

1

Generative Model Zoo (part II)

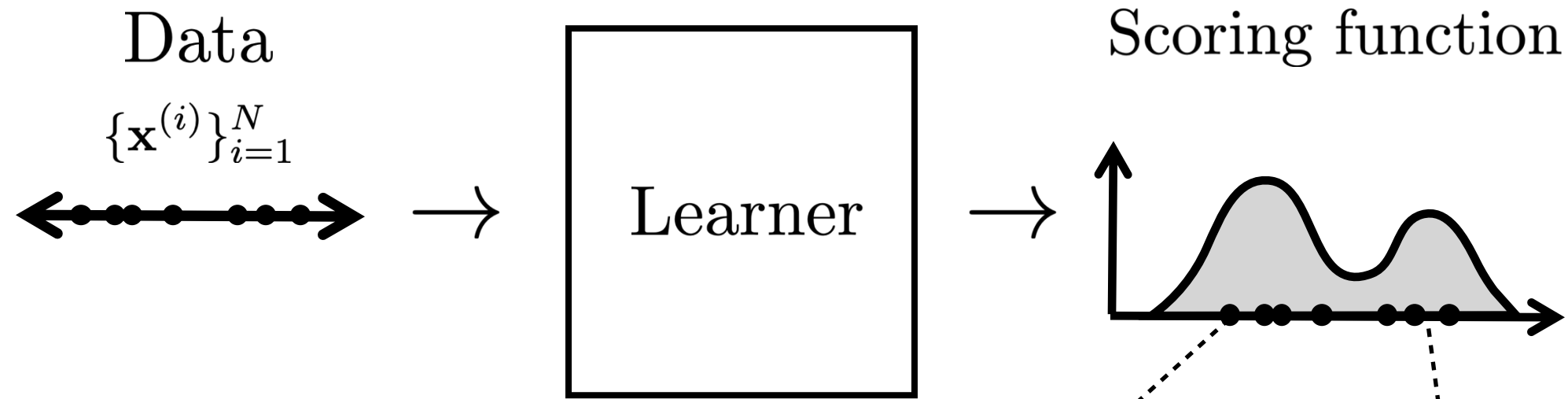
Jun-Yan Zhu

16-726 Learning-based Image Synthesis, Spring 2025

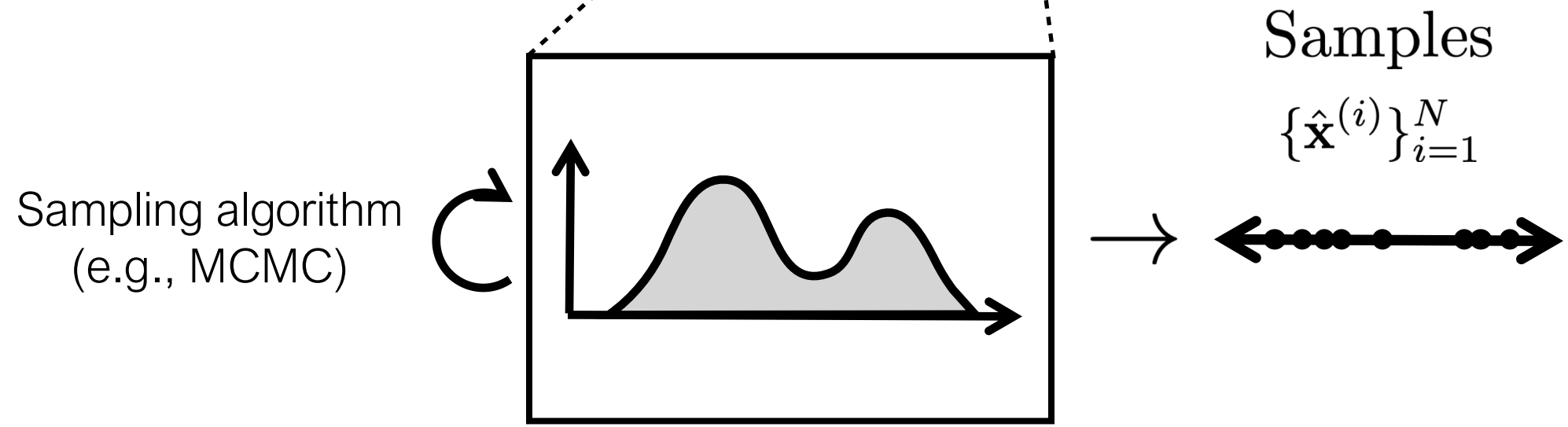
Explicit Density Models

e.g., likelihood, energy, "score function"

