

Lecture 8: Autoencoders

“Deep Learning”, Spring 2017: Lecture 8, “Autoencoders”

Learning representation – no supervision

How do we know what is relevant and what is not?

Main idea of unsupervised representation learning: a new representation should be “smaller” but allow to recover the original one.

Why smaller? (because we need to eliminate irrelevant, because we need something small)

Recap: PCA

We now want to pick the d optimal components

$$\{v_1, \dots, v_d\} \subset \mathbb{R}^N$$

Let $\tilde{x} = [\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_N]$

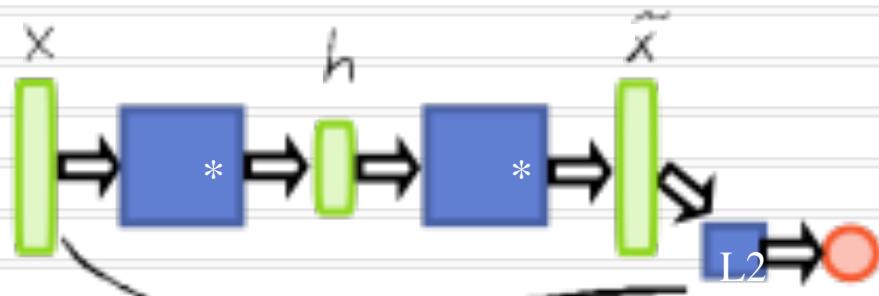
be the reconstruction of the original vectors.

We want to minimize the reconstruction error:

$$\sum_{i=1}^n \|x_i - \tilde{x}_i\|_2^2 \rightarrow \min$$

$$\|\mathbf{x} - \tilde{\mathbf{x}}\|_F^2 \rightarrow \min$$

Neural network version (autoencoder)



PCA

$$\begin{aligned} h &= u x \\ \tilde{x} &= w h \\ m \cdot n \cdot \|x - \tilde{x}\|_2^2 & \\ u & \\ s.t. \quad & \langle u, u \rangle = \delta_{ij} \end{aligned}$$

Autoencoder:

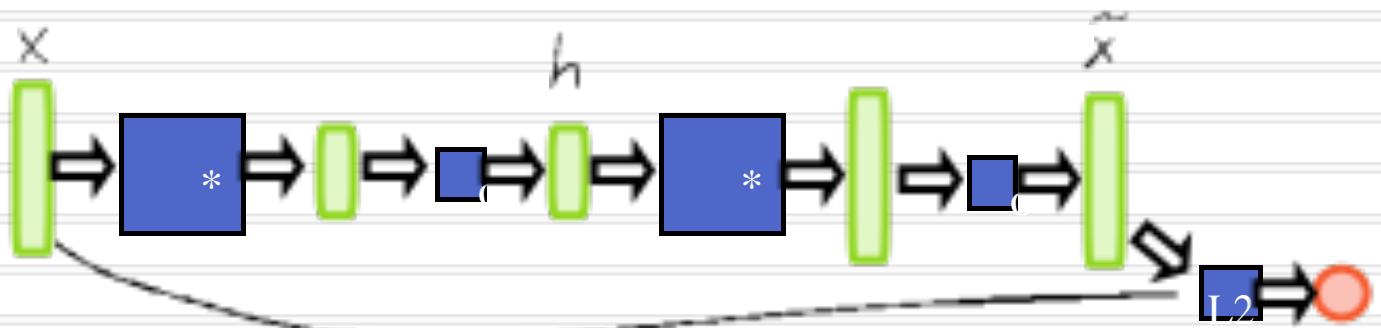
$$\begin{aligned} h &= w x \\ \tilde{x} &= w h \\ m \cdot n \cdot \|x - \tilde{x}\|_2^2 & \\ w & \end{aligned}$$

[Bengio et al. 07]

