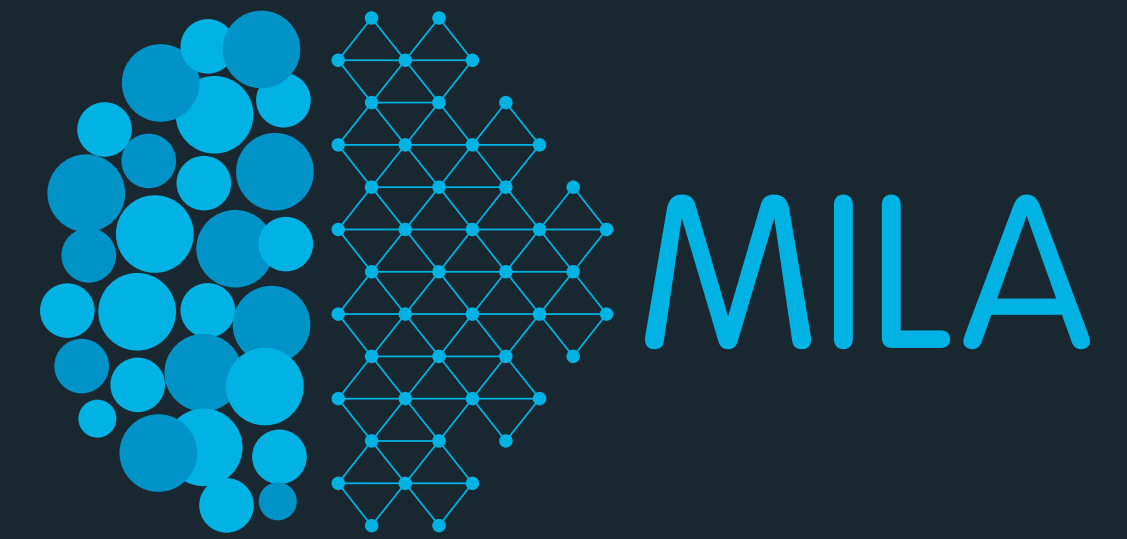


Institut
des algorithmes
d'apprentissage
de Montréal



Deep Generative Models

Aaron Courville
MILA, Université de Montréal

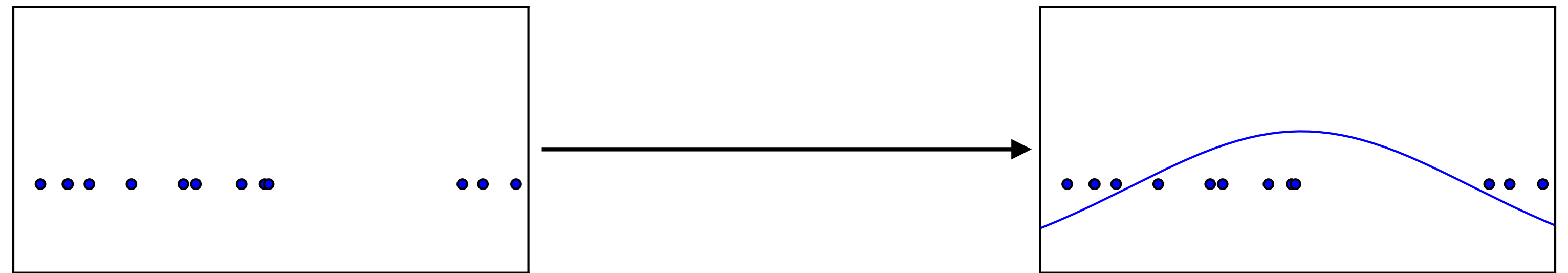
6.S191: Introduction to Deep Learning

MIT, Jan 30th, 2018

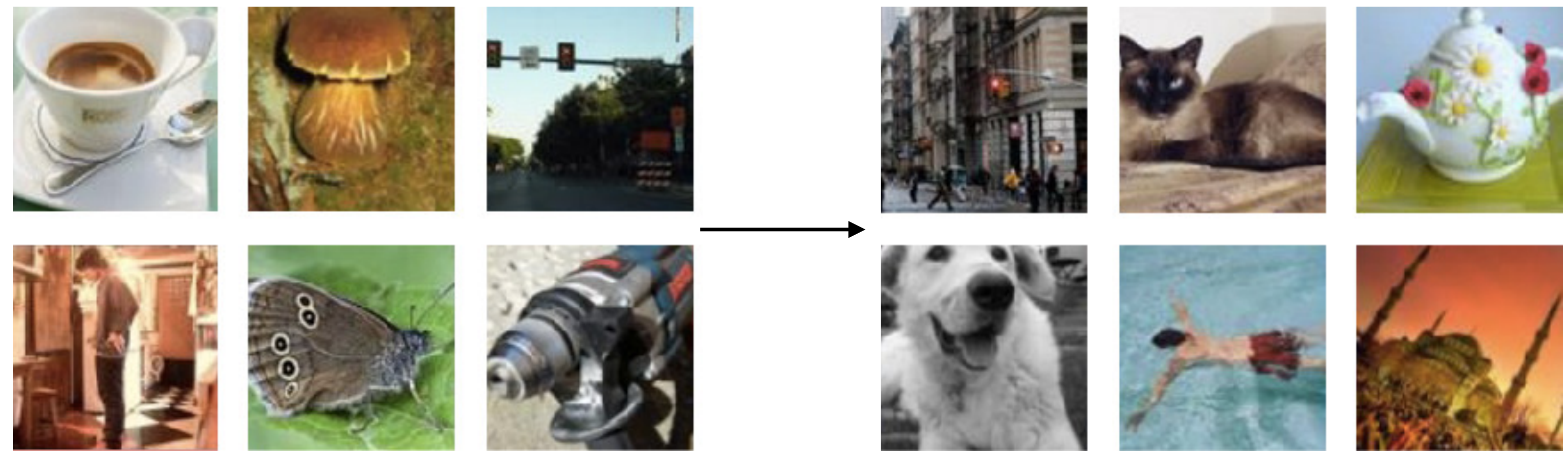
Generative modeling

- Generative models take training samples from some data distribution and learn a model that represents that distribution.

- Density estimation:



- Sample generation:



Training examples

Model samples

Why generative models?



- Many tasks require structured output
 - Eg. Machine translation

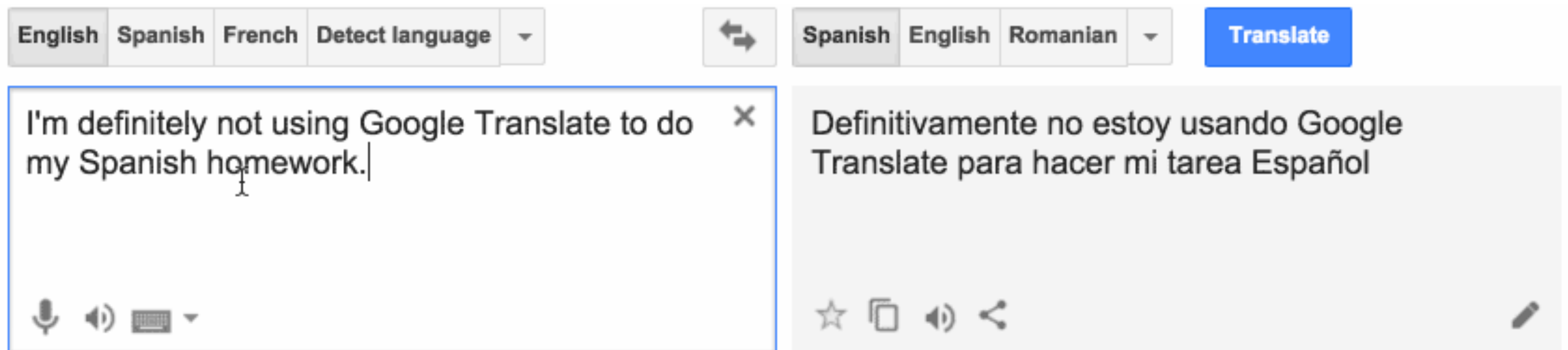
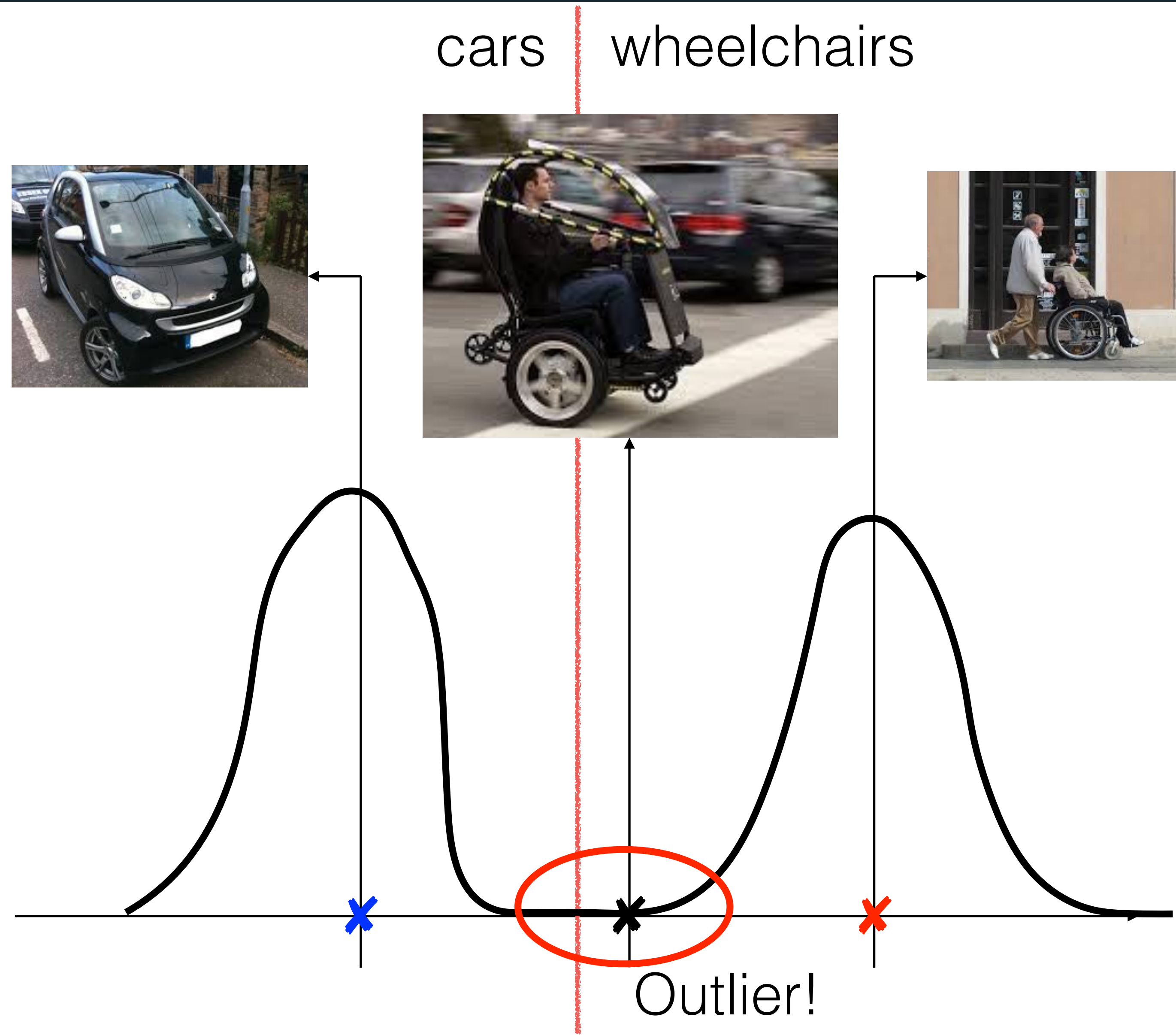


image credit: Adam Geitgey blog (2016) *Machine Learning is Fun Part 5: Language Translation with Deep Learning and the Magic of Sequences*

Why Generative Models? Outlier detection

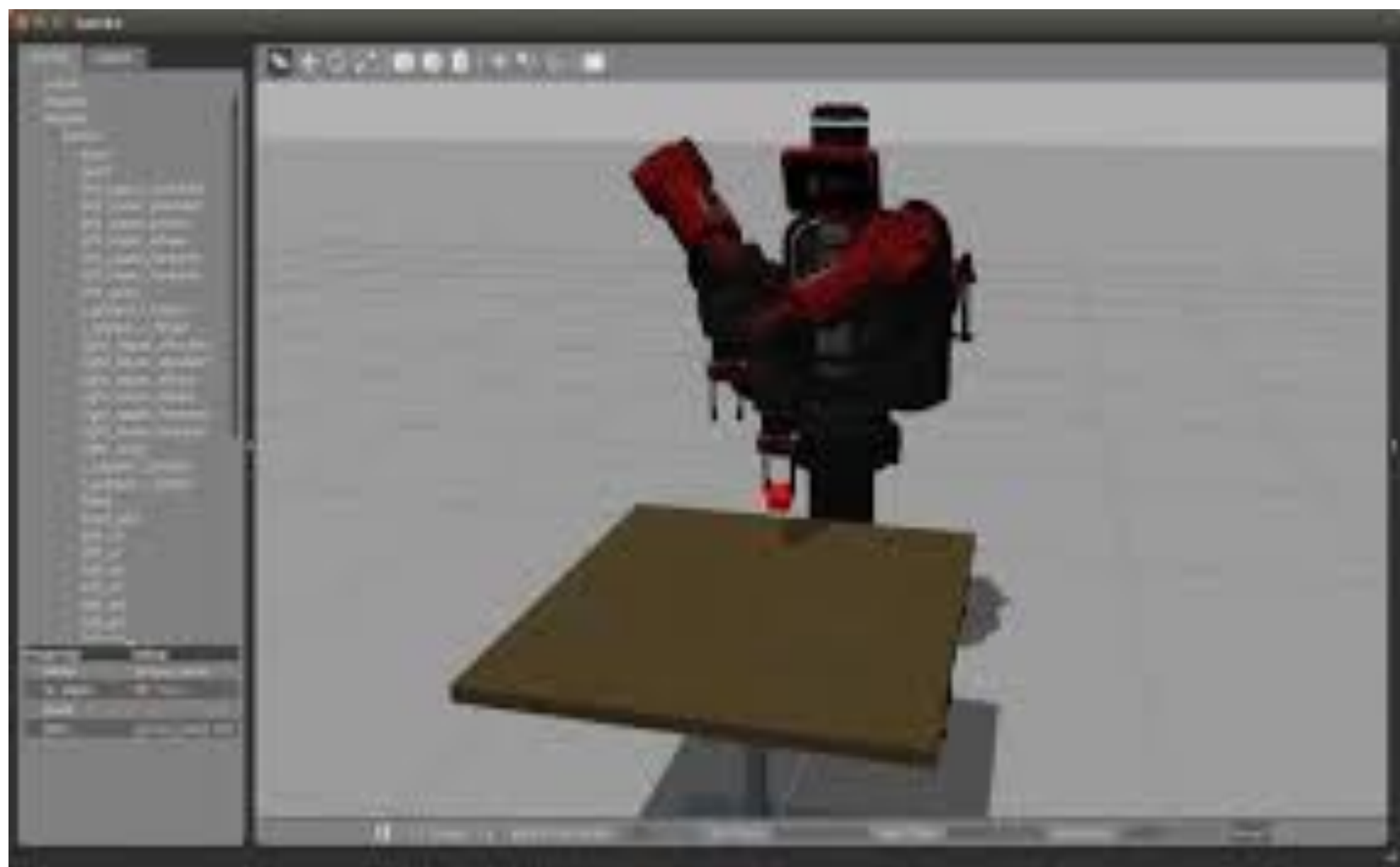
- Large-scale deployment of CNN-based perception systems is becoming a reality.
- How do we detect when we encounter something new or rare (i.e. not appearing in the training data)?
- **Goal:** detect these outliers (anomalies) to avoid dangerous misclassification.
- **Strategy:** Leverage generative models of the training distribution to detect outliers.



Why Generative Models? Generation for Simulation



- Supports Reinforcement Learning for Robotics: Make simulations sufficiently realistic that learned policies can readily transfer to real-world application



Generative model



Photo from IEEE Spectrum

Deep Generative Models: Outline



Autoregressive models

- Deep NADE, PixelRNN, PixelCNN, WaveNet, Video Pixel Network, etc.