Training

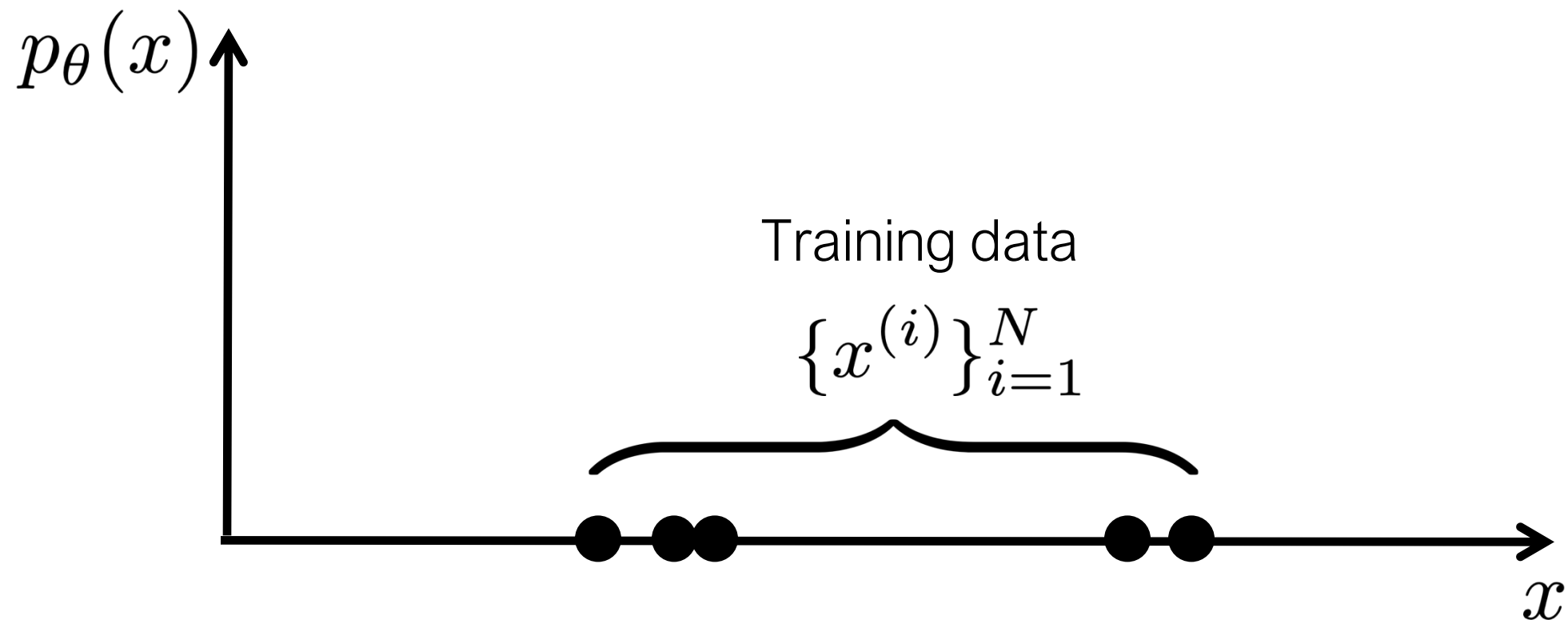


Sampling



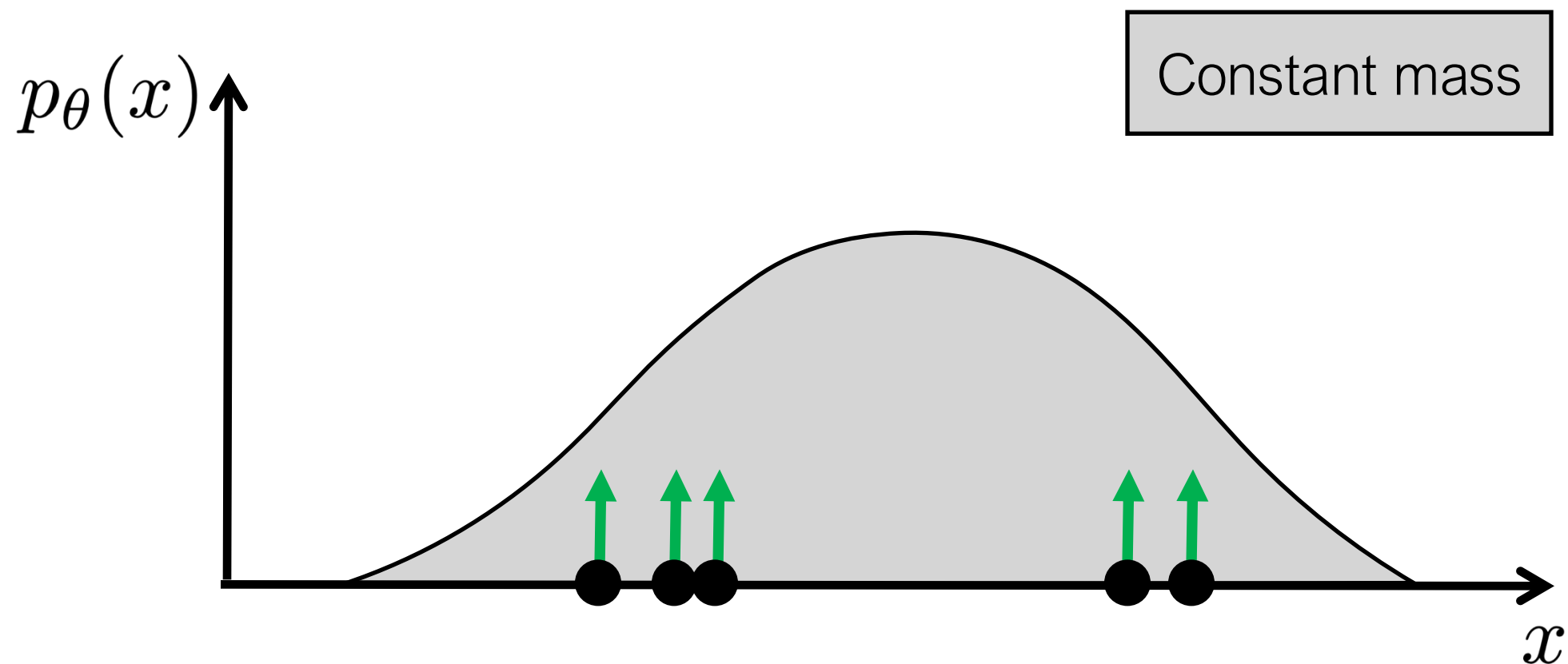
Explicit Density Models

$$p_{\theta} : \mathcal{X} \rightarrow [0, \infty)$$



Explicit Density Models

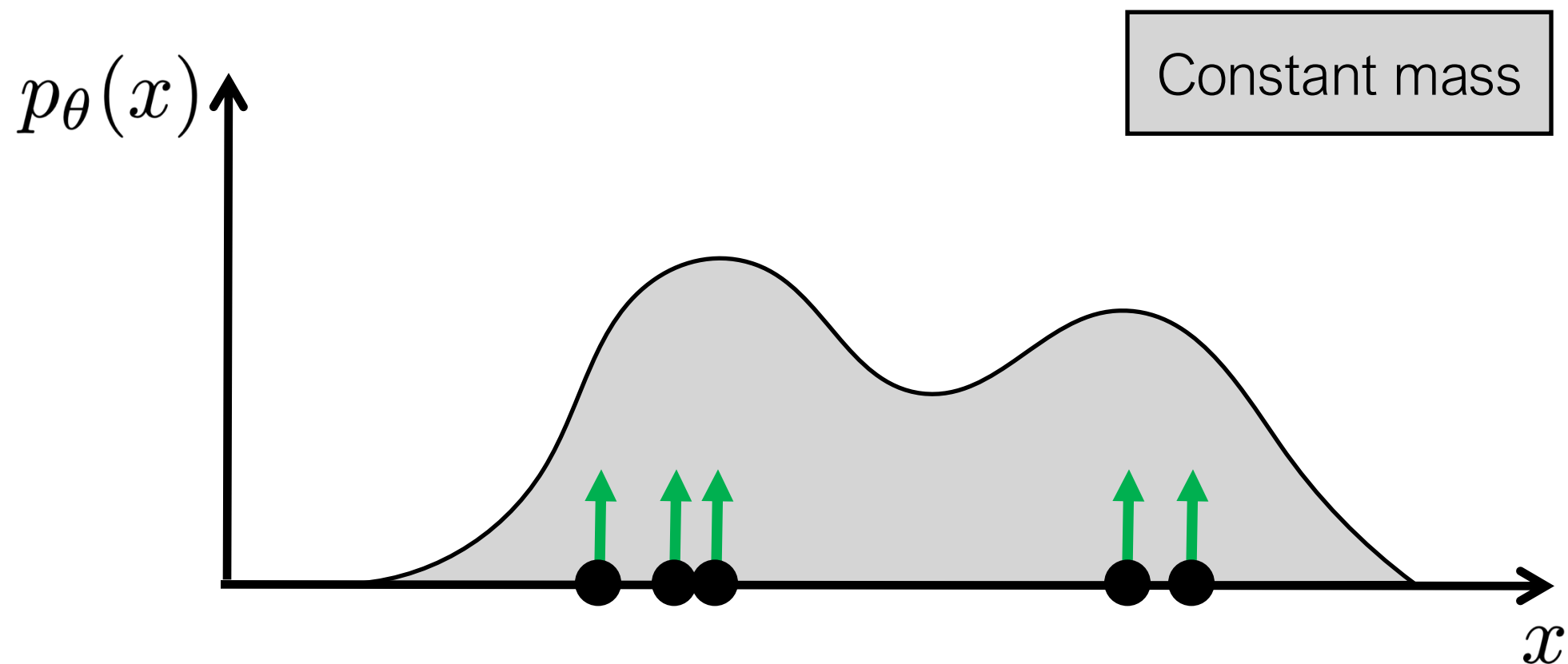
$$p_\theta : \mathcal{X} \rightarrow [0, \infty)$$