Autoencoder: adding non-linearity

encoder

decoder

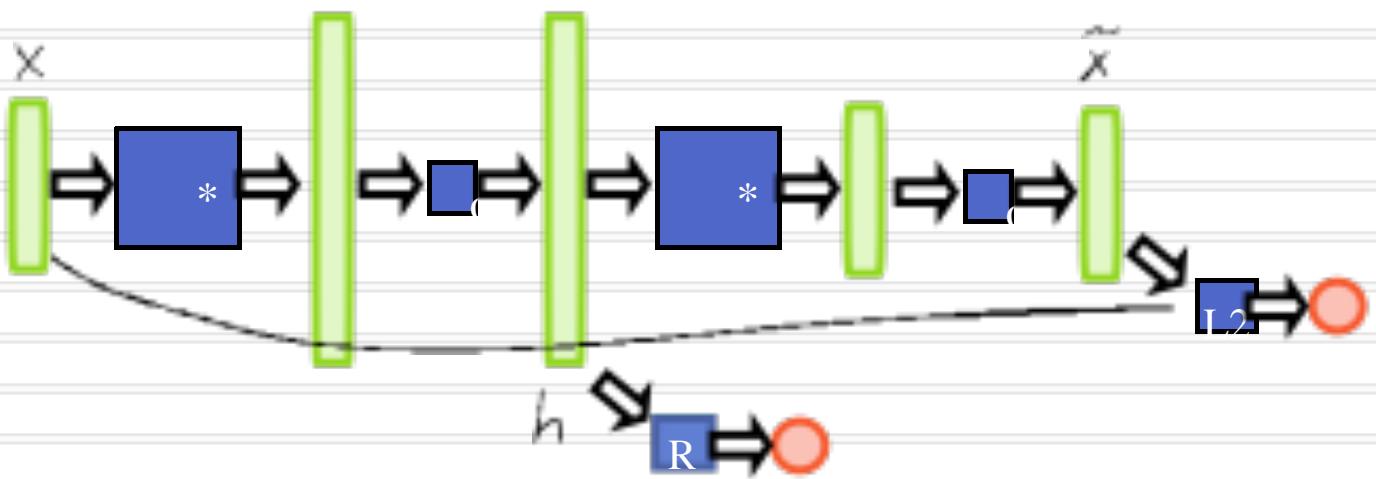


$$h = \sigma(wx)$$

$$\tilde{x} = \sigma(w h)$$

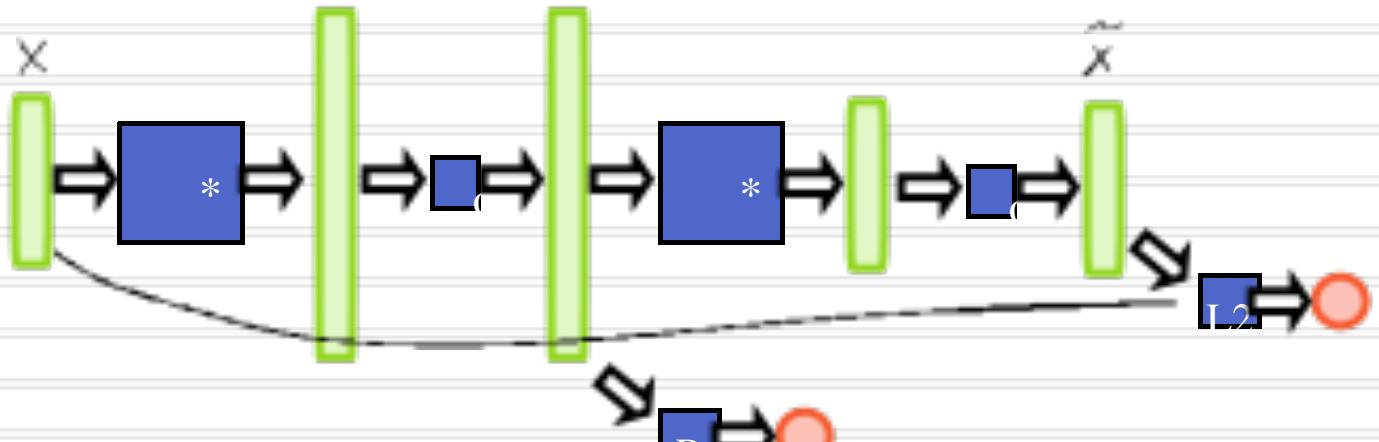
$$\min_{w, w'} \|x - \tilde{x}\|_2^2$$

Learning bigger/richer representation



- Naïve approach will learn an identity function
- Solution 1: regularize hidden representation

Sparse autoencoders



$$h = 6(w \times)$$

$$\tilde{x} = 6(w \cdot h)$$

$$\text{min}_{w, w'} \frac{1}{2} \|x - \tilde{x}\|_2^2 + \lambda R(h)$$

$$\sum_{i=1}^n \|h_i\|$$

Sparse autoencoders

$$\min_{\mathbf{w}, \mathbf{h}} \frac{1}{m} \frac{1}{n} \frac{1}{2} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 + \lambda R(\mathbf{h})$$

estimating from e.g. minibatch

$$R(\mathbf{h}) = \sum_i KL(B_{\text{real}}(p_i) \parallel B_{\text{model}}(p_i))$$
$$\sum_i p_i \log \frac{p_i}{\hat{p}_i} + (1 - p_i) \log \frac{1 - p_i}{1 - \hat{p}_i}$$

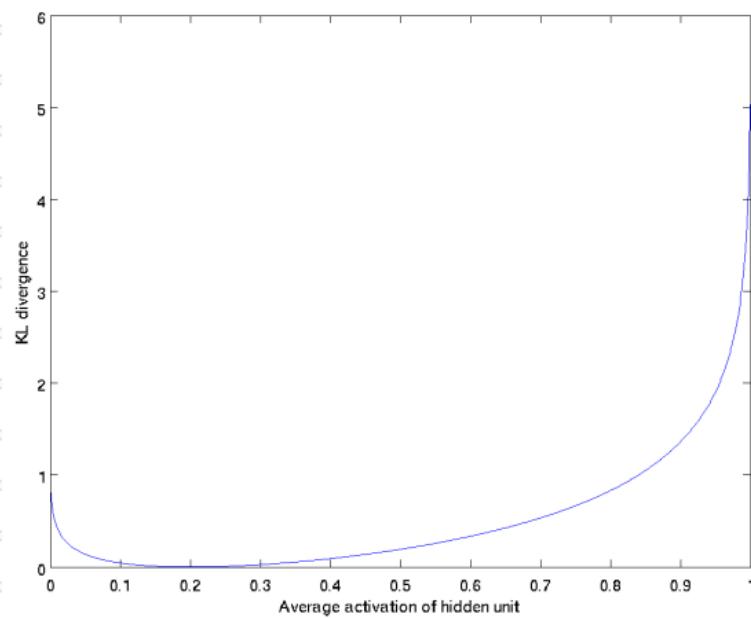
$$\frac{\partial R}{\partial p_i} = \frac{p_i}{\hat{p}_i} - \frac{1 - p_i}{1 - \hat{p}_i}$$

[Ng
2011]

Sparse autoencoders

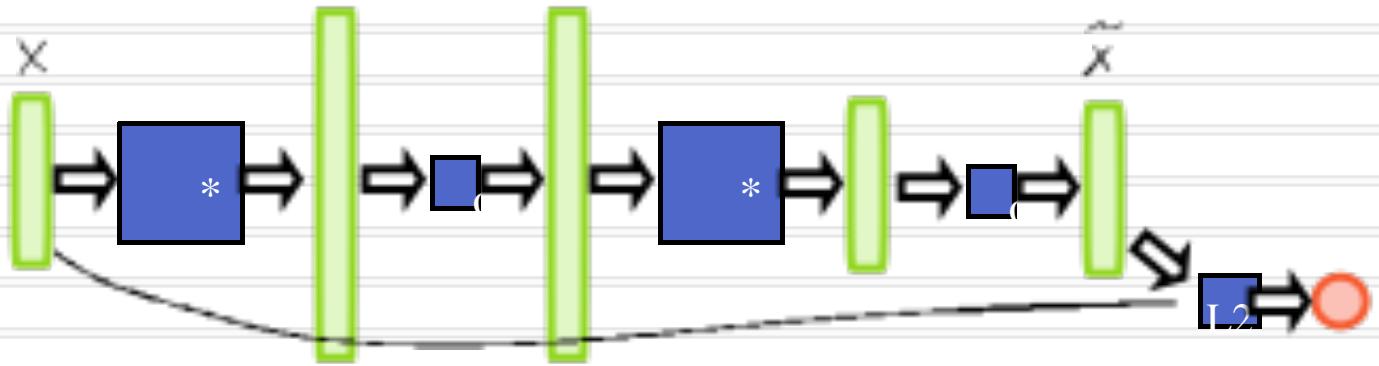
$$R(t) = \sum_i KL(\text{Bernoulli}(p_i) \parallel \text{Bernoulli}(\hat{p}_i))$$

$$\sum_i p_i \log \frac{p_i}{\hat{p}_i} + (1 - p_i) \log \frac{1 - p_i}{1 - \hat{p}_i}$$



