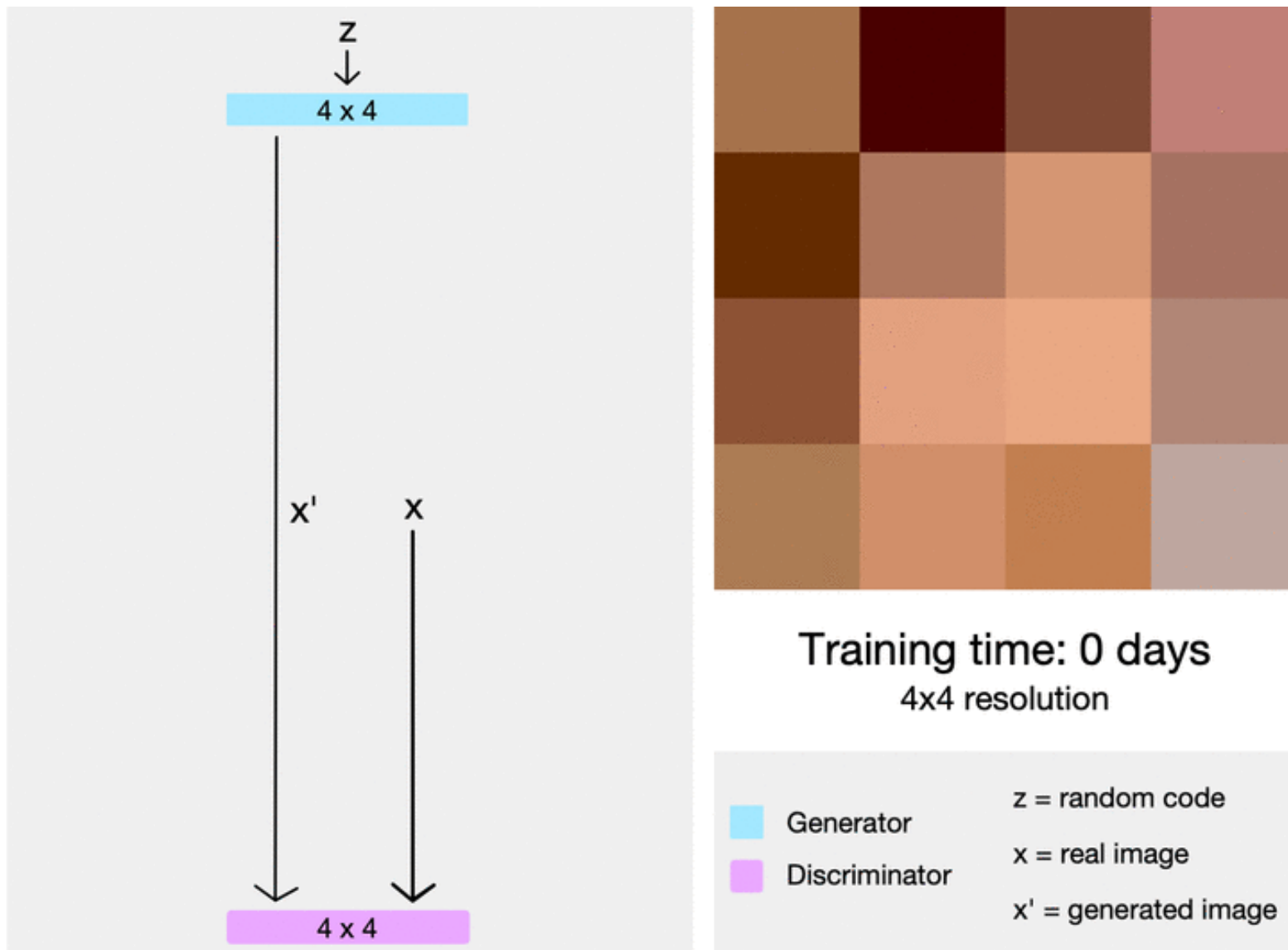
- GANs can over-train



- Loss

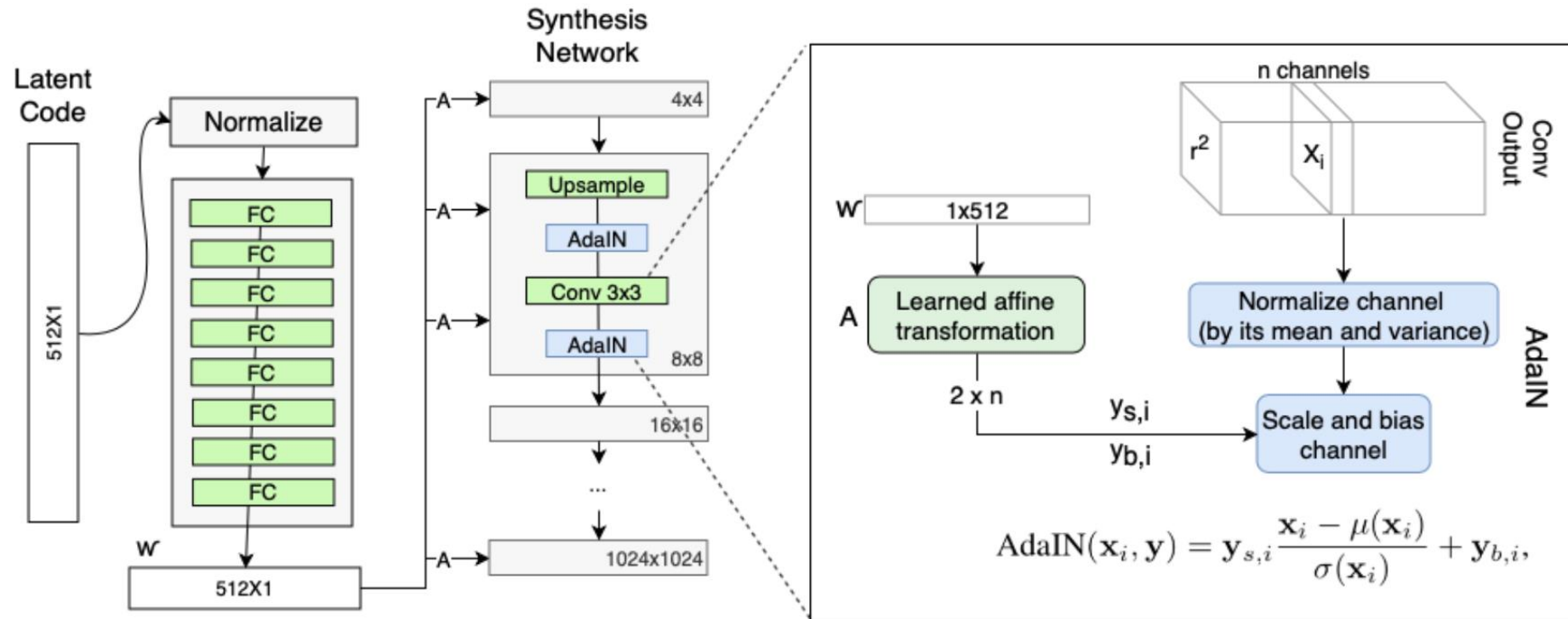
Progressive Grow





Progressive Growing of GAN, Karras et al., Feb2018

Style Modules (AdaIN)



The generator's Adaptive Instance Normalization (AdaIN)

Results

Source A: gender, age, hair length, glasses, pose



Source B:
everything
else

Result of combining A and B