$$\int_x p_\theta(x) dx = 1$$

Explicit Density Models

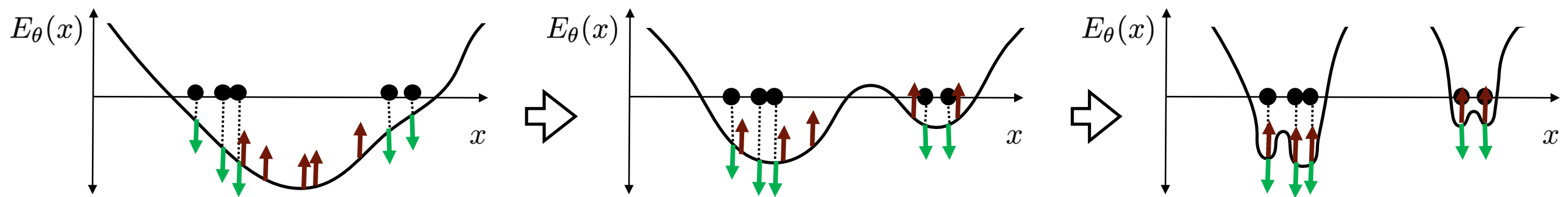
$$p_\theta : \mathcal{X} \rightarrow [0, \infty)$$



$$\int_x p_\theta(x) dx = 1$$

Energy-based models

At convergence, green (data) and red (model) samples are identical and model update (green-red) cancels out