[Ng
2011]

Learning bigger representation

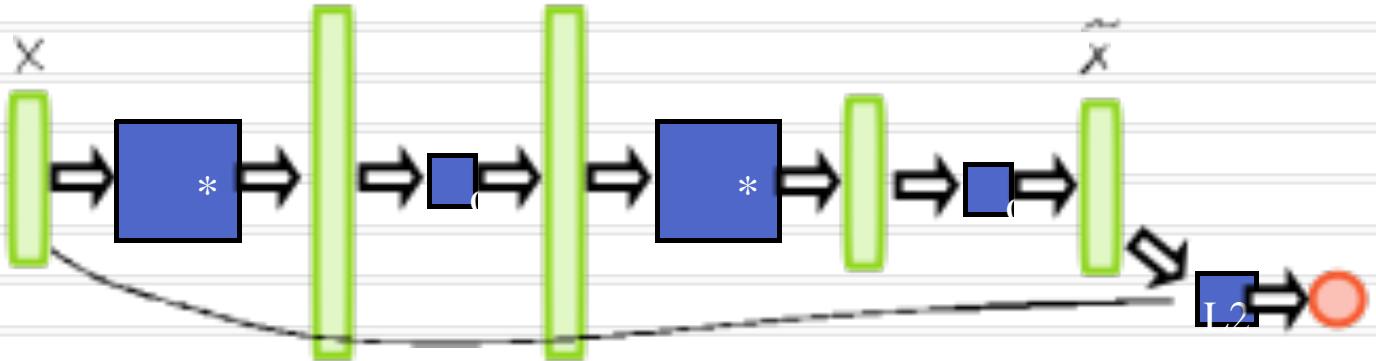


Naïve approach will learn an identity function

Solution 2: corrupt the input and reconstruct uncorrupted

[Vincent et al. 2010]

Denoising autoencoder



$$h = \sigma(w(x + n))$$

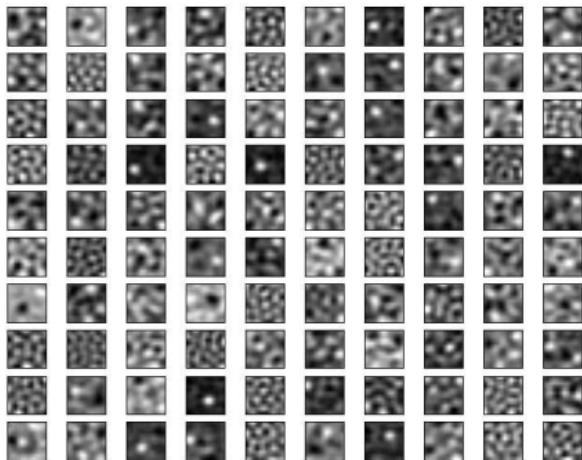
$$\tilde{x} = \sigma(w h)$$

$$m_n \sum_{w, w'} \|x_i - \tilde{x}_i\|_2^2$$

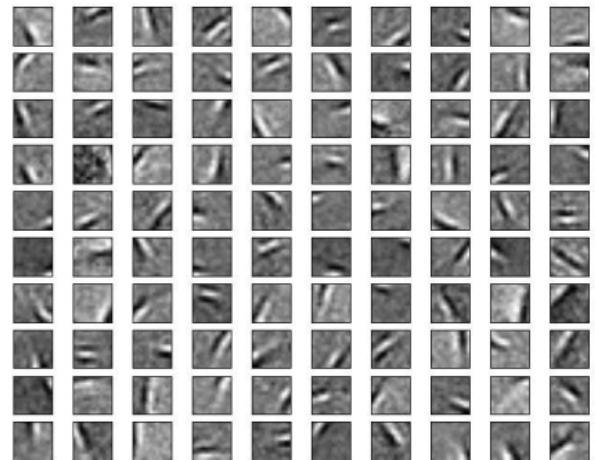
[Vincent et al. 2010]

Denoising autoencoder

12x12 patches, 200 hidden units

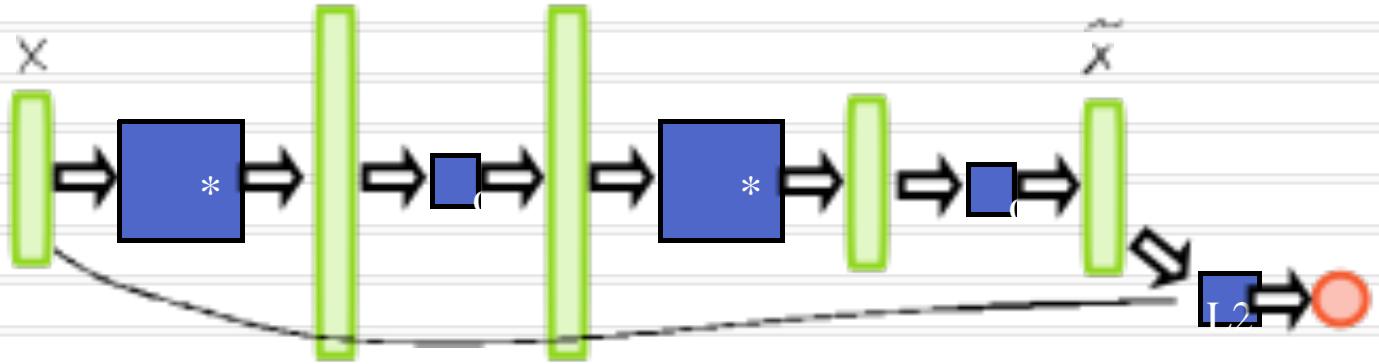


Autoencoder with
weight-decay



Denoising
autoencoder
(Gaussian noise)
[Vincent et al. 2010]

Contractive autoencoders



$$h = \sigma(wx)$$

$$\hat{x} = \sigma(w h)$$

$$\min_{w, w'} \|x - \hat{x}\|_2^2 + \lambda \left\| \frac{\partial h(x)}{\partial x_j} \right\|_1^2$$