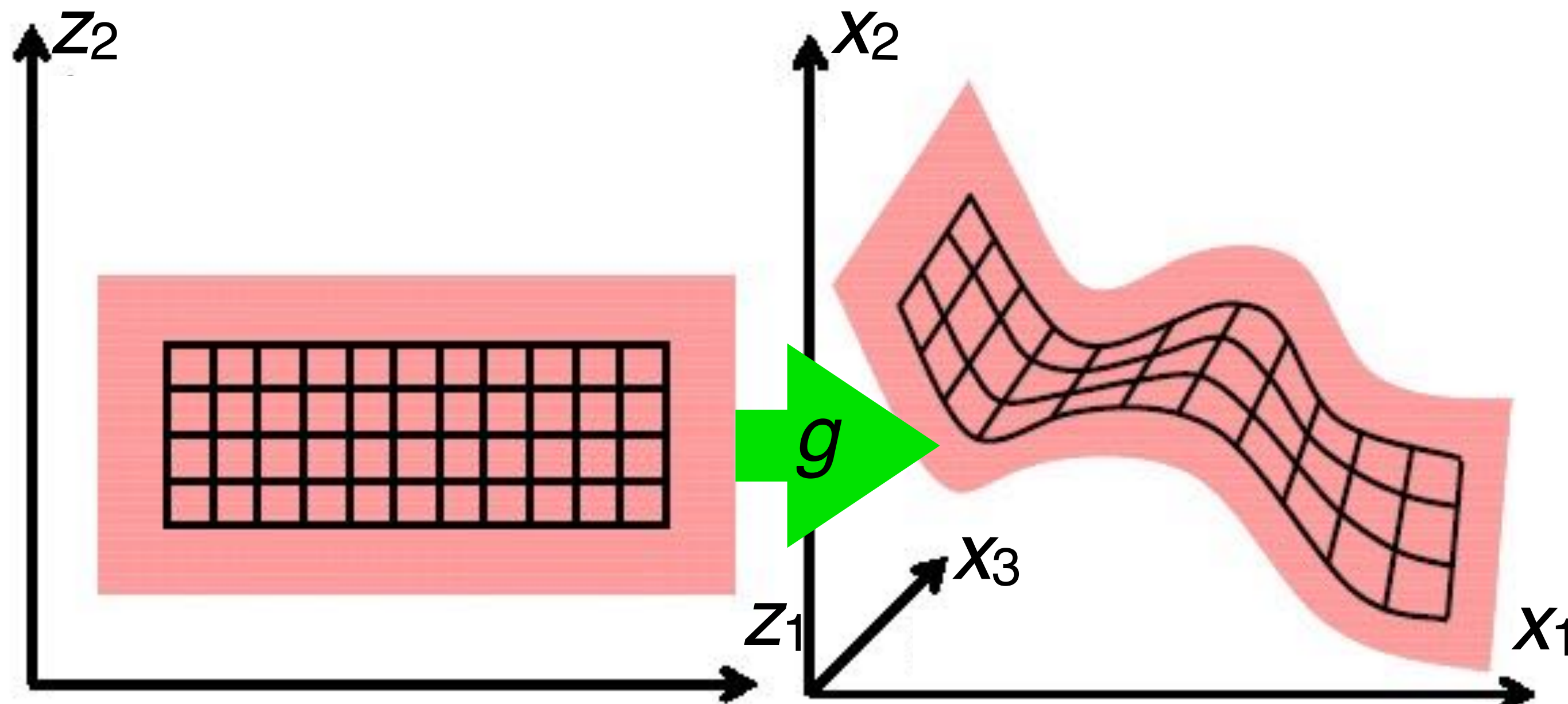
Latent variable models

- Variational Auto encoders
- Generative Adversarial Networks

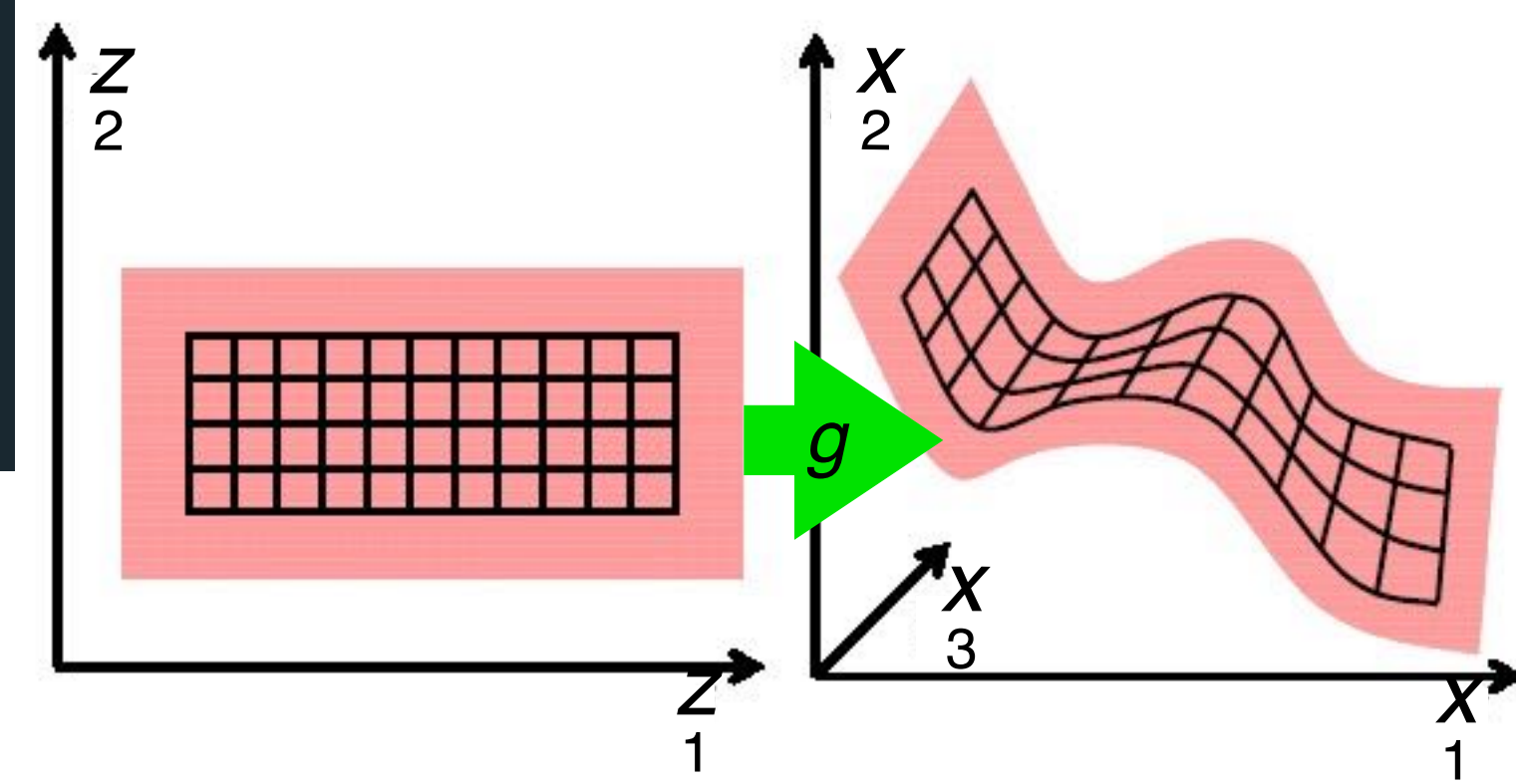
our focus today

Latent Variable Models

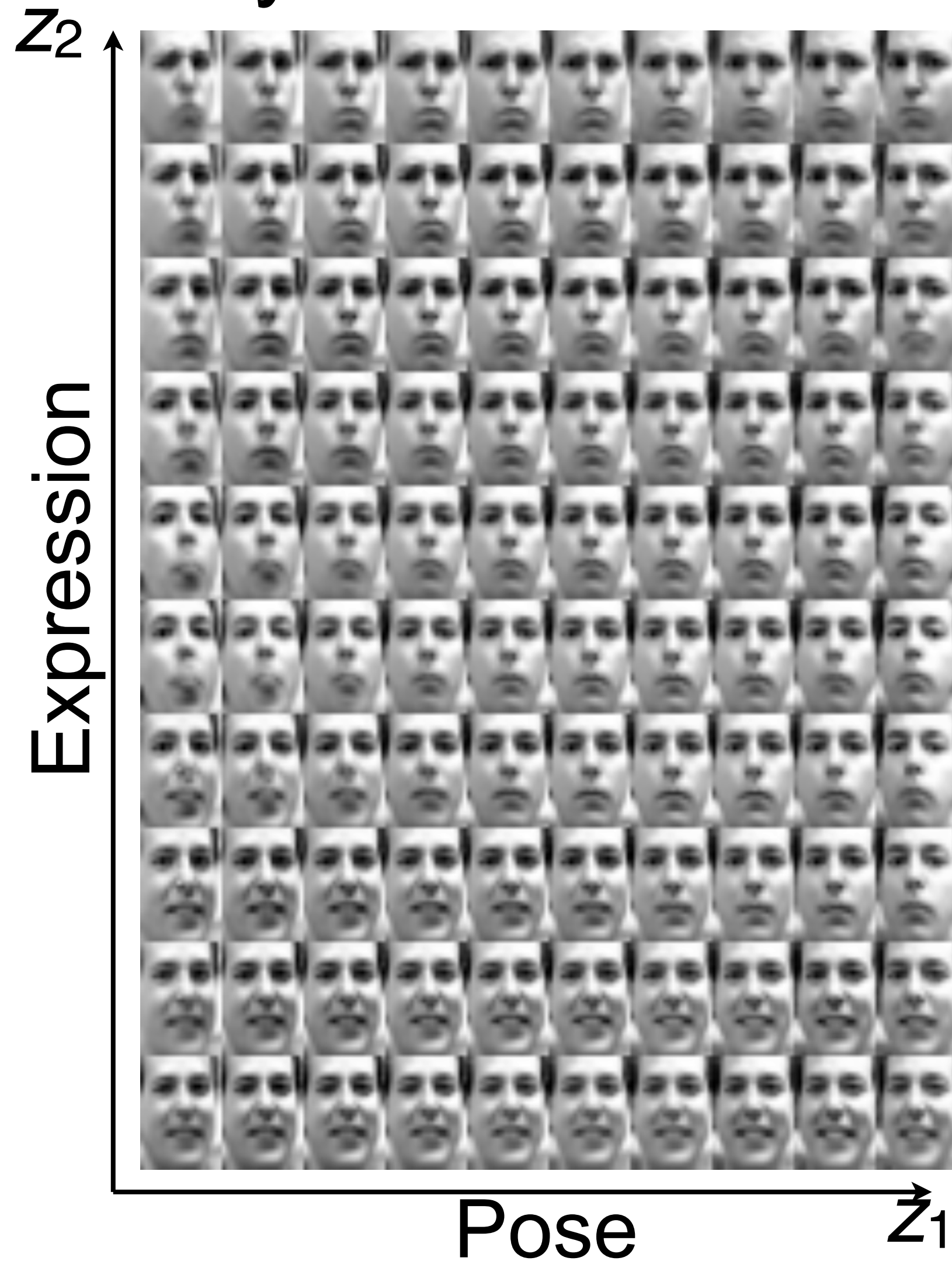
- The Variational Autoencoder model:
 - Kingma and Welling, *Auto-Encoding Variational Bayes*, *International Conference on Learning Representations (ICLR)* 2014.
 - Rezende, Mohamed and Wierstra, *Stochastic back-propagation and variational inference in deep latent Gaussian models*. ICML 2014.



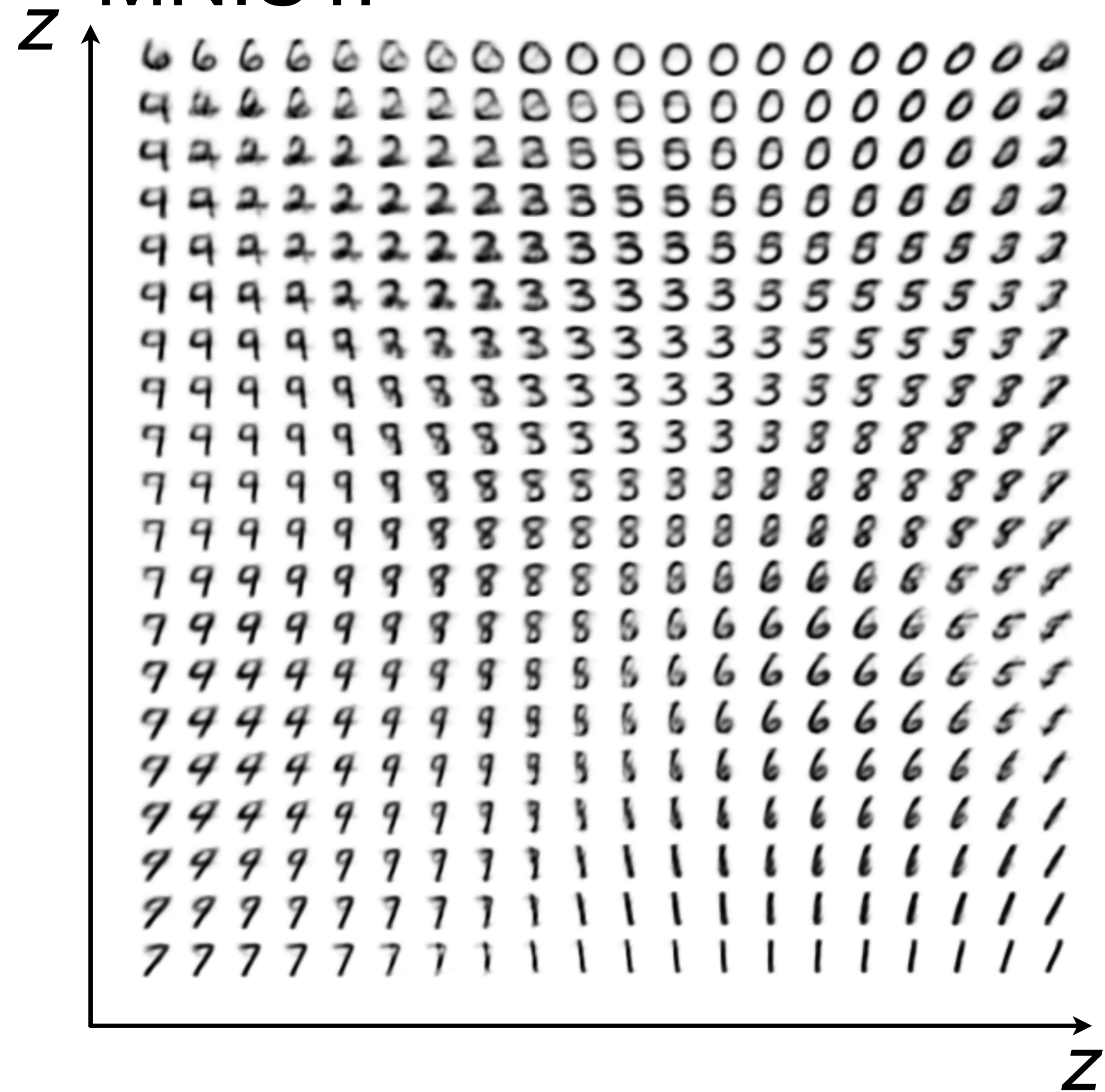
Latent Variable Models



Frey Faces:



MNIST:



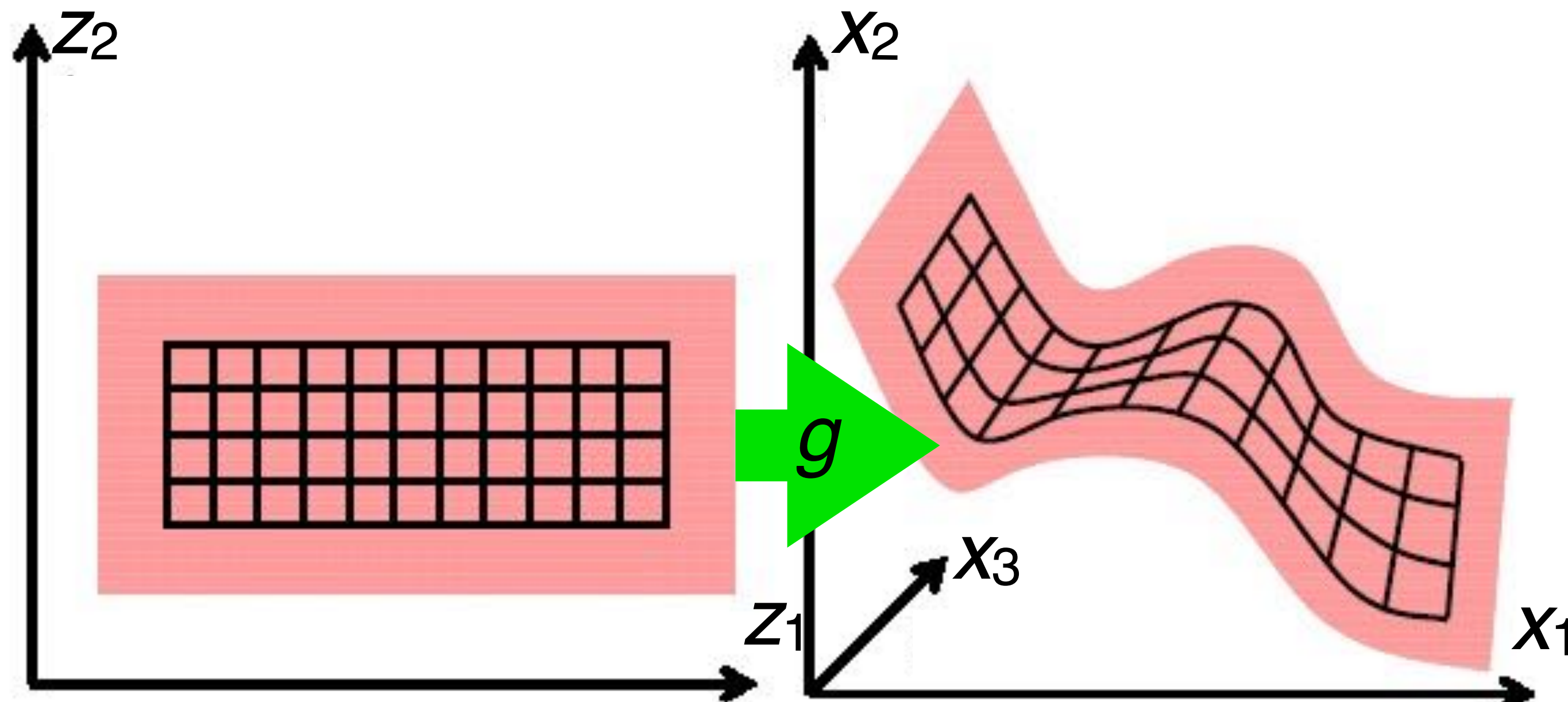
Latent Variable Models

- **latent variable model**: learn a mapping from some latent variable z to a complicated distribution on x .