Energy-based models — learning algorithm

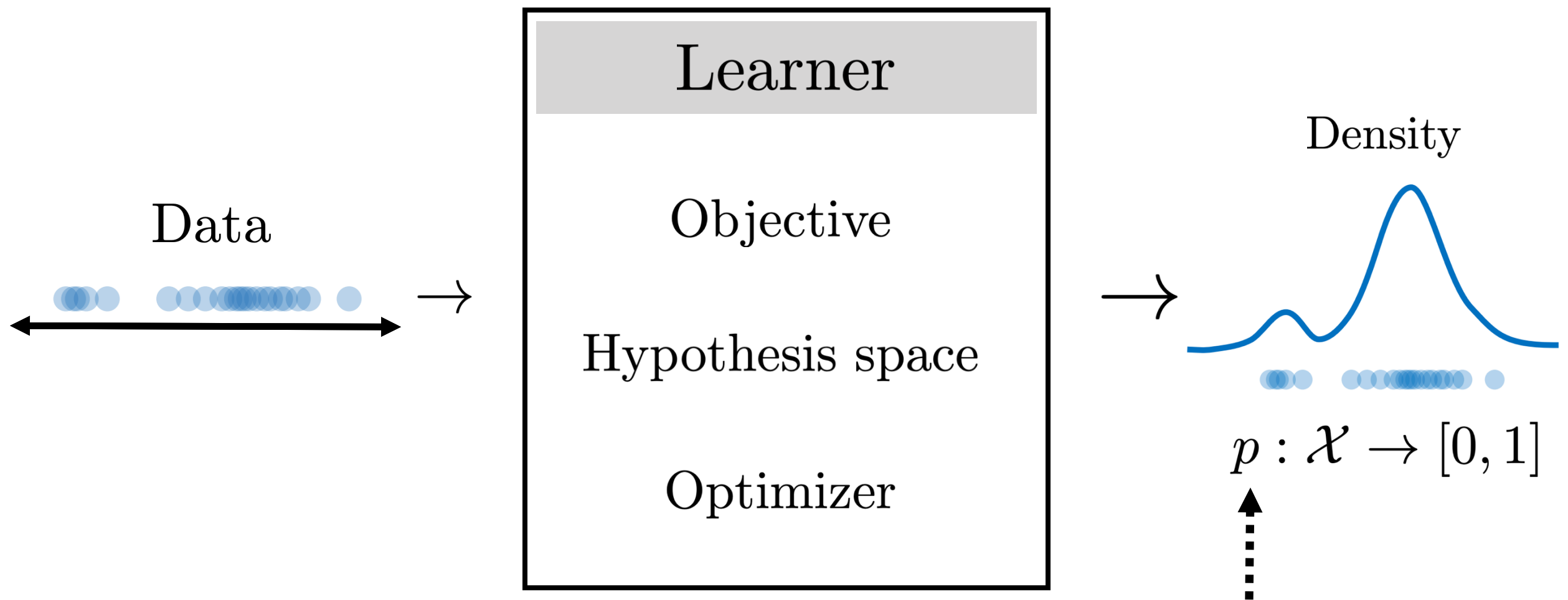
$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log p_{\theta}(\mathbf{x})] &= \nabla_{\theta} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[\log \frac{e^{-E_{\theta}(\mathbf{x})}}{Z(\theta)} \right] \\ &= \boxed{-\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\nabla_{\theta} E_{\theta}(\mathbf{x})]} - \underbrace{\boxed{\nabla_{\theta} \log Z(\theta)}}\end{aligned}$$

How to measure this?

Energy-based models — learning algorithm

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log p_{\theta}(\mathbf{x})] &= \nabla_{\theta} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[\log \frac{e^{-E_{\theta}(\mathbf{x})}}{Z(\theta)} \right] \\ &= \boxed{-\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\nabla_{\theta} E_{\theta}(\mathbf{x})]} - \boxed{\nabla_{\theta} \log Z(\theta)} \\ &= \boxed{-\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\nabla_{\theta} E_{\theta}(\mathbf{x})]} + \boxed{\mathbb{E}_{\mathbf{x} \sim p_{\theta}} [\nabla_{\theta} E_{\theta}(\mathbf{x})]} \\ &\approx \boxed{-\frac{1}{N} \sum_{i=1}^N \nabla_{\theta} E_{\theta}(\mathbf{x}^{(i)})} + \boxed{\frac{1}{N} \sum_{i=1}^N \nabla_{\theta} E_{\theta}(\hat{\mathbf{x}}^{(i)})} \\ &\quad \mathbf{x}^{(i)} \sim p_{\text{data}} \qquad \hat{\mathbf{x}}^{(i)} \sim p_{\theta}\end{aligned}$$

Learning an Explicit Density Model



Integral of probability density function needs to be 1 \longrightarrow Normalized distribution
(some models output unnormalized *energy functions*)

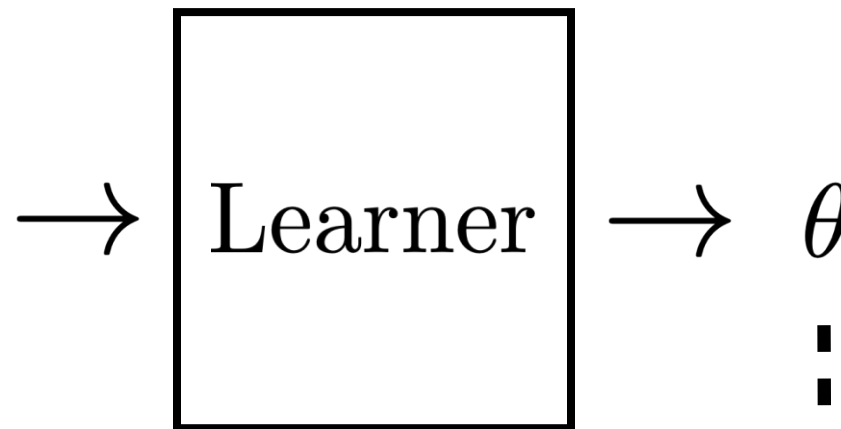
[figs modified from: http://introtodeeplearning.com/materials/2019_6S191_L4.pdf]

Useful for abnormality/outlier detection (detect unlikely events)