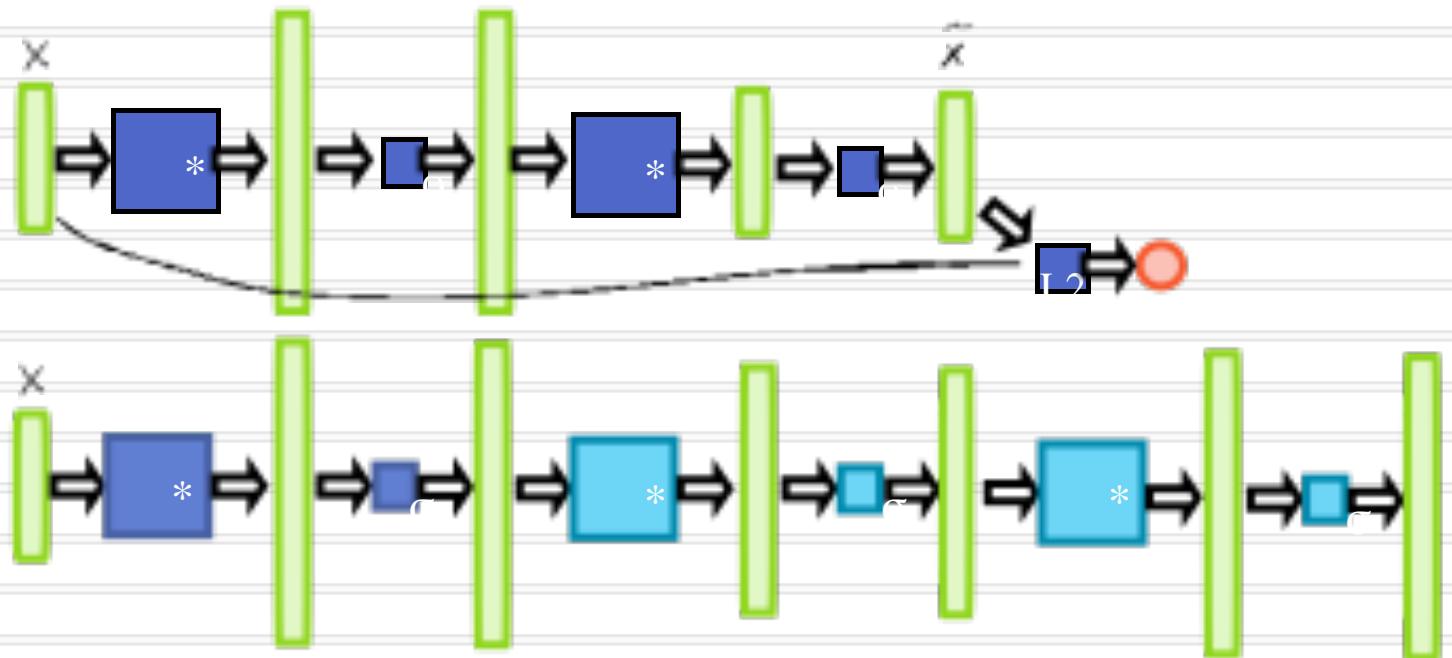
$$\frac{\partial h}{\partial x_j} = (h_j \cdot (1 - h_j)) \cdot w_j$$

encoder

[Rifai et al.

2011]

Training deep autoencoders

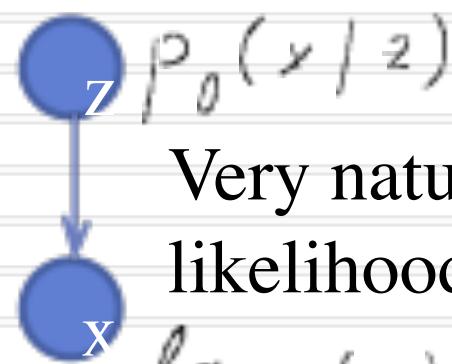


Historically: layerwise pretraining

Most modern implementations: end-to-end

Commonly including convolutional and up-convolutional layers

Maximum likelihood



Very natural idea: train θ by maximum likelihood:

$$\log p_\theta(x) - \log \int_z p_\theta(x, z) dz \rightarrow \max$$

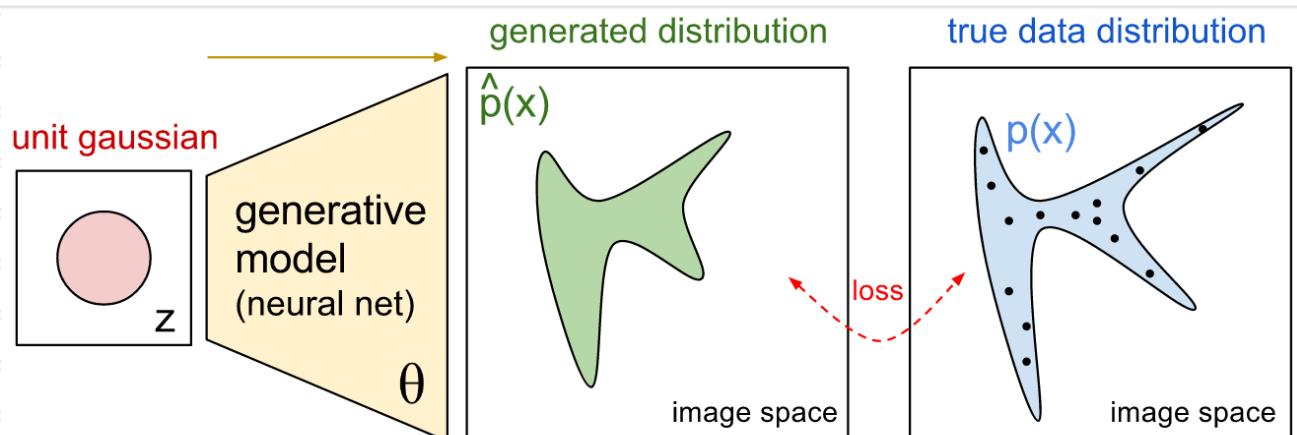
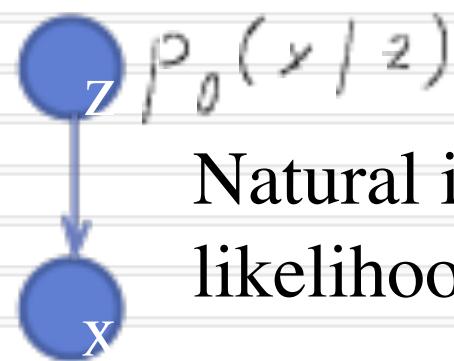