$$p(x) = \int p(x, z) dz \quad \text{where } p(x, z) = p(x | z)p(z)$$

Prior $p(z) = \text{something simple}$ $p(x | z) = g(z)$

- Can we learn to decouple the true **explanatory factors** underlying the data distribution?
E.g. separate identity and expression in face images



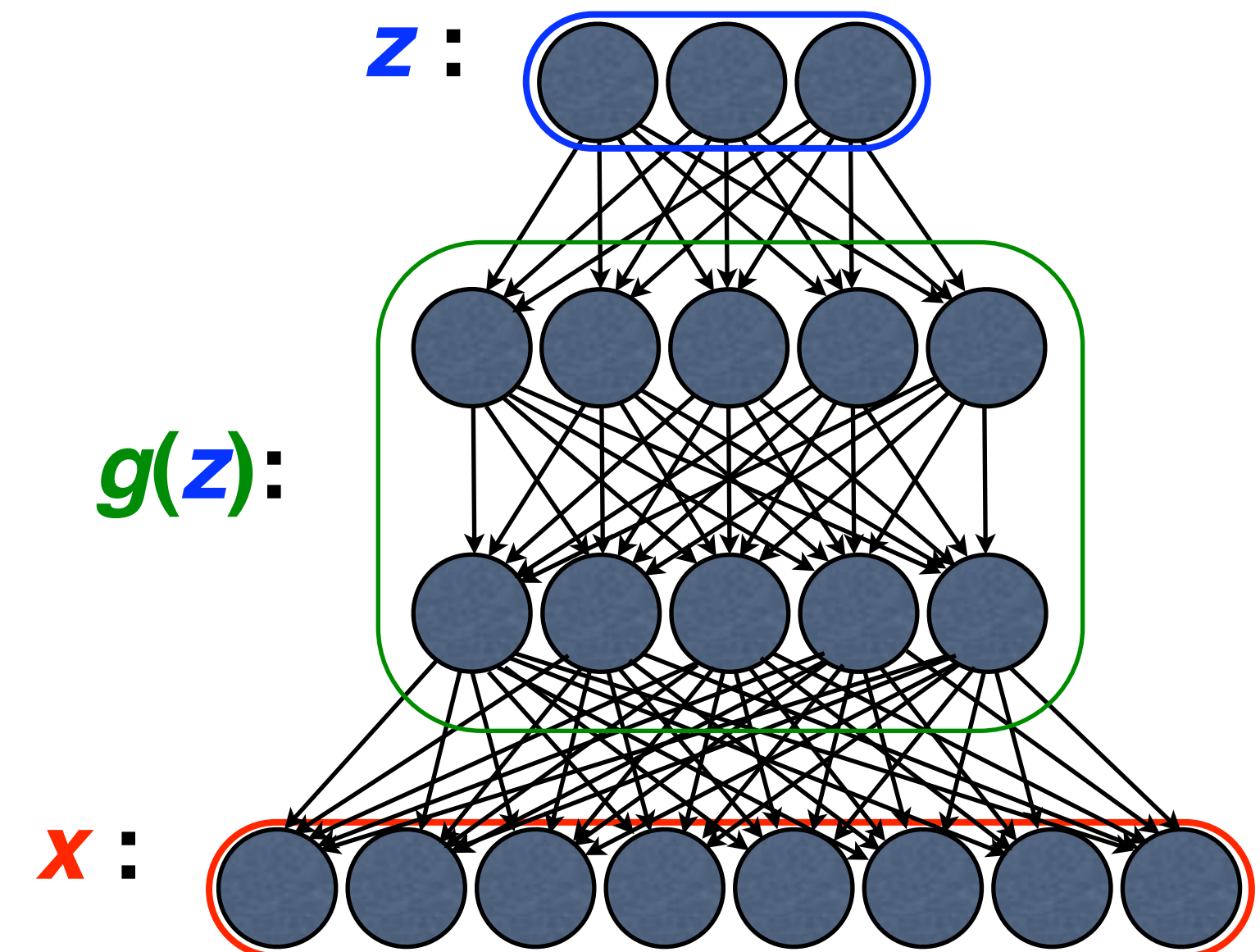
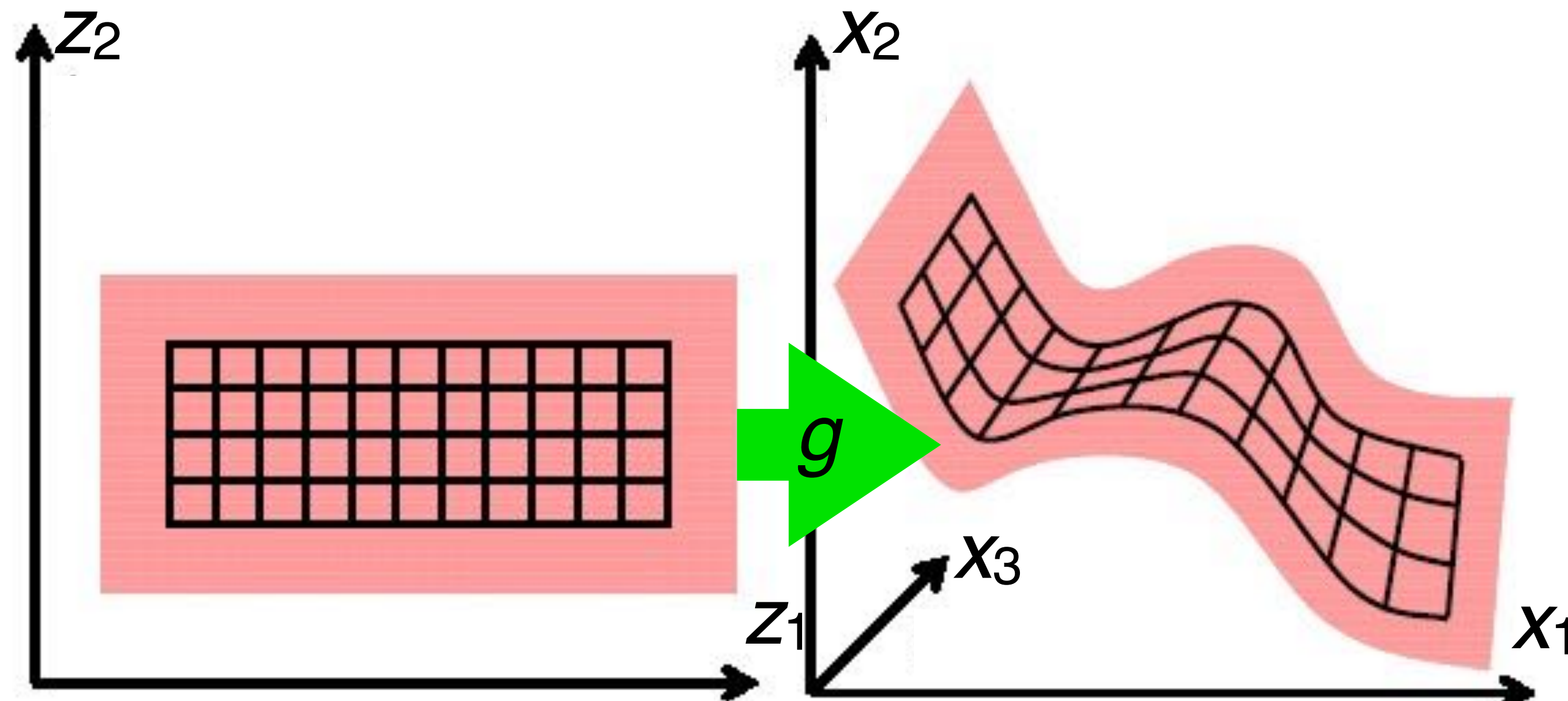
Latent Variable Models

- **latent variable model**: learn a mapping from some latent variable z to a complicated distribution on x .

$$p(x) = \int p(x, z) dz \quad \text{where } p(x, z) = p(x | z)p(z)$$

Prior $p(z) = \text{something simple}$ $p(x | z) = g(z)$

- Can we learn to decouple the true **explanatory factors** underlying the data distribution?
E.g. separate identity and expression in face images



Variational Auto-Encoder (VAE)



- Where does z come from? — The classic DAG problem.
- The VAE approach: introduce an inference machine $q_\phi(z | x)$ that **learns** to approximate the posterior $p_\theta(z | x)$.
- Define a **variational lower bound** on the data likelihood: $\log p_\theta(x) \geq \mathcal{L}(\theta, \phi, x)$

$$\begin{aligned}\log p_\theta(x) &\geq \log p_\theta(x) - D_{\text{KL}}[q_\phi(z|x) \| p_\theta(z|x)] \\ &= \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x) + \log p_\theta(z|x) - \log q_\phi(z|x)] \\ &= \mathcal{L}(\theta, \phi, x) \\ &= \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z) + \log p_\theta(z) - \log q_\phi(z|x)] \\ &= \underbrace{-D_{\text{KL}}[q_\phi(z|x) \| p_\theta(z)]}_{\text{regularization term}} + \underbrace{\mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)]}_{\text{reconstruction term}}\end{aligned}$$

- What is $q_\phi(z | x)$?

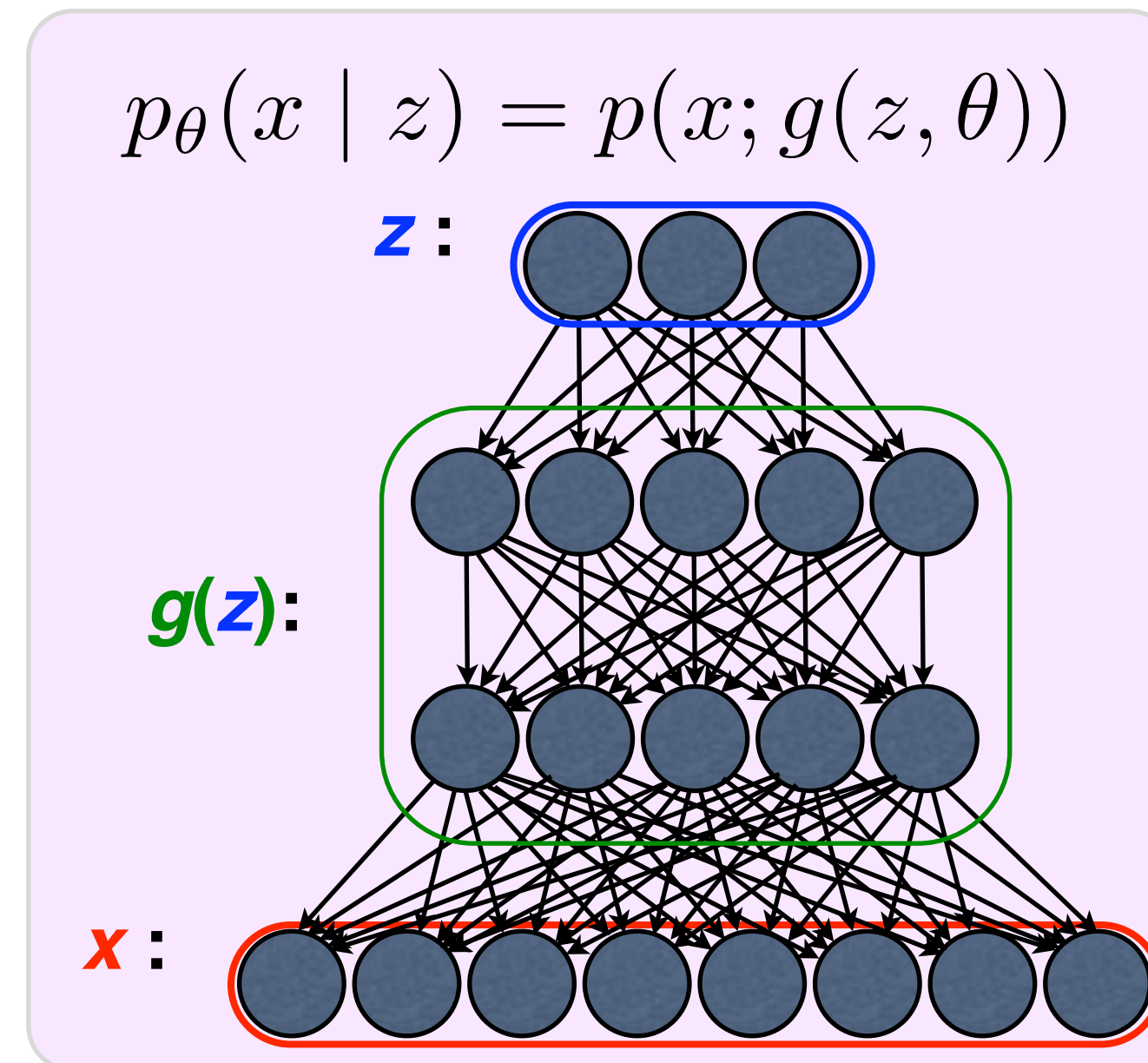
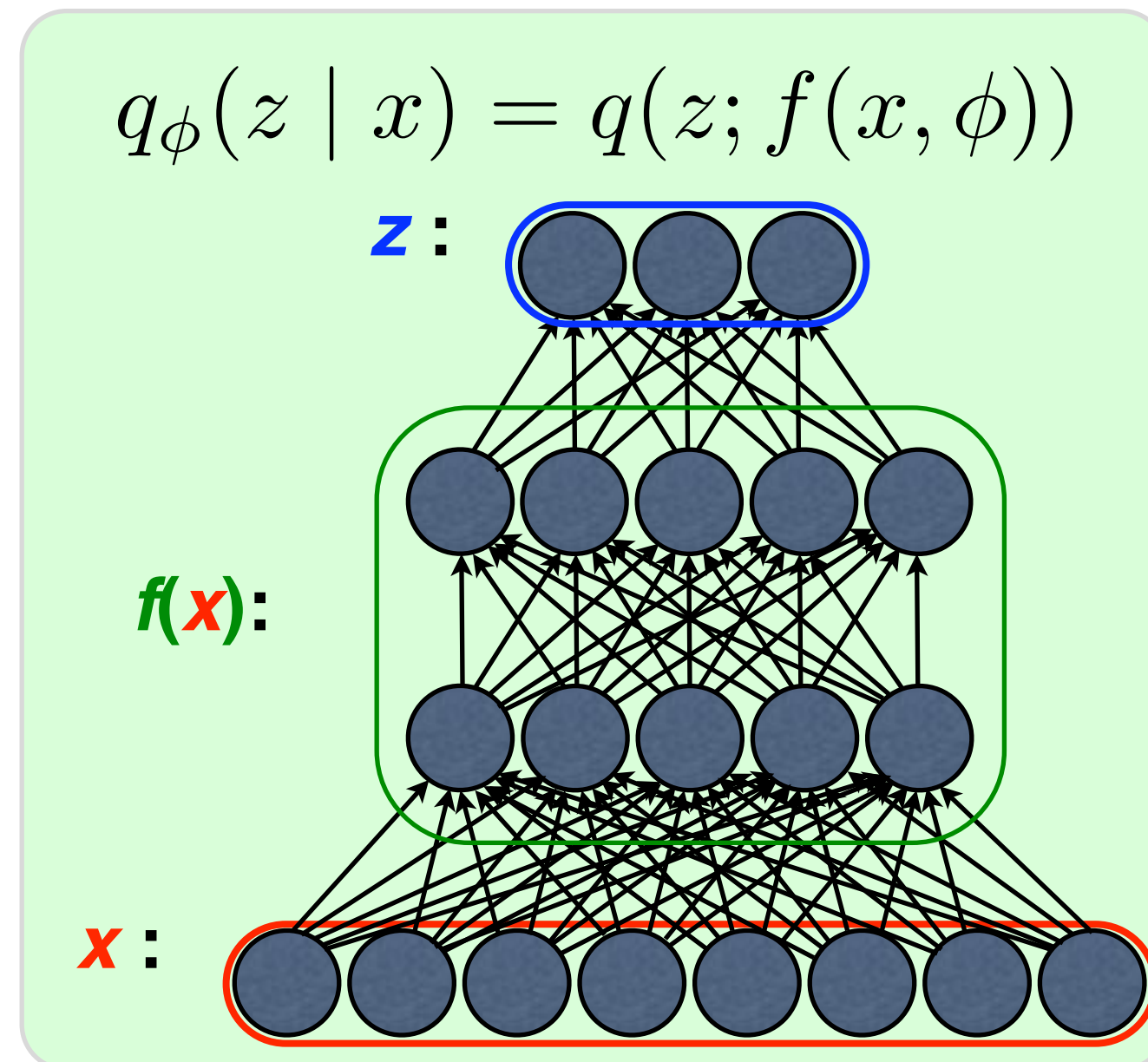
regularization term **reconstruction term**