Implicit Generative Models

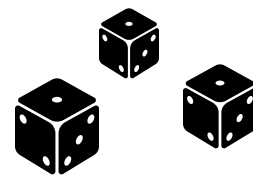
Data

Training

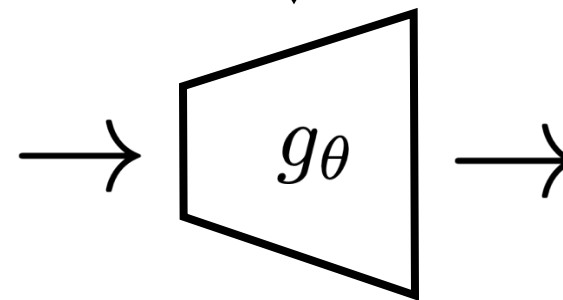


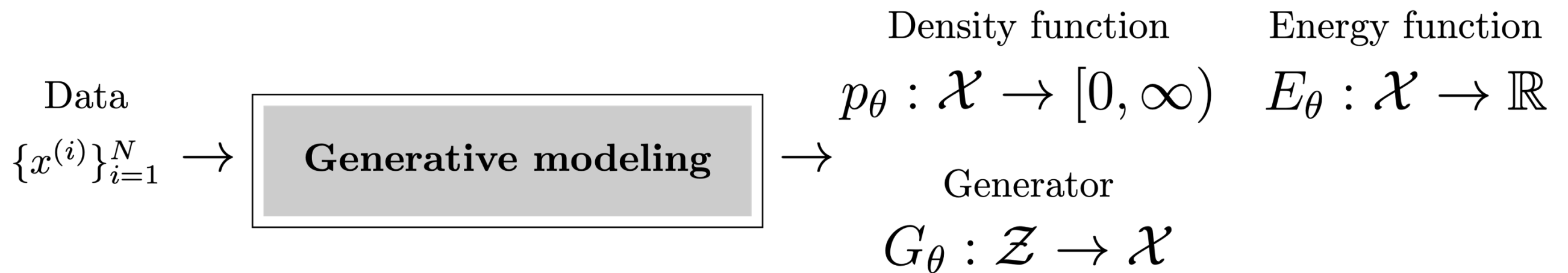
Samples

Sampling



z





You can represent the data generating process directly or indirectly

Case study #1: Fitting a Gaussian to data

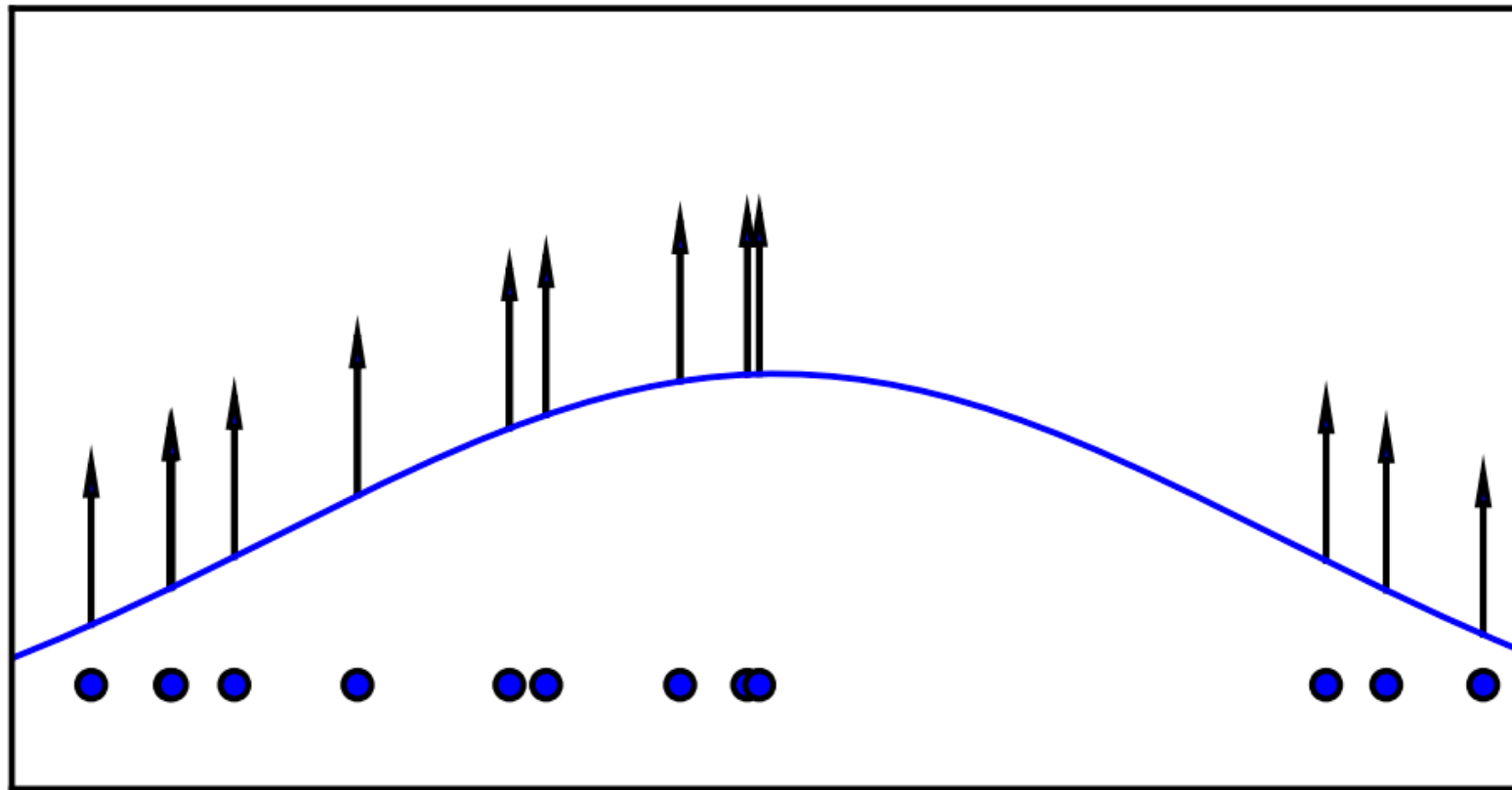


fig from [Goodfellow, 2016]

Max likelihood objective

$$\max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} [\log p_{\theta}(x)]$$