Image credit: OpenAI blog

Maximum likelihood via VAE



Natural idea: train θ by maximum likelihood.

$$\log p_\theta(x) = \log \int_z p_\theta(x, z) dz \rightarrow \max$$

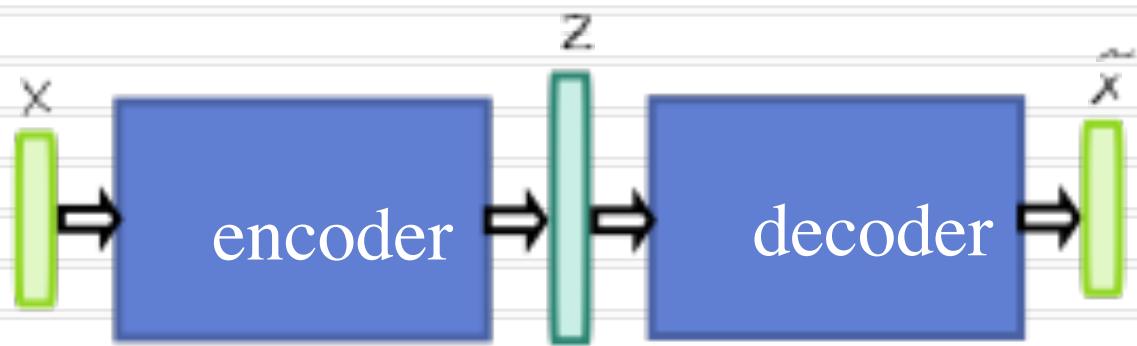
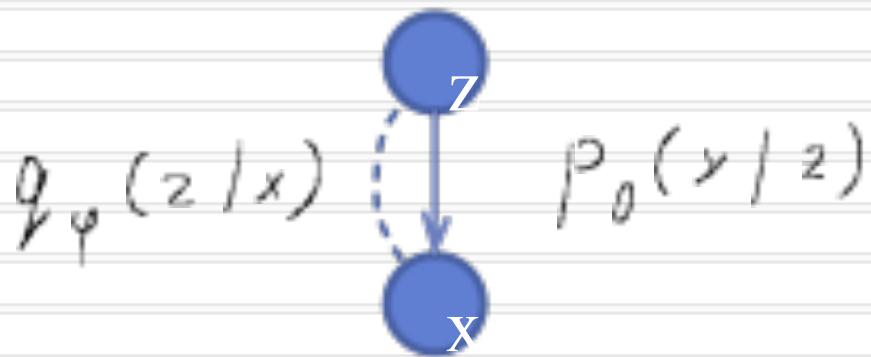
The integral is intractable.

VAE solution: introduce an *encoder*

$$q_\varphi(z|x) \approx p(z|x)$$

Derive an approximation to ML objective

Variational autoencoder as a graphical model

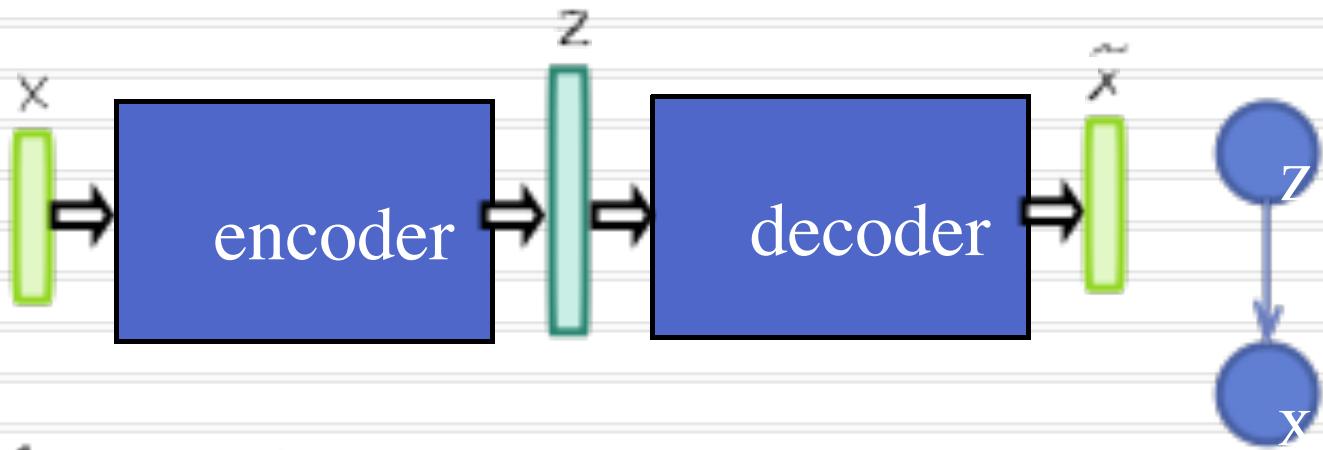


[Kingma & Welling
14]

Variational lower bound

$$\begin{aligned} \log p(x) - \log \int_z p(x, z) dz &= \\ - \log \int_z p(x, z) \frac{q(z|x)}{q(z|x)} dz &= \\ = \log E_{q(z|x)} \frac{p(x, z)}{q(z|x)} &= \\ \log E_{q(z|x)} \frac{p(x|z) p(z)}{q(z|x)} & \geq \\ \geq E_{q(z|x)} \log p(x|z) - E_{q(z|x)} \log \frac{p(z)}{q(z|x)} &= \\ - E_{q(z|x)} \log p(x|z) - KL(q(z|x) || p(z)) & \end{aligned}$$