VAE Inference model

- The VAE approach: introduce an inference model $q_\phi(z | x)$ that learns to approximate the intractable posterior $p_\theta(z | x)$ by optimizing the variational lower bound:

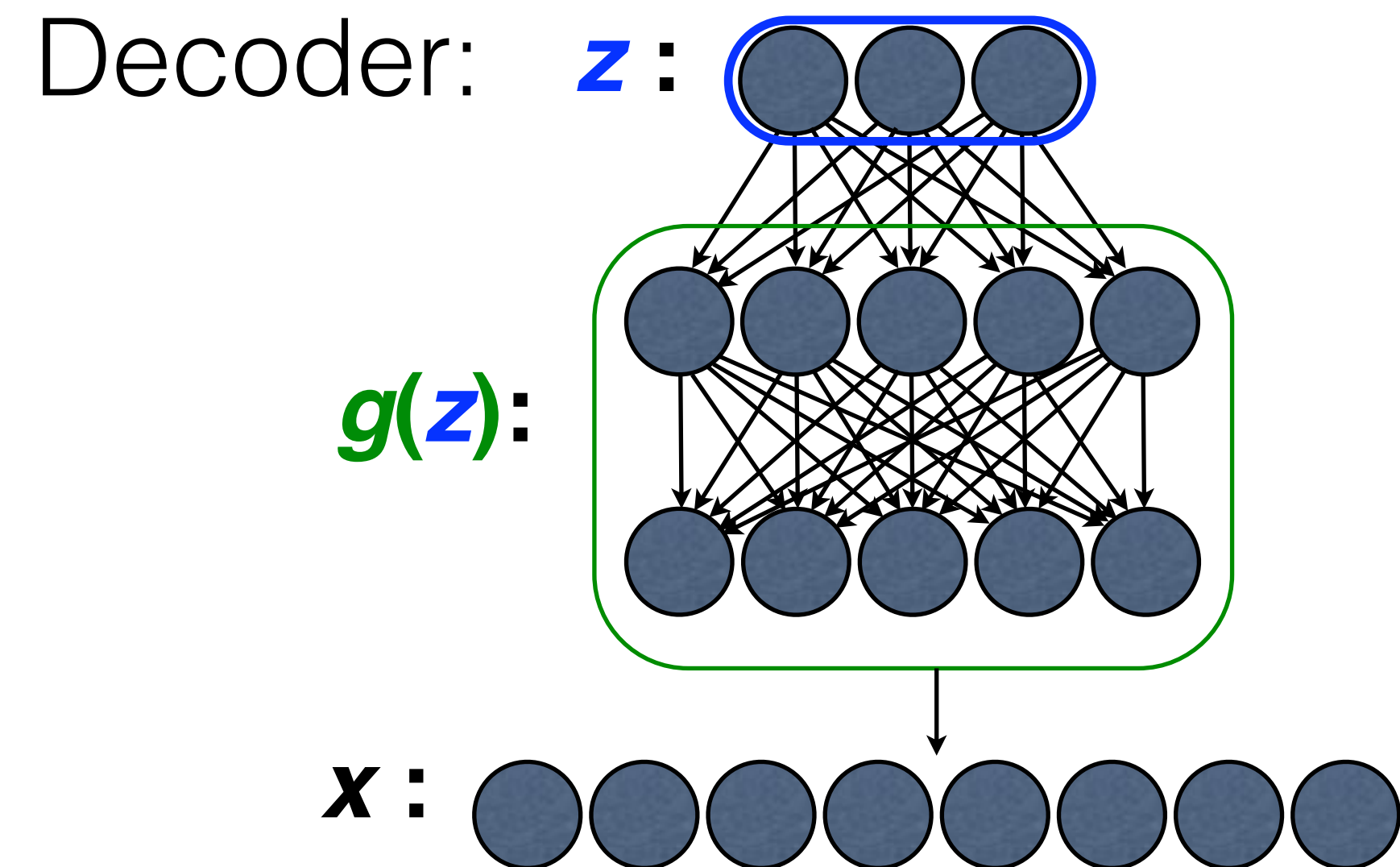
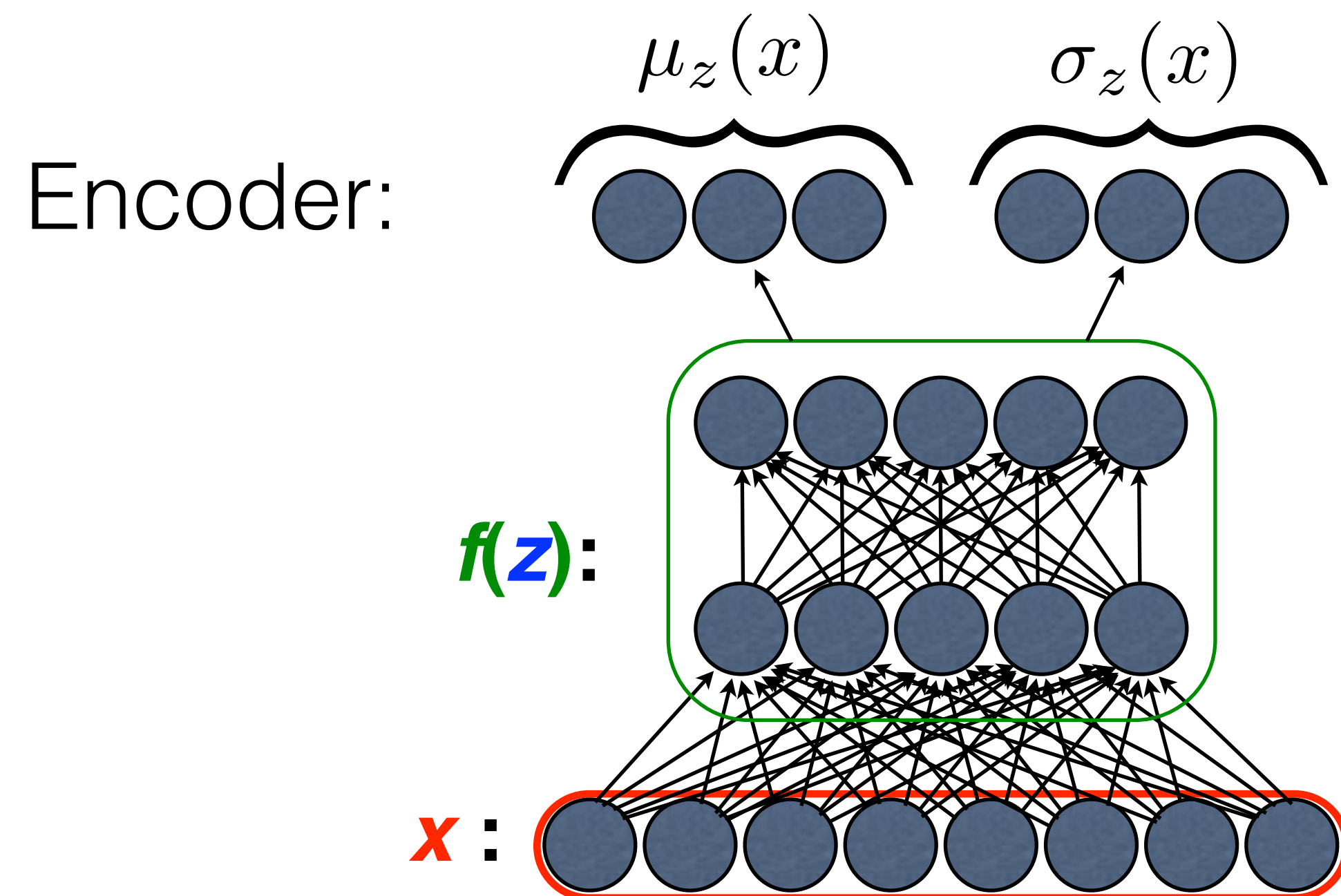
$$\mathcal{L}(\theta, \phi, x) = -D_{\text{KL}}(q_\phi(z | x) || p_\theta(z)) + \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)]$$

- We parameterize $q_\phi(z | x)$ with another neural network:



Reparametrization trick

- Adding a few details + one really important trick
- Let's consider \mathbf{z} to be real and $q_\phi(z | x) = \mathcal{N}(z; \mu_z(x), \sigma_z(x))$
- Parametrize \mathbf{z} as $z = \mu_z(x) + \sigma_z(x)\epsilon_z$ where $\epsilon_z = \mathcal{N}(0, 1)$

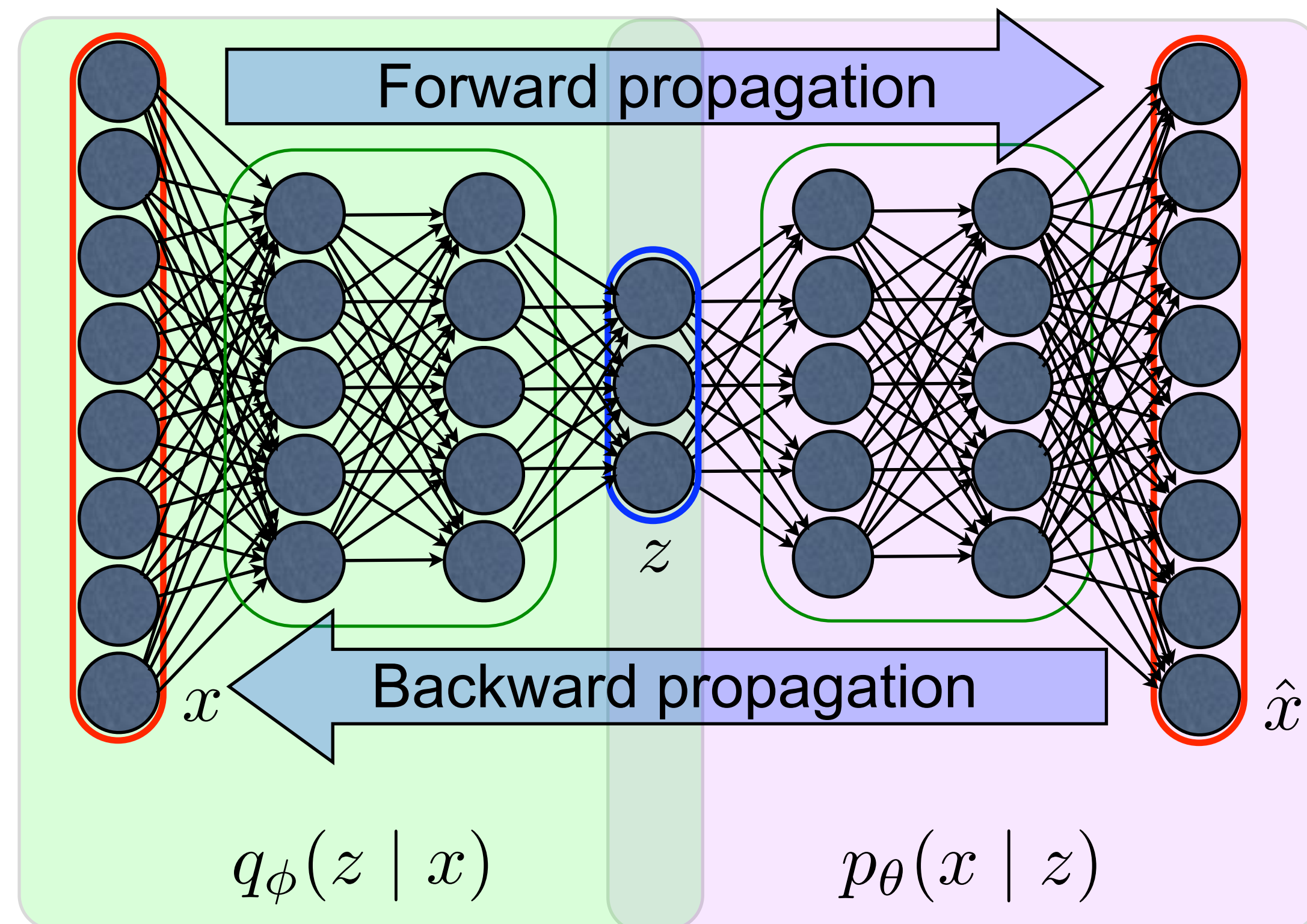


Training with backpropagation!



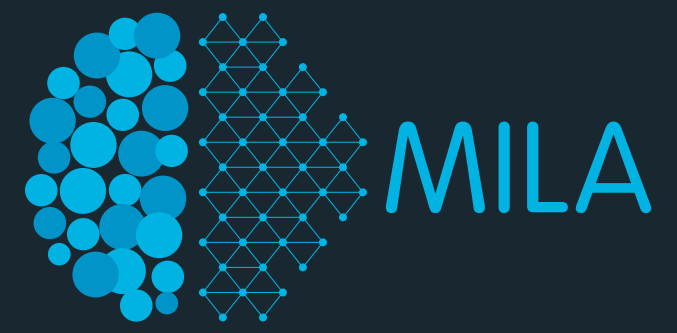
- Due to a **reparametrization** trick, we can simultaneously train both the generative model $p_{\theta}(x | z)$ and the inference model $q_{\phi}(z | x)$ by optimizing the variational bound using gradient **backpropagation**.

Objective function: $\mathcal{L}(\theta, \phi, x) = -D_{\text{KL}}(q_{\phi}(z | x) || p_{\theta}(z)) + \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x | z)]$