Considering only Gaussian fits

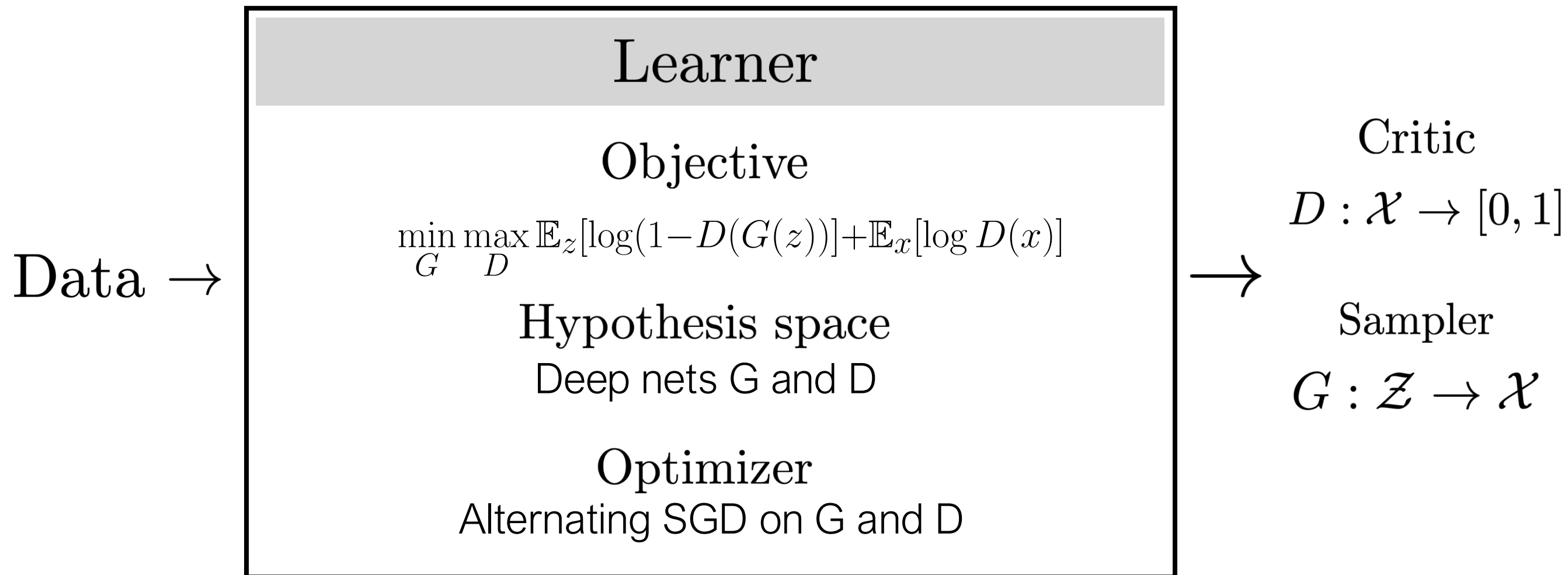
$$p_{\theta}(x) = \mathcal{N}(x; \mu, \sigma)$$

$$\theta = [\mu, \sigma]$$

Closed form optimum:

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i \quad \sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

Case study #2: Generative Adversarial Network



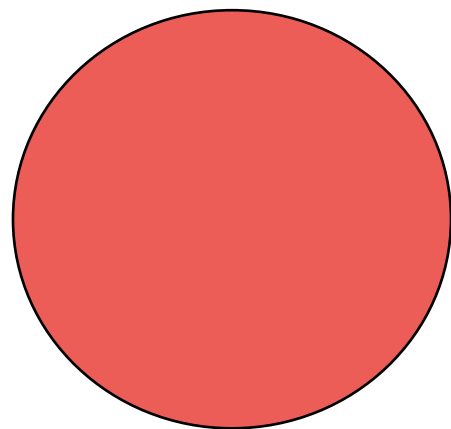
Variational Autoencoder (VAE)

Variational Autoencoders (VAEs)

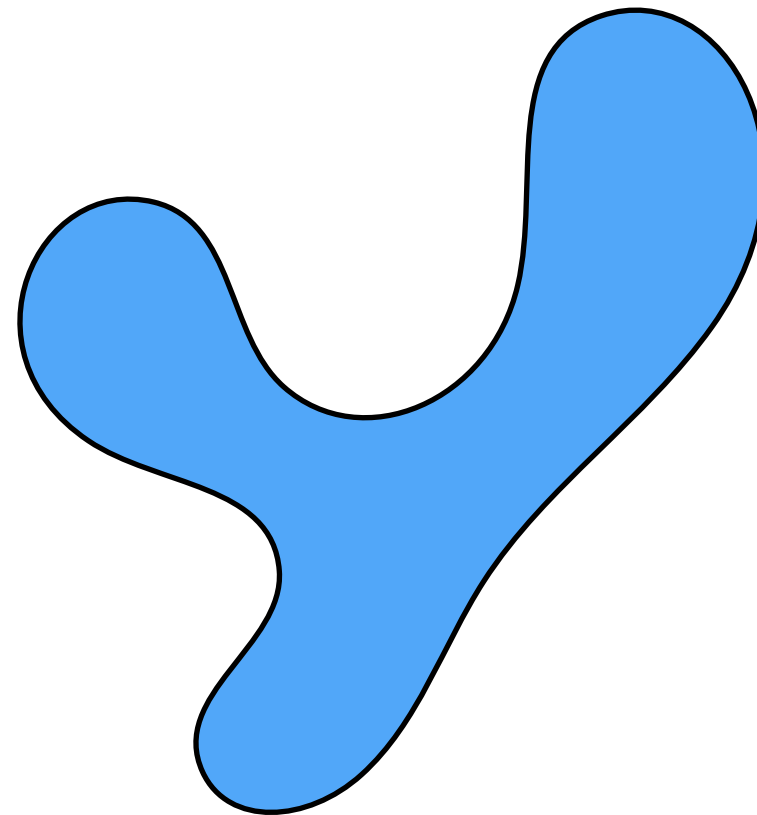
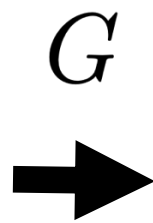
[Kingma & Welling, 2014; Rezende, Mohamed, Wierstra 2014]