Variational lower-bound



$$\log p(x) \geq$$

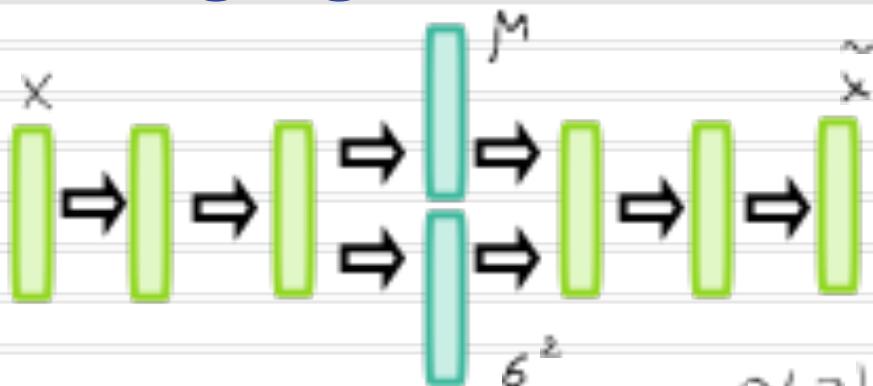
$$KL(q_{\varphi}(z|x) || p(z)) + E_{q_{\varphi}(z|x)} [\log p_{\theta}(x|z)]$$

regularization

[Kingma & Welling

~"denoising auto-encoder"

Minimizing regularization term



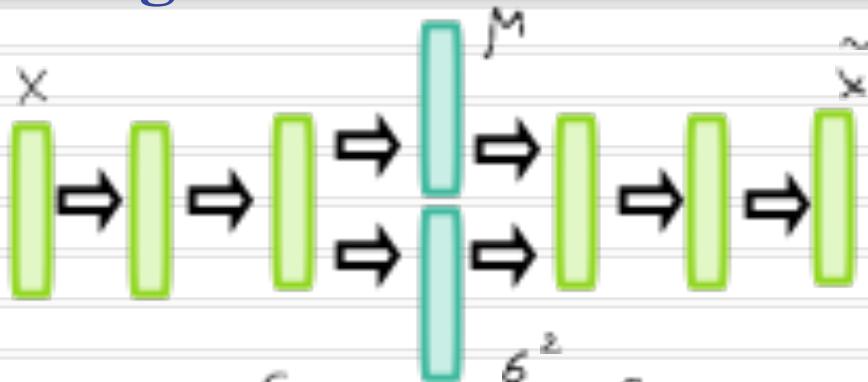
$$p(z) = \mathcal{N}(0, 1)$$

$$KL(g_{\phi}(z|x) || p(z))$$

$$= \frac{1}{2} \sum \left(\mu_j^2 + \sigma_j^2 - 1 - \log \sigma_j^2 \right)$$

[Kingma & Welling
14]

Optimizing reconstruction term



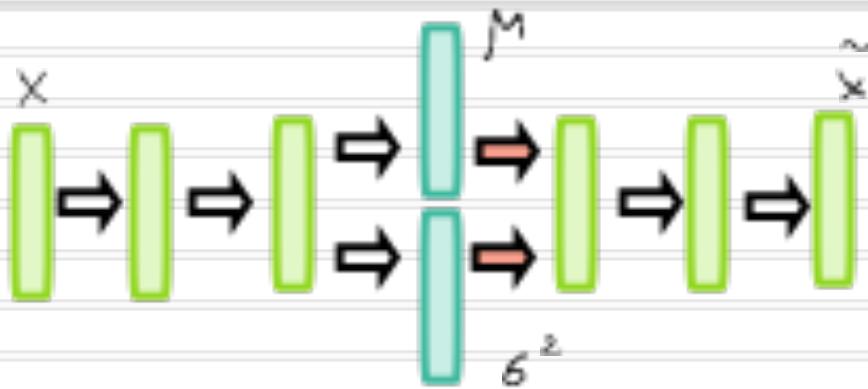
$$E_{\mathcal{D}_p(z|x)} [\log p_\theta(\tilde{x}|z)] \rightarrow \max$$

$$p_\theta(x|z) \propto \exp(-\|x - \tilde{x}\|^2)$$

$$E_{\mathcal{D}_p(z|x)} \|d_\theta(z) - x\|_2^2 \rightarrow \min$$

[Kingma & Welling
14]

Reparameterization trick



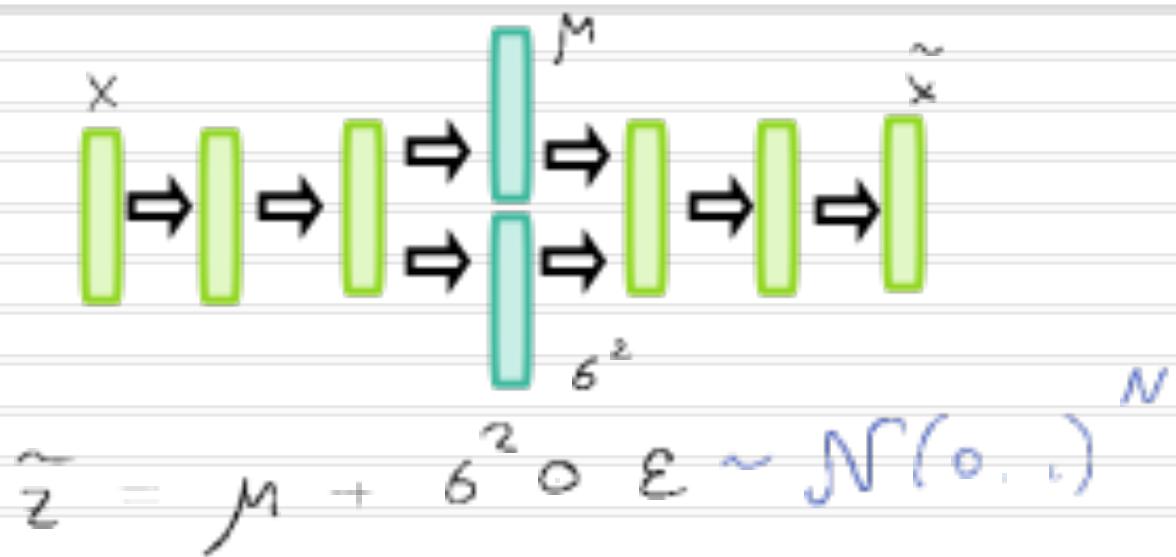
$$L_{\mathcal{G}_\theta}(z|x) \parallel d_\theta(z) - x \parallel_2^2 \rightarrow \min$$

$$\tilde{z} = \mu + \sigma^2 \circ \varepsilon \sim \mathcal{N}(0, \cdot)$$

$$L = \parallel d_\theta(\tilde{z}) - x \parallel_2^2 \rightarrow \min$$

[Kingma & Welling
14]