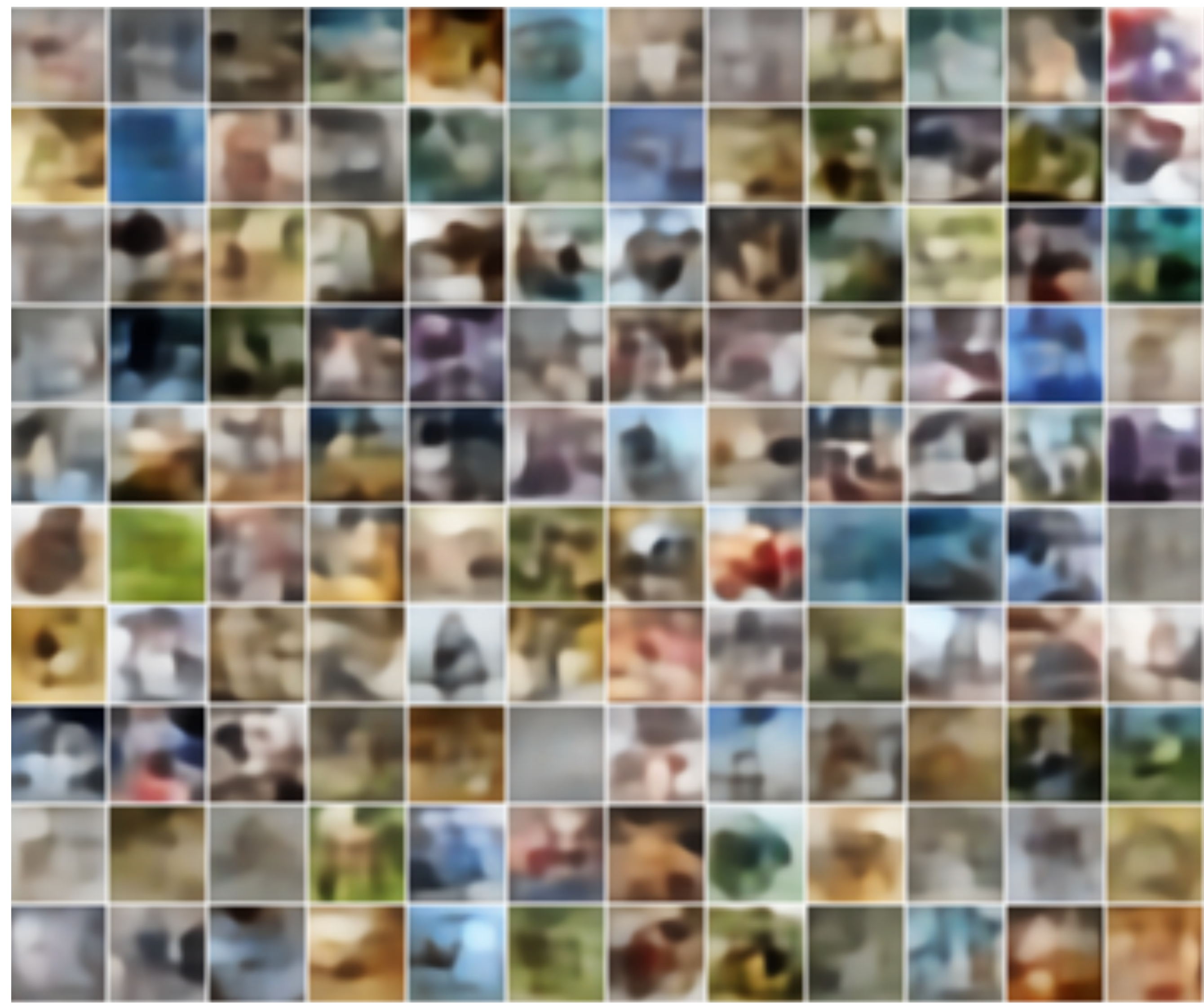


vanilla VAE samples

Impressive ...
... at the time

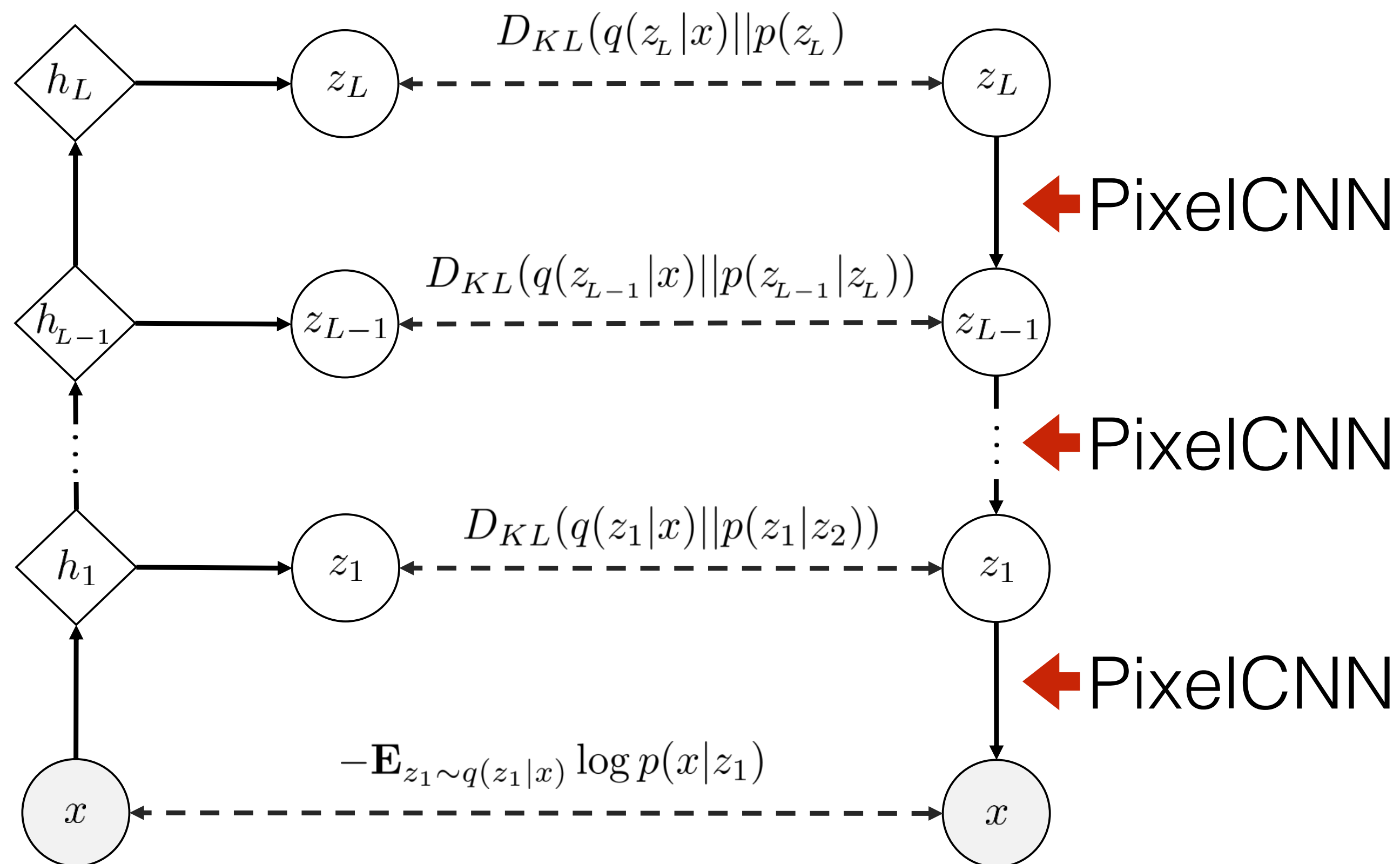


Labelled Faces in the Wild (LFW)



ImageNet (small)

- Uses a PixelCNN in the VAE decoder to help avoid the blurring caused by the standard VAE assumption of independent pixels.

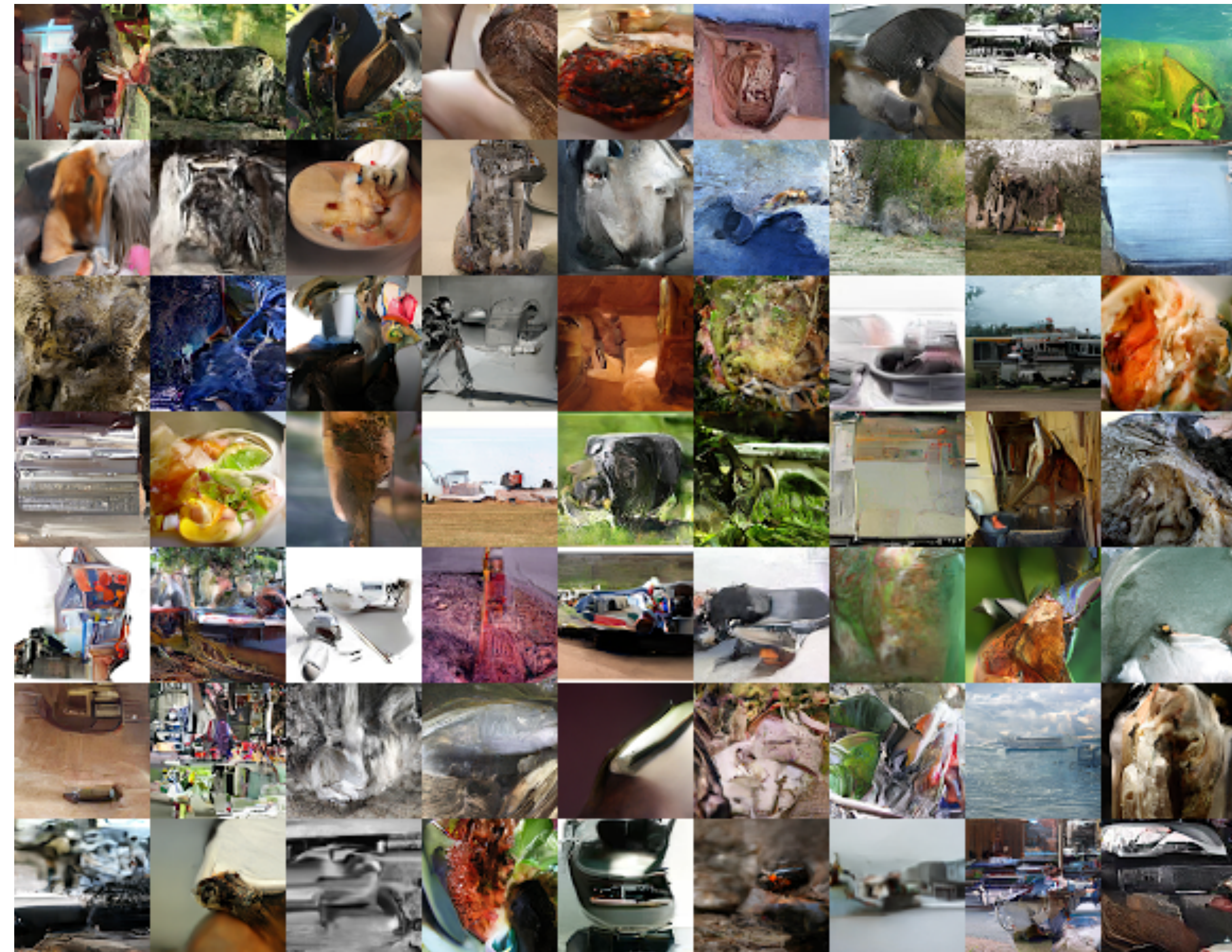


PixelVAE Samples

(Gulrajani et al. 2017)

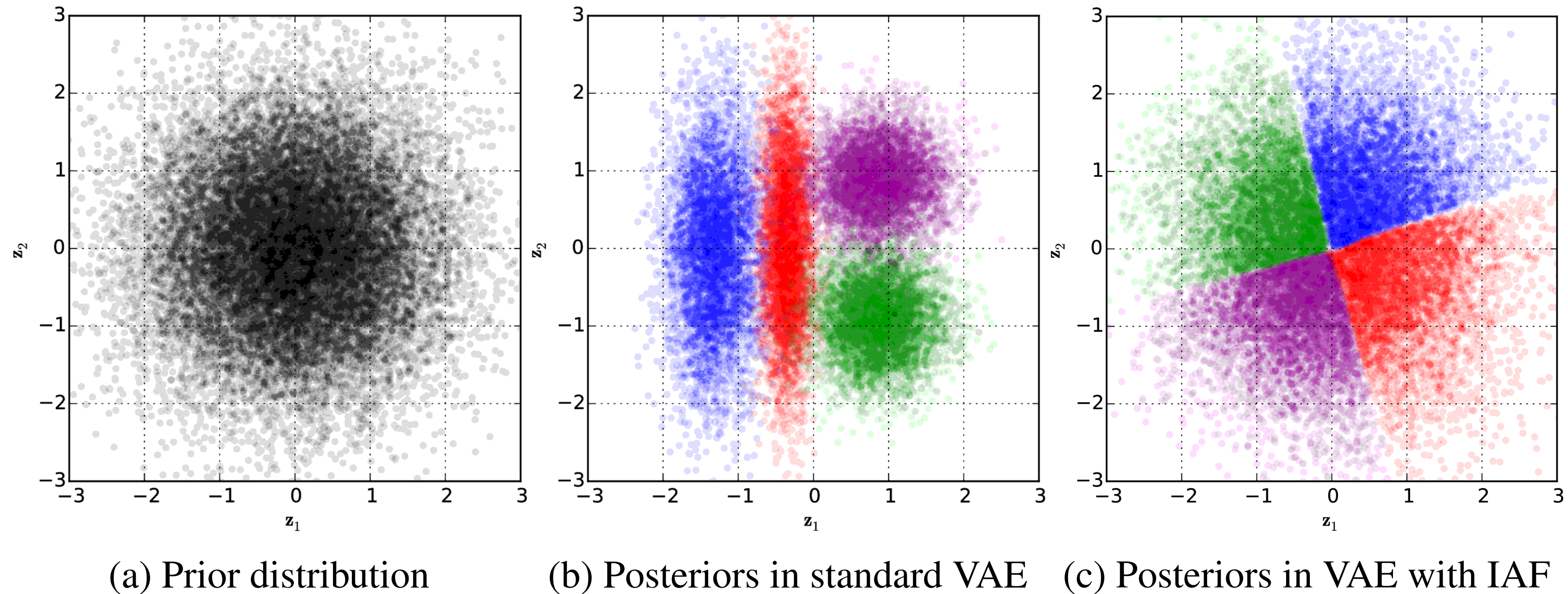


LSUN bedroom scenes (64x64)



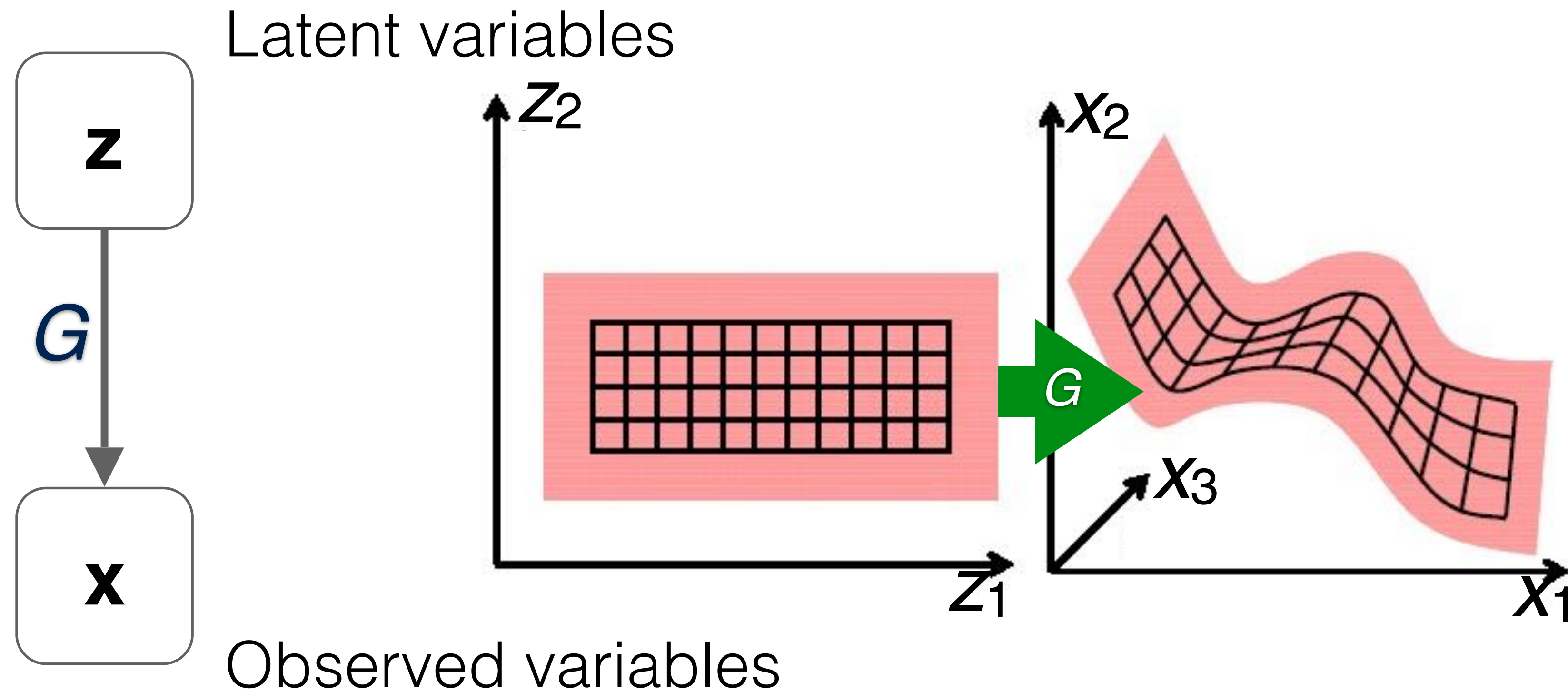
ImageNet (64x64)

Inverse Autoregressive Flow (Kingma et al., NIPS 2016)

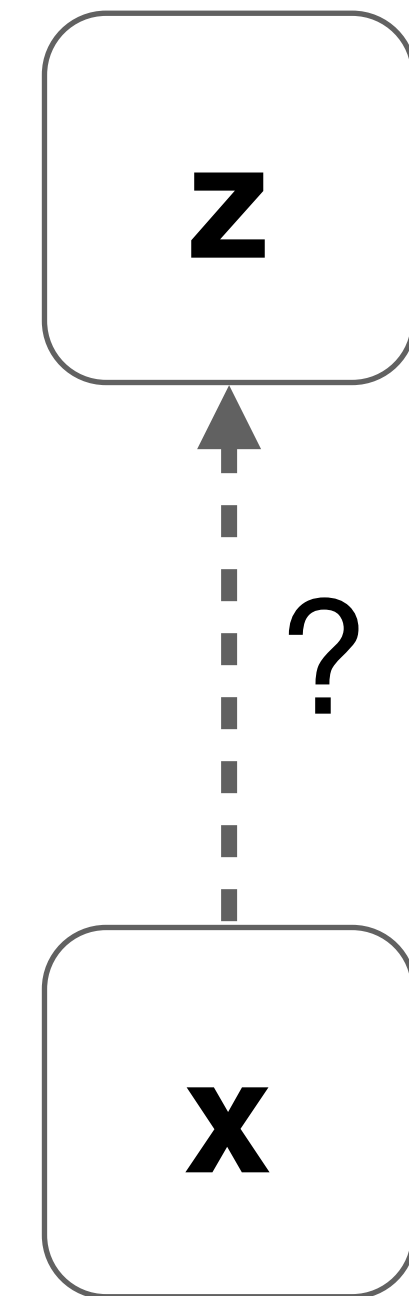


- Standard VAE posteriors are factorized - limiting how well they can (marginally) fit the prior.
- IAF greatly improves the flexibility of the posterior distributions, and allows for a much better fit between the posteriors and the prior.