Prior distribution

Target distribution

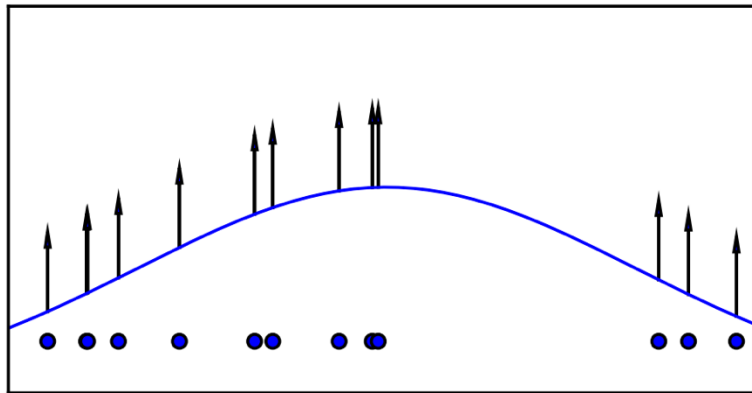


$p(z)$



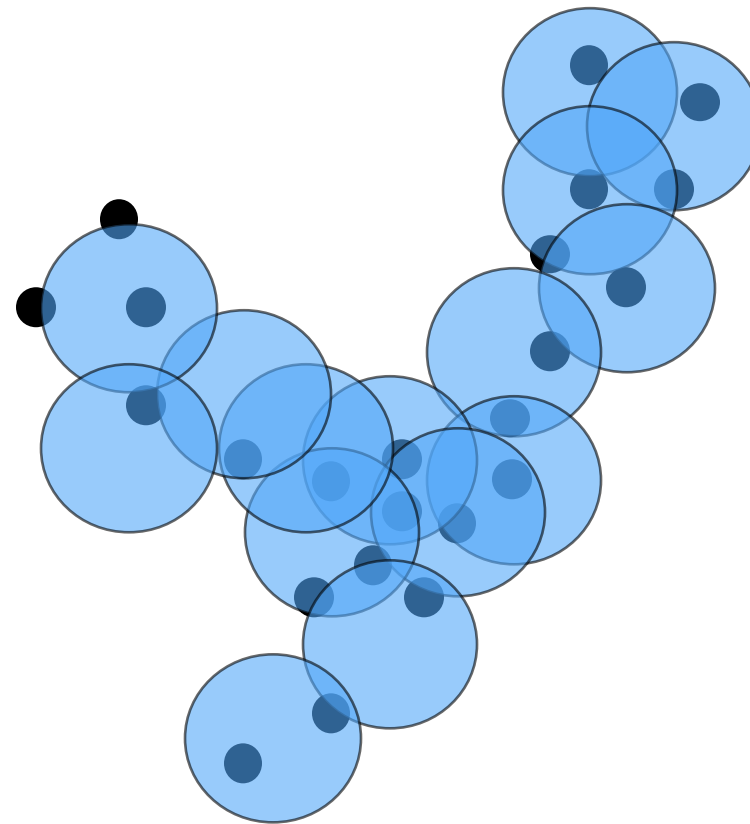
$p(x)$

Mixture of Gaussians



$$p_{\theta}(x) = \sum_{i=1}^k w_i \mathcal{N}(x; \mu_i, \Sigma_i)$$

Target distribution



$x \sim p_{\text{data}}(x)$

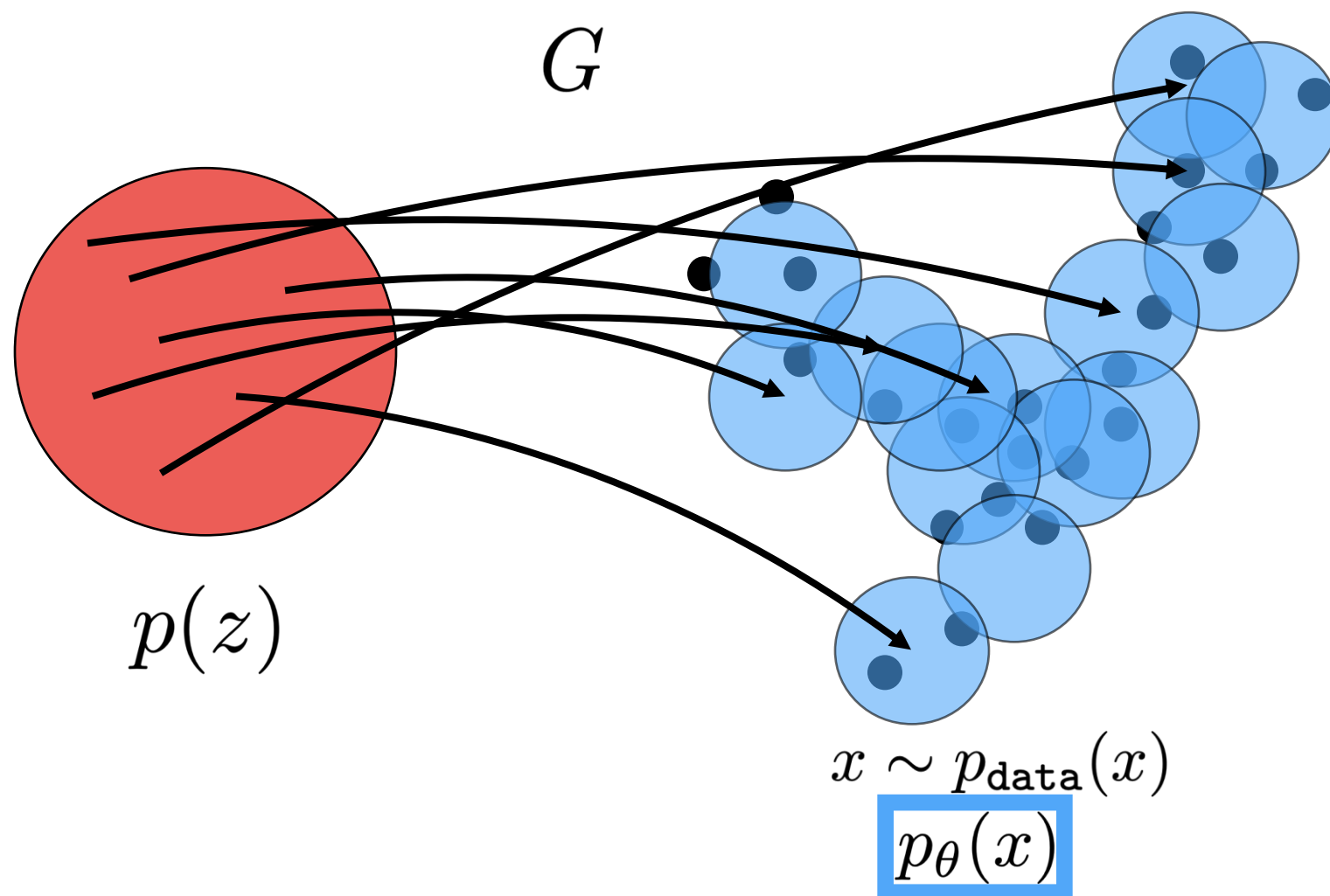
$p_{\theta}(x)$

Variational Autoencoders (VAEs)

[Kingma & Welling, 2014; Rezende, Mohamed, Wierstra 2014]

Prior distribution

Target distribution



Density model:

$$p_\theta(x) = \int p(x|z; \theta) p(z) dz$$