Reparameterization trick

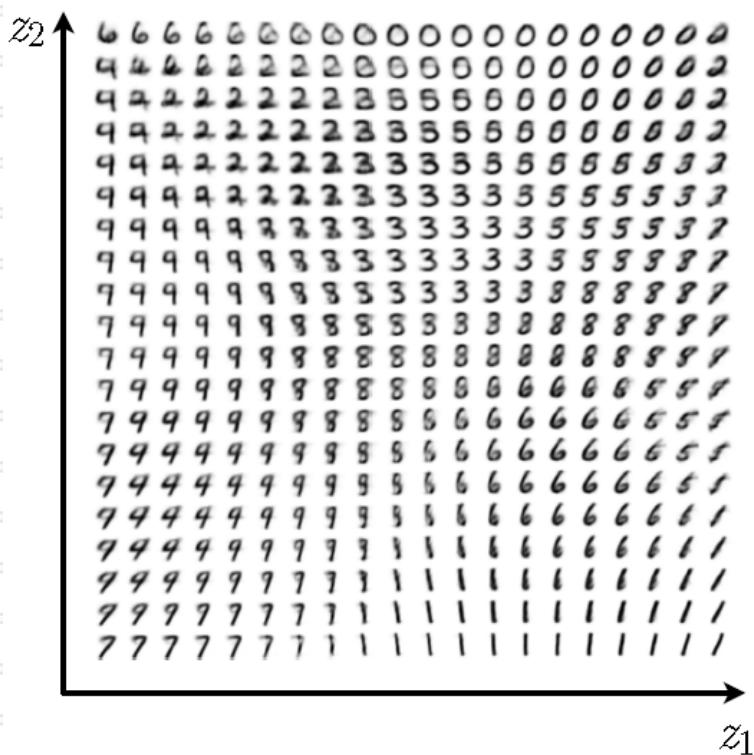


$$\frac{dL}{d\mu} = \frac{dL}{d\tilde{z}} = \frac{dL}{d\tilde{z}^2} = \frac{dL}{d\tilde{z}} \circ \epsilon$$

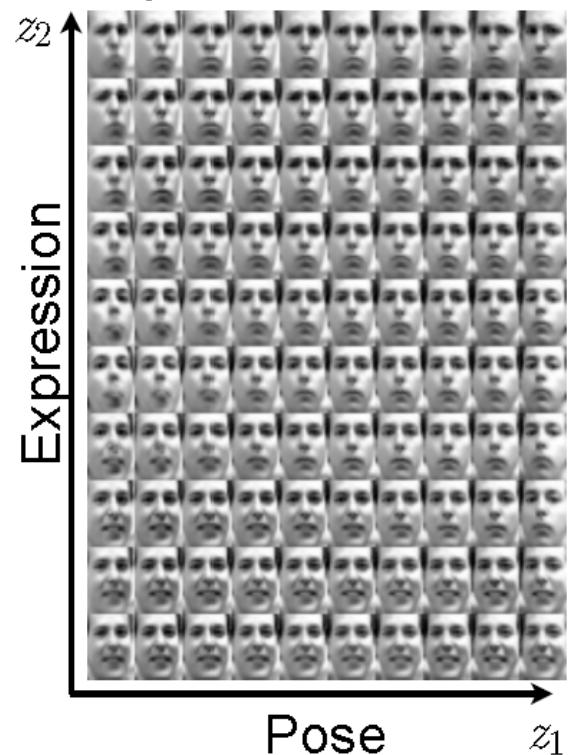
[Kingma & Welling
14]

VAE-learned manifolds

MNIST:



Frey Face dataset:



[Kingma & Welling
14]

VAE-learned manifolds

6 6 1 7 8 1 4 8 2 8
9 6 0 3 9 6 0 3 1 9
3 9 9 1 3 6 9 1 7 9
8 9 0 8 6 9 1 9 6 3
8 2 3 3 3 3 1 3 3 6
6 9 9 8 6 1 6 6 6 6
9 5 2 6 6 5 1 8 9 9
9 9 8 1 8 1 2 8 2 3
0 4 6 1 2 3 2 0 8 8
9 7 5 9 9 3 4 8 5 1

(a) 2-D latent space

(b) 5-D latent space

2 8 3 1 3 8 5 7 3 8
8 3 8 2 7 9 3 3 3 8
8 5 9 9 4 3 9 5 1 6
1 9 8 8 8 3 3 1 9 7
2 7 3 6 4 3 0 2 0 3
5 9 7 0 5 9 2 8 4 5
6 9 4 3 6 2 8 5 5 2
8 4 9 0 8 0 7 0 6 6
7 4 3 6 3 0 3 6 0 1
2 1 2 0 4 7 1 0 6 0

(c) 10-D latent space

(d) 20-D latent space

[Kingma & Welling
14]