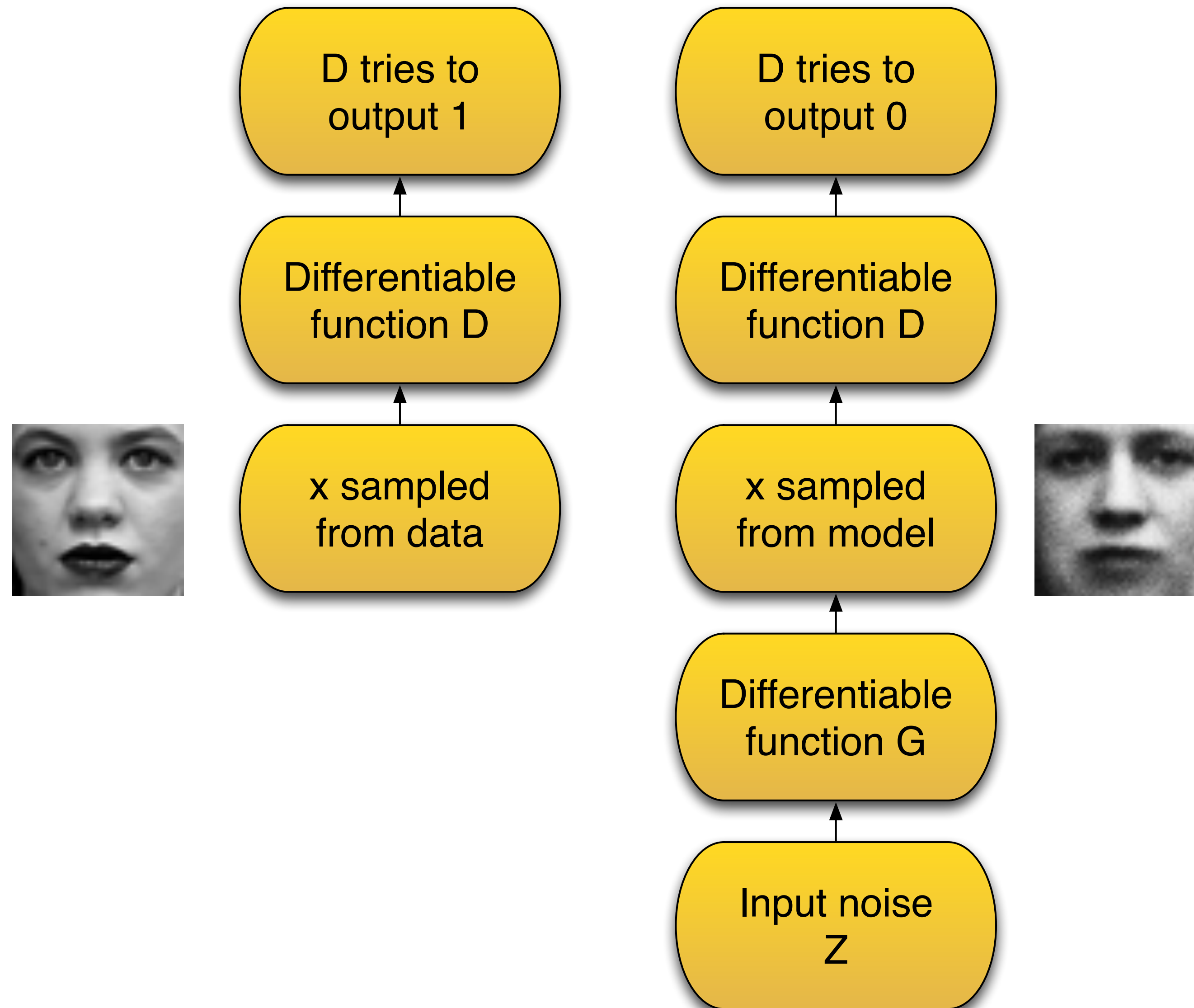
Another way to train a latent variable model?



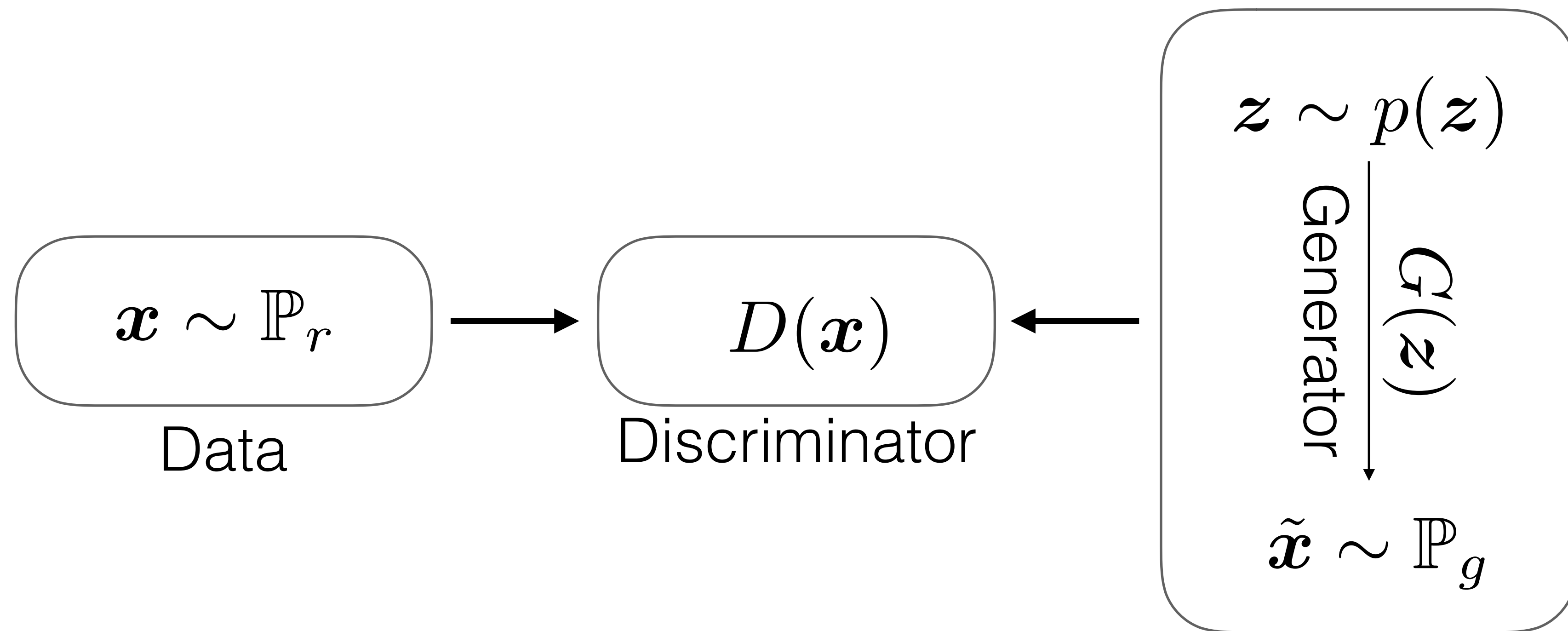
inference



Generative Adversarial Networks



Generative Adversarial Networks



GAN Objective



- Formally, express the game between discriminator D and generator G with the minimax objective:

$$\min_G \max_D \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r} [\log(D(\mathbf{x}))] + \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g} [\log(1 - D(\tilde{\mathbf{x}}))].$$

where:

- \mathbb{P}_r is the data distribution
- \mathbb{P}_g is the model distribution implicitly defined by:

$$\tilde{\mathbf{x}} = G(\mathbf{z}), \quad \mathbf{z} \sim p(\mathbf{z})$$

- the generator input \mathbf{z} is sampled from some simple noise distribution, (e.g. uniform or Gaussian).

- Optimal (nonparametric) discriminator:

$$D^*(\mathbf{x}) = \frac{p_r(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})}$$

- Under an ideal discriminator, the generator minimizes the Jensen-Shannon divergence between \mathbb{P}_r and \mathbb{P}_g .

$$\text{JS}(\mathbb{P}_r \parallel \mathbb{P}_g) = \text{KL} \left(\mathbb{P}_r \parallel \frac{\mathbb{P}_r + \mathbb{P}_g}{2} \right) + \text{KL} \left(\mathbb{P}_g \parallel \frac{\mathbb{P}_r + \mathbb{P}_g}{2} \right)$$

$$\text{where } \text{KL}(\mathbb{P}_r \parallel \mathbb{P}_g) = \int \log \left(\frac{p_r(x)}{p_g(x)} \right) p_r(x) d\mu(x)$$

GAN Theory ... in practice



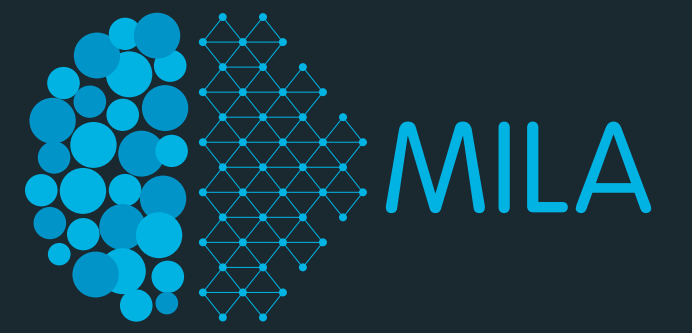
- The minimax objective leads to vanishing gradients as the discriminator saturates.
- In practice, Goodfellow et al (2014) advocate the heuristic training objective:

$$\max_D \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r} [\log(D(\mathbf{x}))] + \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g} [\log(1 - D(\tilde{\mathbf{x}}))].$$

$$\max_G \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g} [\log(D(\tilde{\mathbf{x}}))].$$

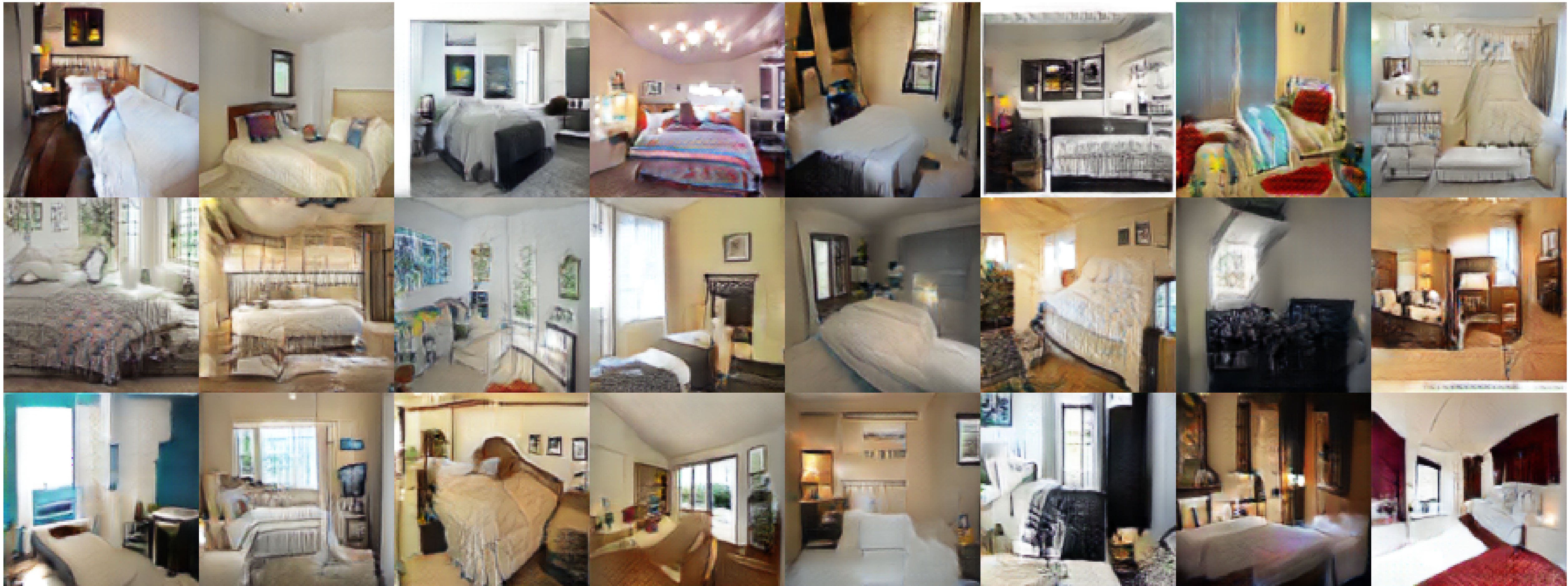
- ▶ However, this modified loss function can still misbehave in the presence of a good discriminator.

GAN samples



Least-Squares GAN

Xudong Mao, Qing Li†, Haoran Xie, Raymond Y.K. Lau and Zhen Wang, ArXiv, Feb. 2017



128x128 LSUN bedroom scenes

DCGAN samples (Radford, Metz and Chintala; 2016)



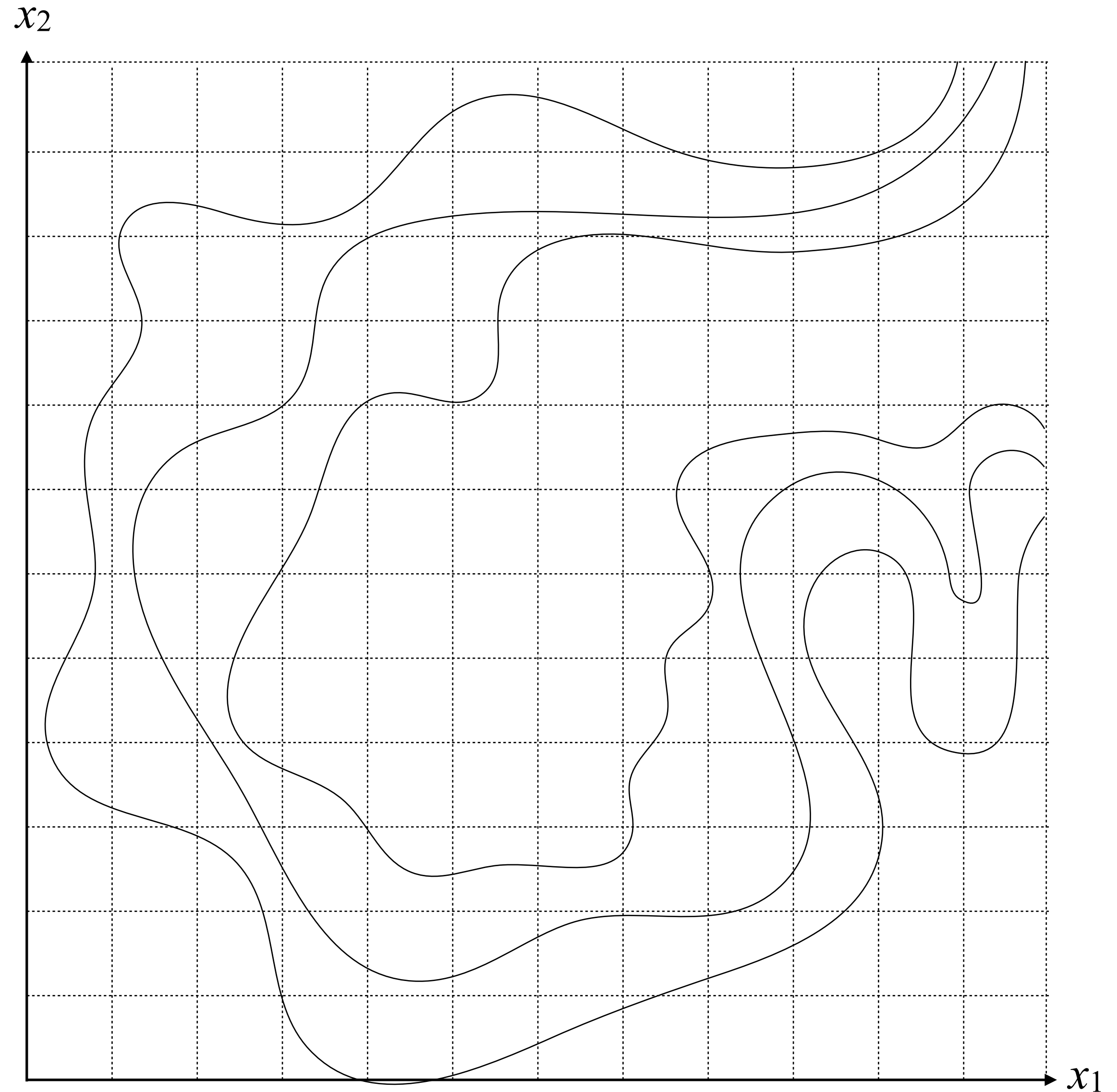
Z-space interpolations



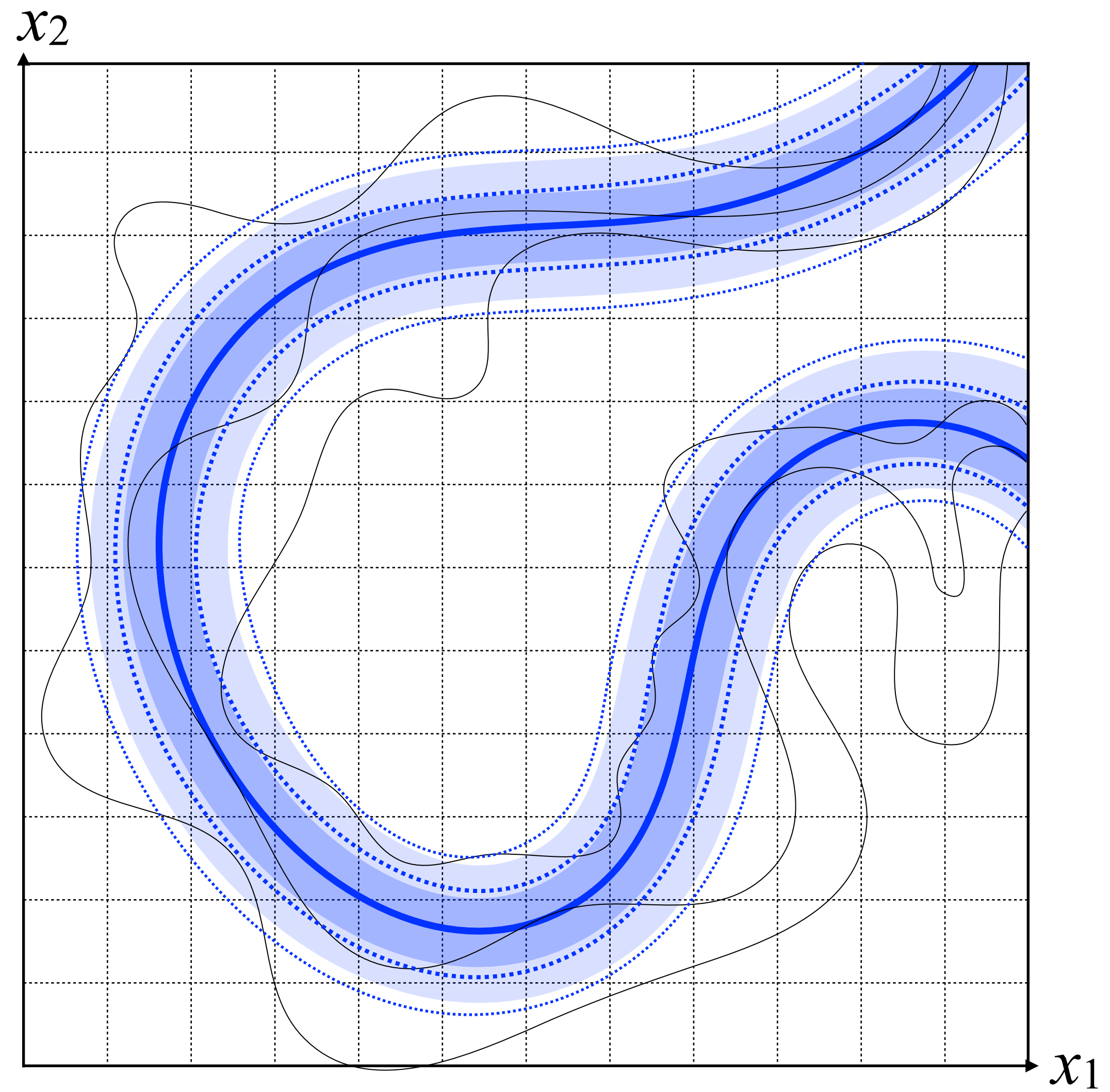
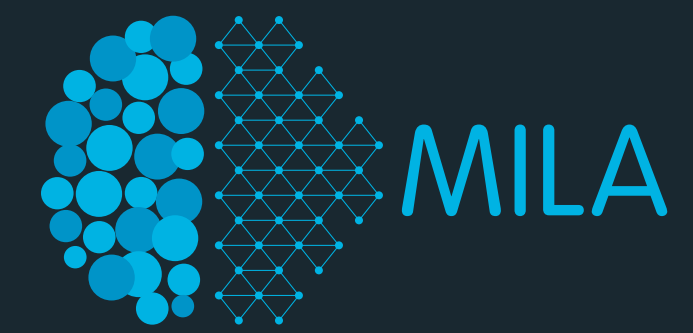
LSUN bedroom scenes

What makes GANs special?

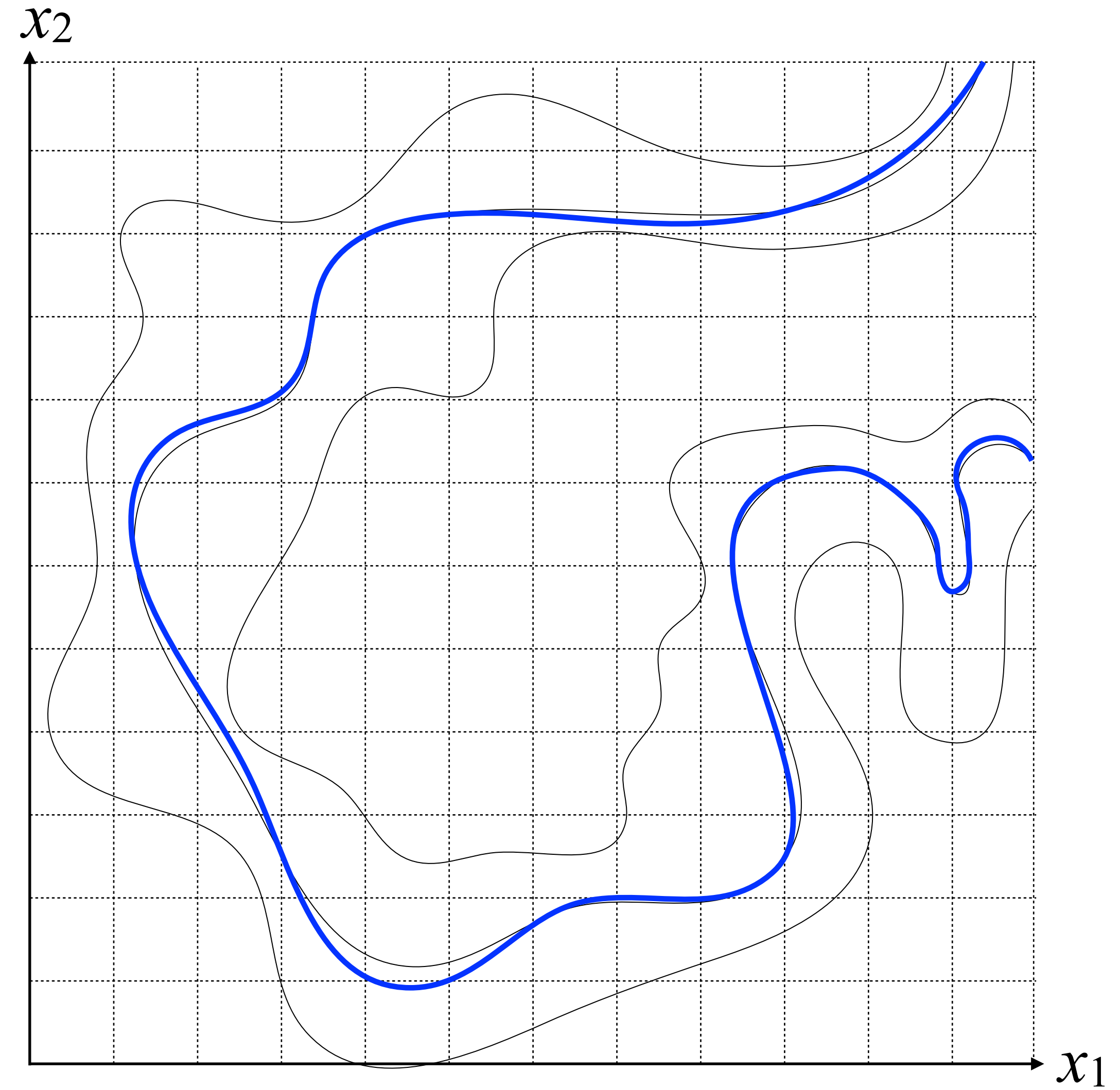
Cartoon of the Image manifold:



What makes GANs special?



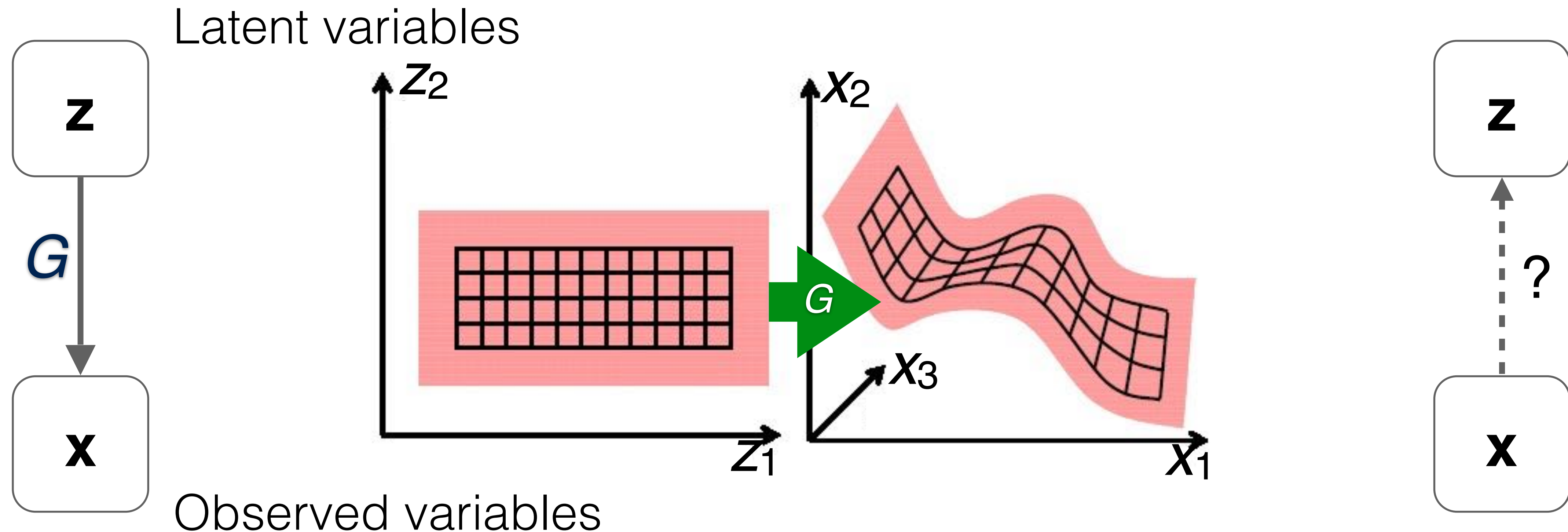
more traditional max-likelihood approach



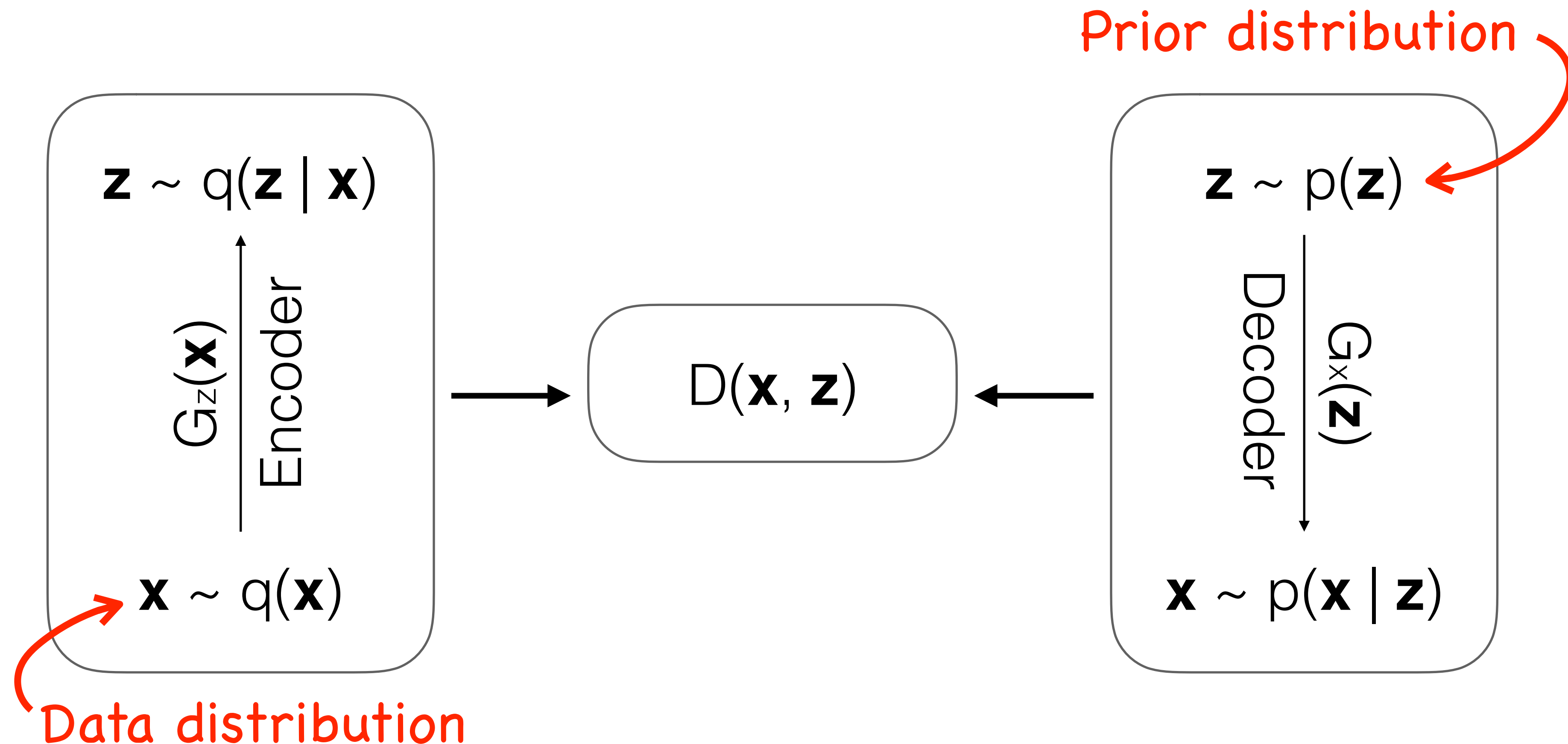
GAN

But what about inference...

- Can we incorporate an inference mechanism into GANs?



ALI / BiGAN: model diagram



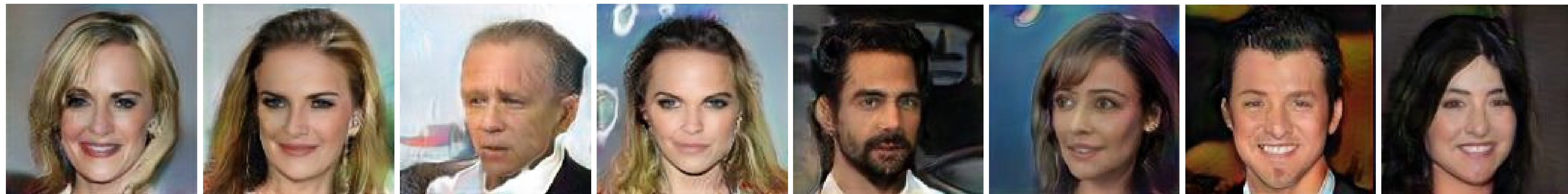
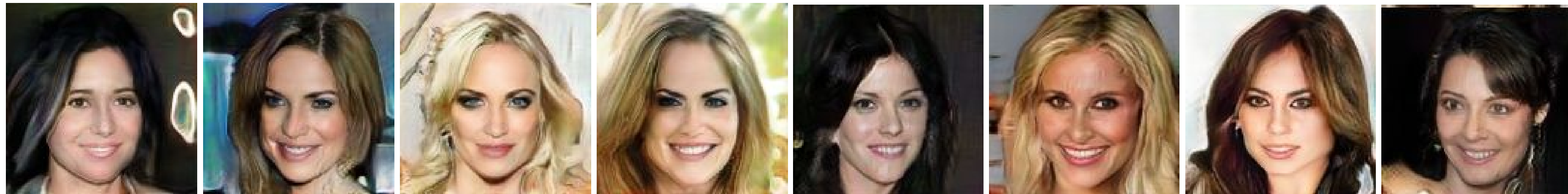
- **ALI:** Vincent Dumoulin, Ishmael Belghazi, Olivier Mastropietro, Ben Poole, Alex Lamb, Martin Arjovsky (2016) *ADVERSARIALLY LEARNED INFERENCE*, arXiv:1606.00704, ICLR 2017

- **BiGAN:** Donahue, Krähenbühl and Darrell (2016), *ADVERSARIAL FEATURE LEARNING*, arXiv:1605.09782, ICLR 2017

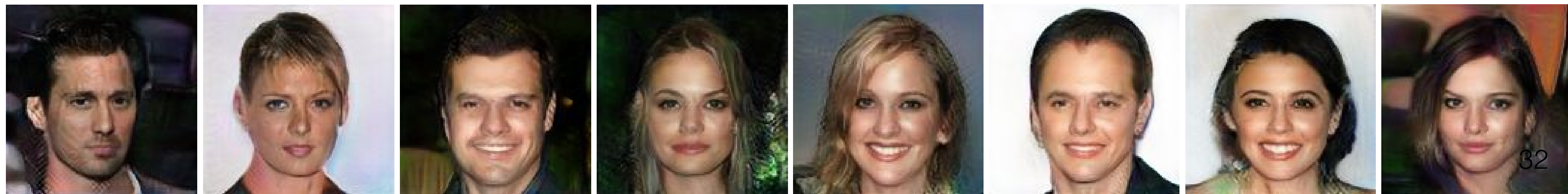
Hierarchical ALI



CelebA-128X128



Model samples



Hierarchical AFI: CelebA-128X128



Data

Recon

Reconstructions given z_1, z_2

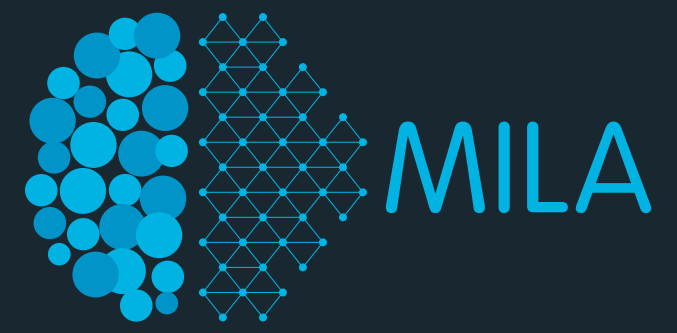
Data

Recon

Reconstructions given z_2

cycleGAN: Adversarial training of domain transformations

(Zhu et al. ICCV 2017)



- CycleGAN learns transformations across domains with unpaired data.
- Combines GAN loss with “cycle-consistency loss”: L1 reconstruction.