Collapsing components

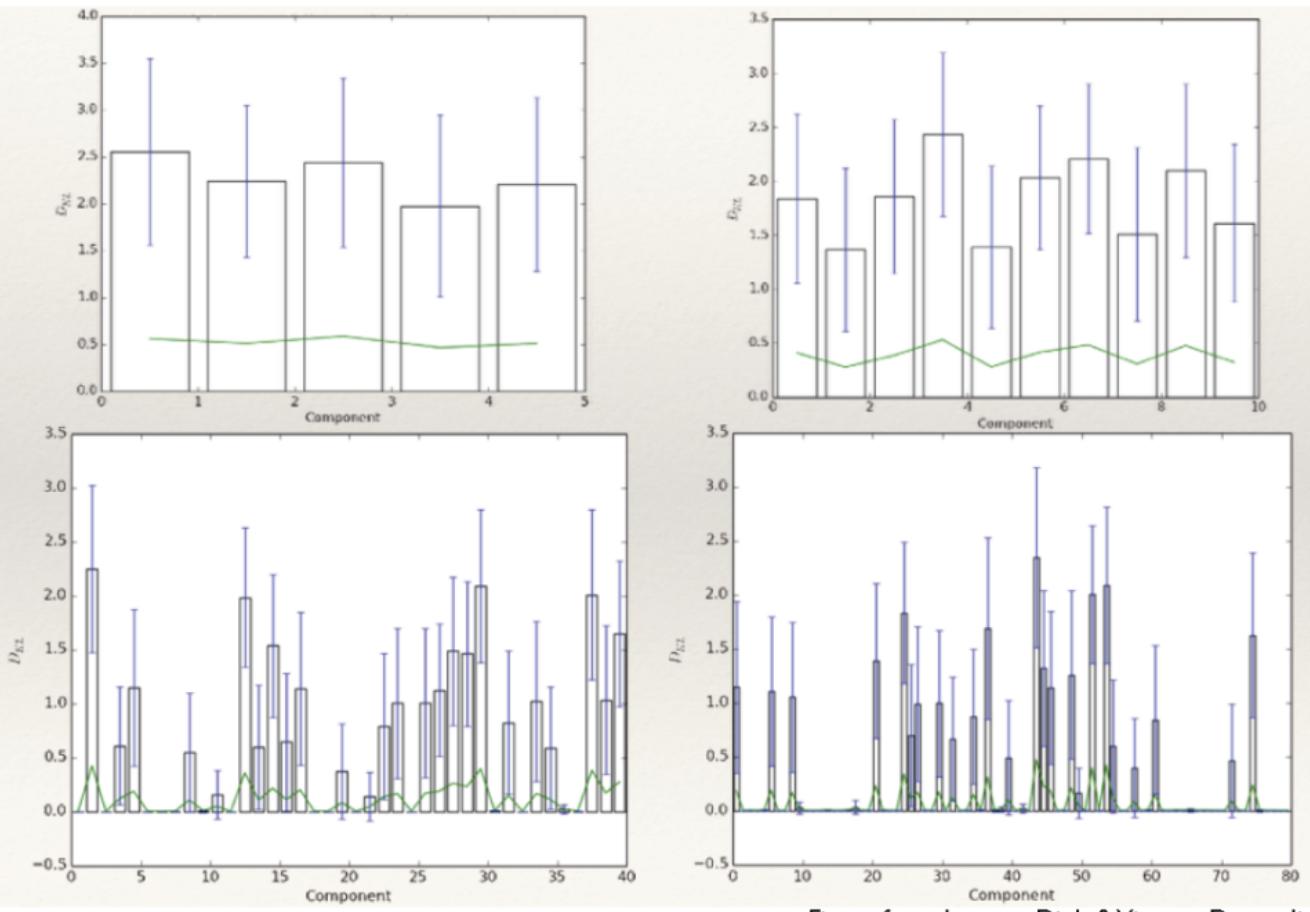


Figure from Laurent Dinh & Vincent Dumoulin

Collapsing components

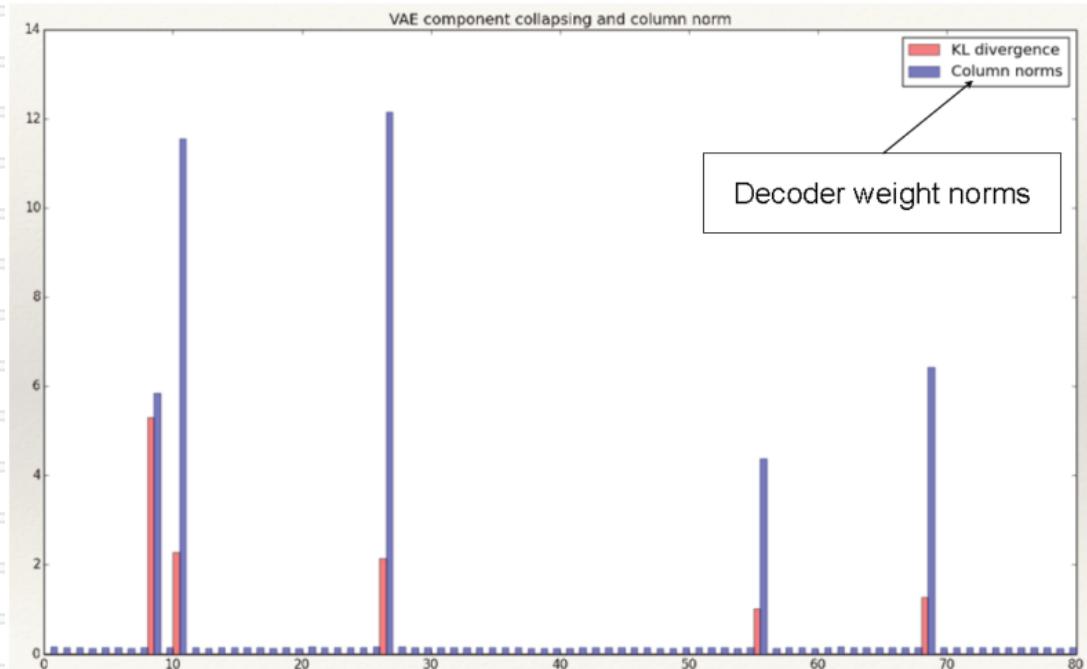


Figure from Laurent Dinh & Vincent Dumoulin

- Components that are shrunk to $N(0,1)$ are not used by decoder

Sample quality comparison

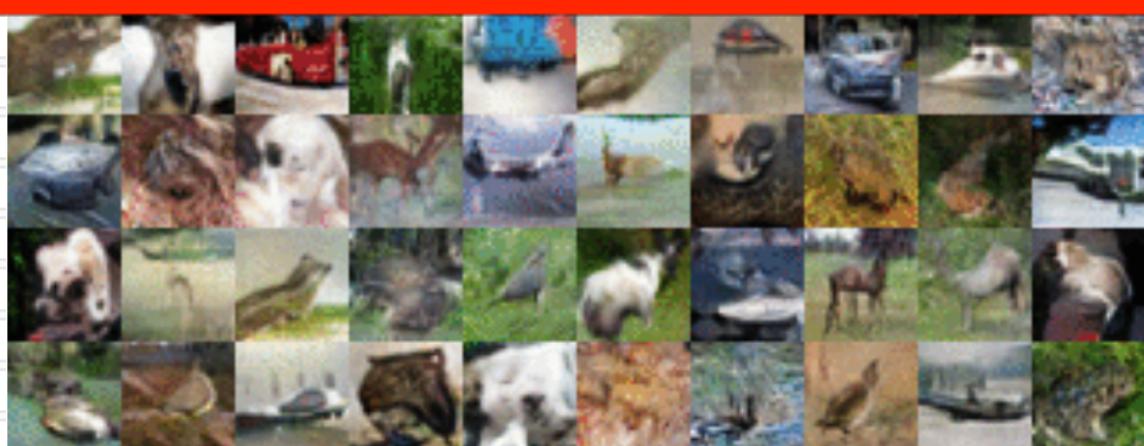


VA

E

GA

N



Source:
OpenAI
blog