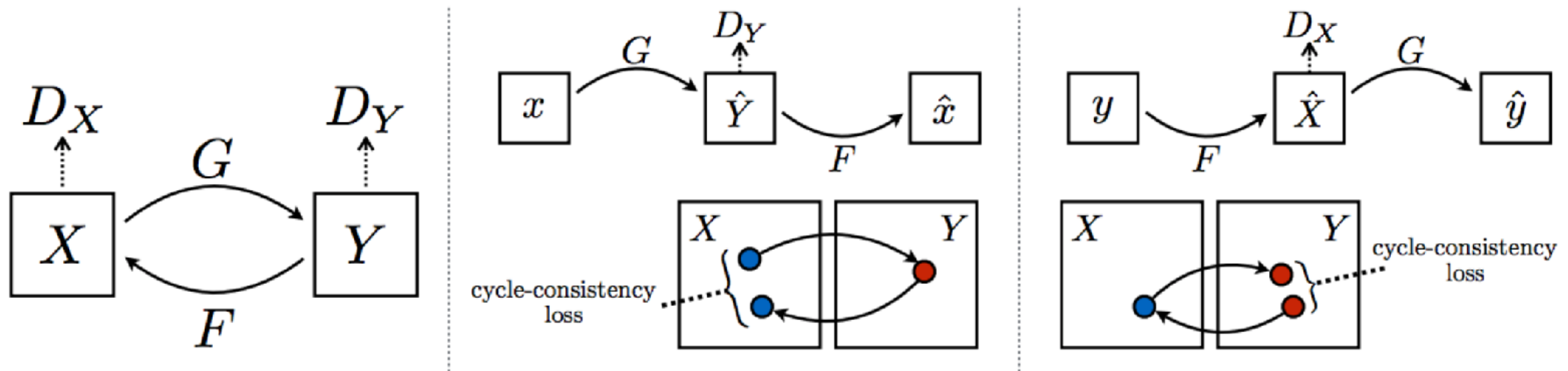


Image credits: Jun-Yan Zhu*, Taesung Park*, Phillip Isola, and Alexei A. Efros. "Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks", in IEEE International Conference on Computer Vision (ICCV), 2017.

CycleGAN for unpaired data

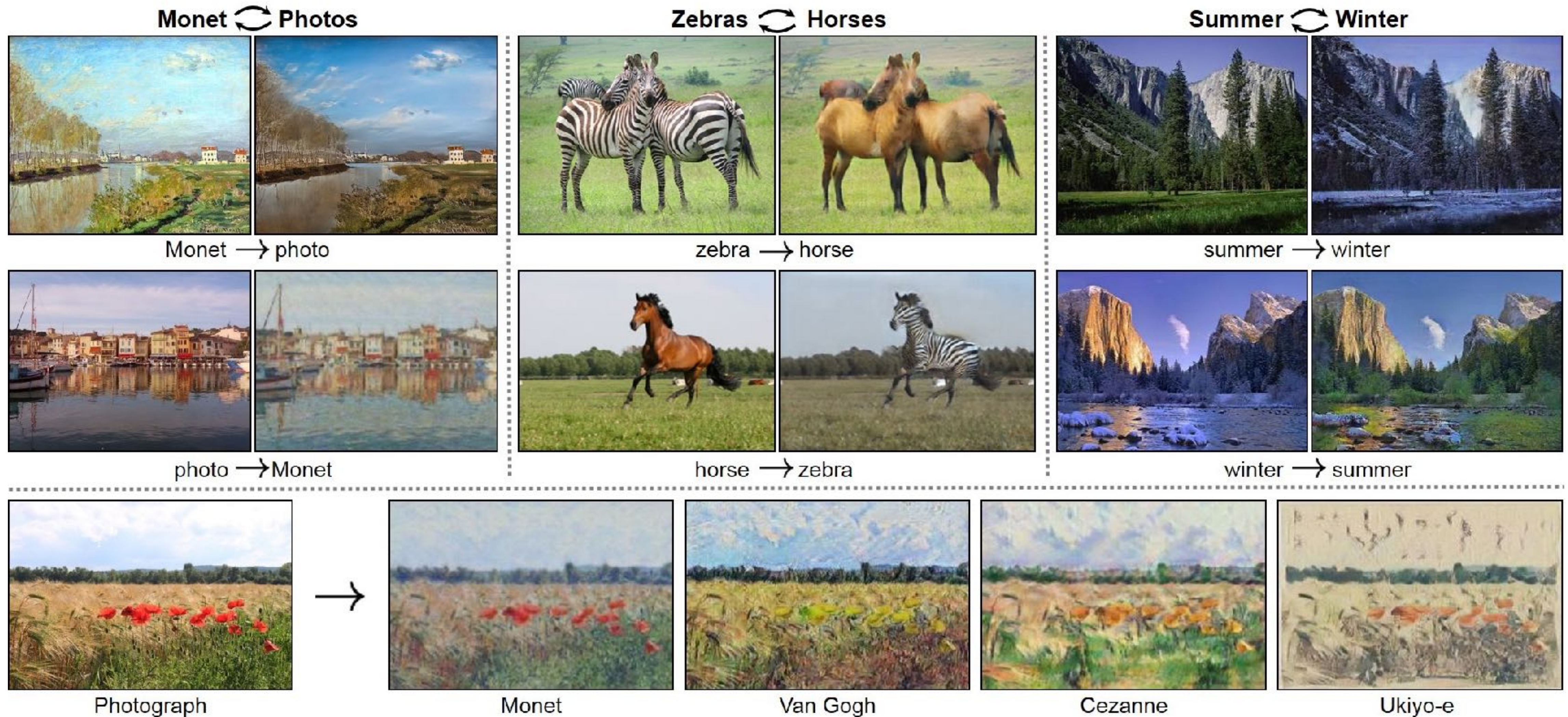
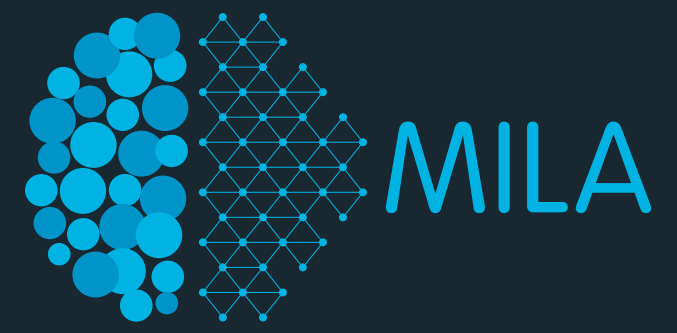


Image credits: Jun-Yan Zhu*, Taesung Park*, Phillip Isola, and Alexei A. Efros. "Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks", in IEEE International Conference on Computer Vision (ICCV), 2017.

PROGRESSIVE GROWING OF GANS FOR IMPROVED QUALITY, STABILITY, AND VARIATION (Kerras et al. from NVIDIA, 2017)



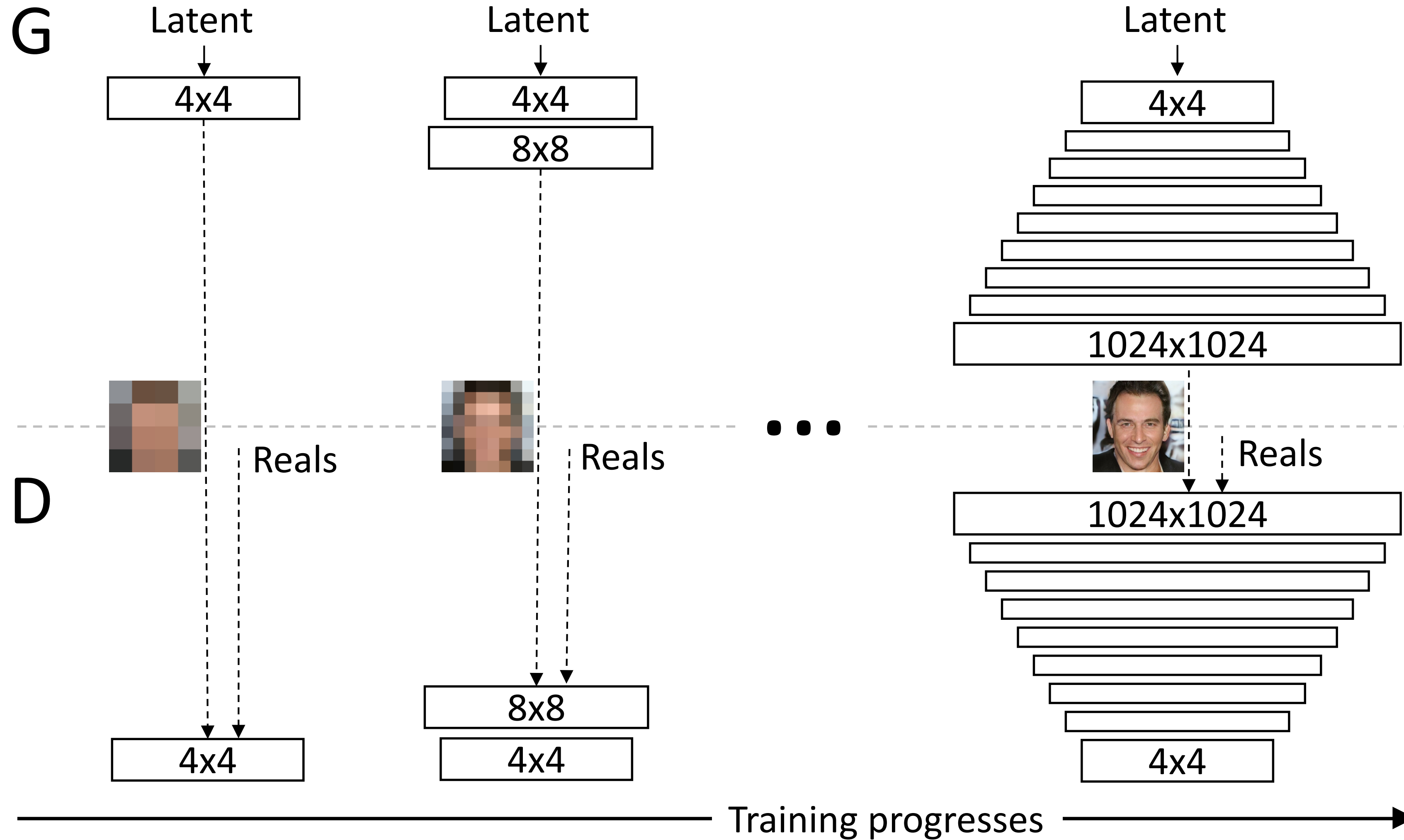
- Recent work from NVIDIA.
- Improves image quality by growing the model size throughout training.
- Samples from a model trained on the CelebA face dataset.



1024x1024 model samples

PROGRESSIVE GROWING OF GANS FOR IMPROVED QUALITY, STABILITY, AND VARIATION

(Kerras et al. from NVIDIA, 2017)



PROGRESSIVE GROWING OF GANS FOR IMPROVED QUALITY, STABILITY, AND VARIATION (Kerras et al. from NVIDIA, 2017)



- Recent work from NVIDIA.
- Improves image quality by growing the model size throughout training.
- Conditional samples from a model trained on the LSUN dataset



POTTEDPLANT

HORSE

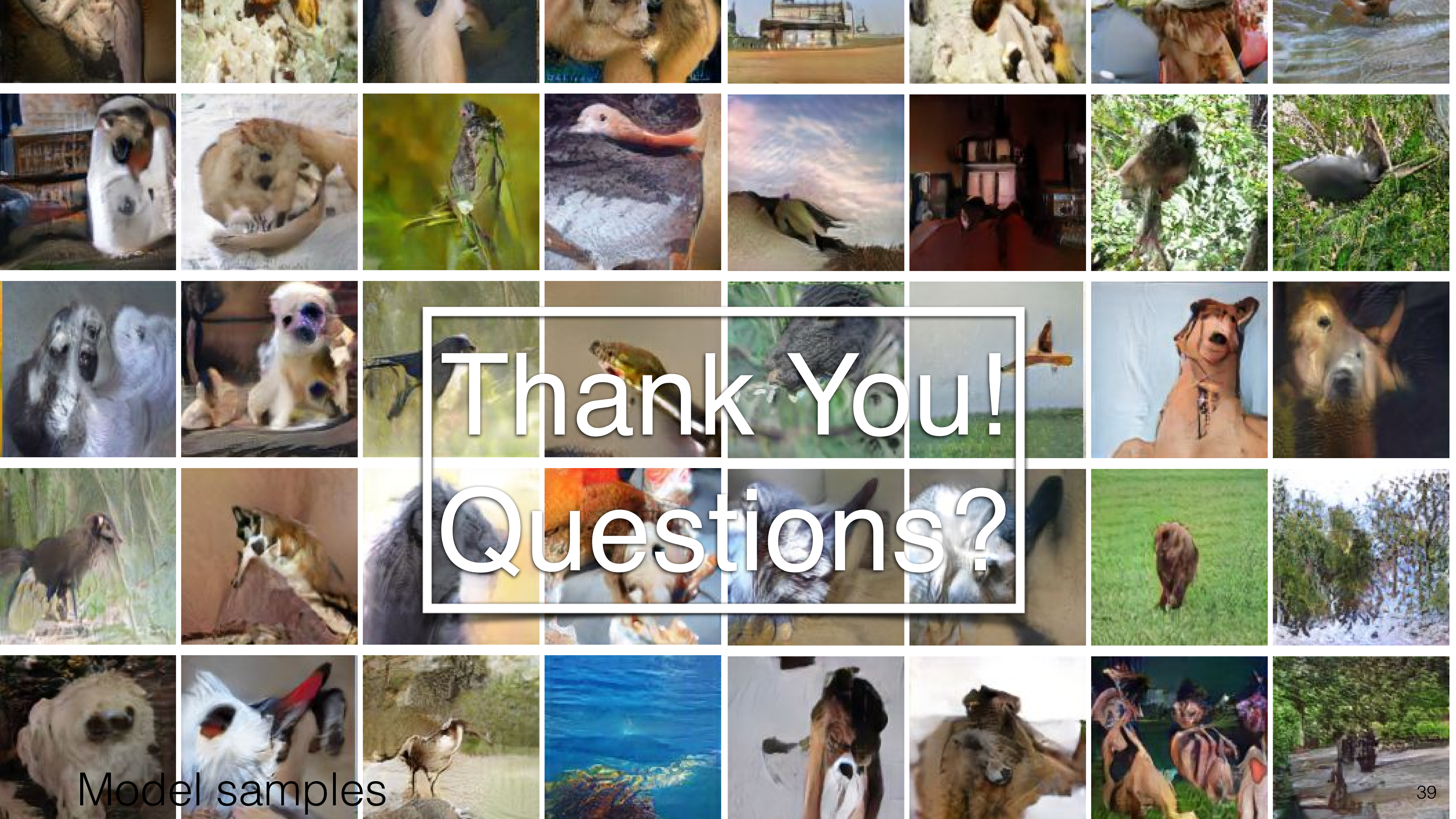
SOFA

BUS

CHURCHOUTDOOR

BICYCLE

TVMONITOR



Thank You!

Questions?

Model samples