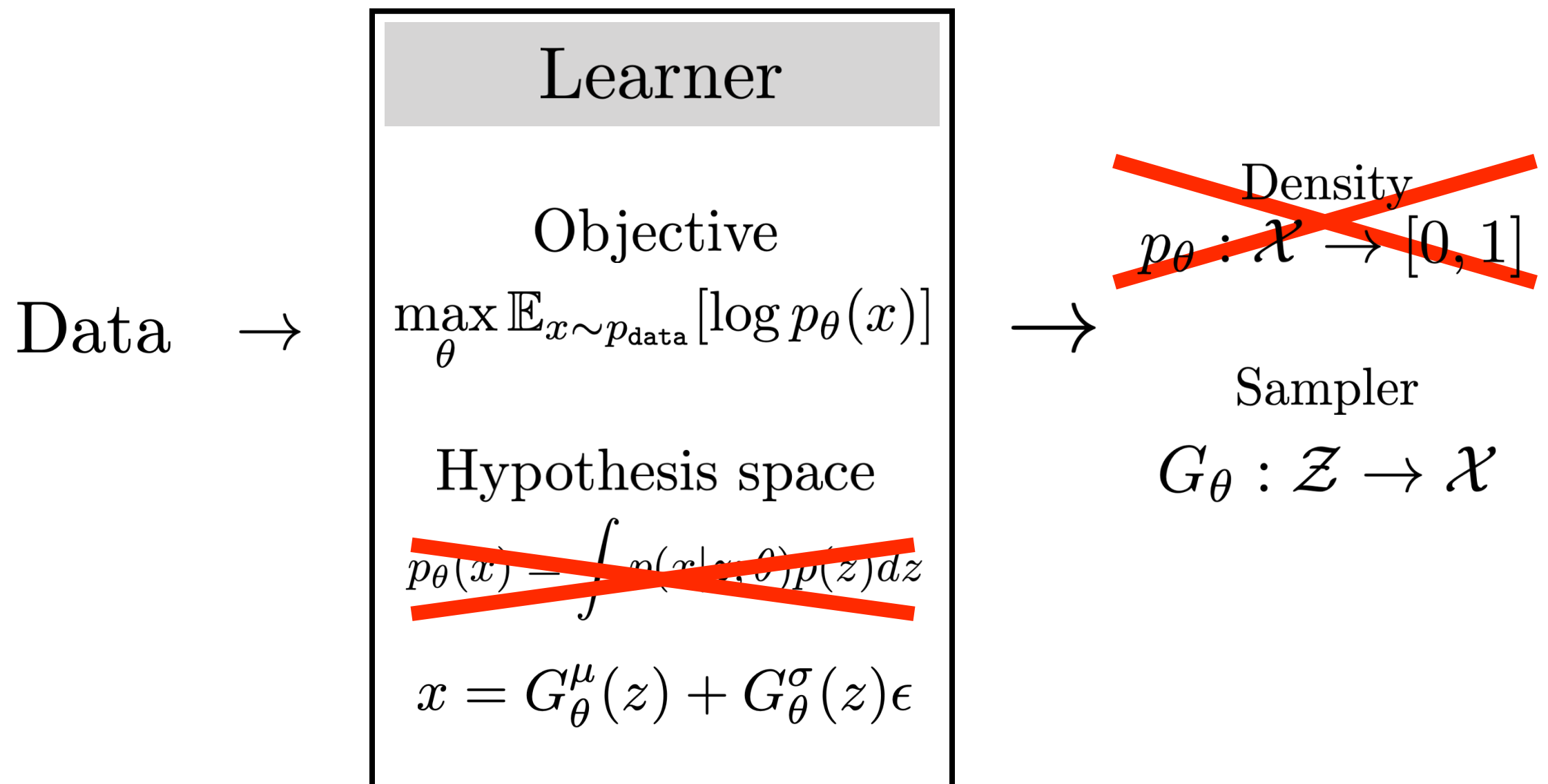
$$p(x|z; \theta) \sim \mathcal{N}(x; G_\theta^\mu(z), G_\theta^\sigma(z))$$

Sampling:

$$z \sim p(z) \quad \epsilon \sim \mathcal{N}(0, 1)$$

$$x = G_\theta^\mu(z) + G_\theta^\sigma(z)\epsilon$$

Variational Autoencoder (VAE)



Variational Autoencoders (VAEs)

Fitting a model to data requires computing $p_{\theta}(x)$

How to compute $p_{\theta}(x)$ efficiently?

$$p_{\theta}(x) = \int p(x|z; \theta)p(z)dz \quad \longleftarrow \quad \text{almost all terms are near zero}$$

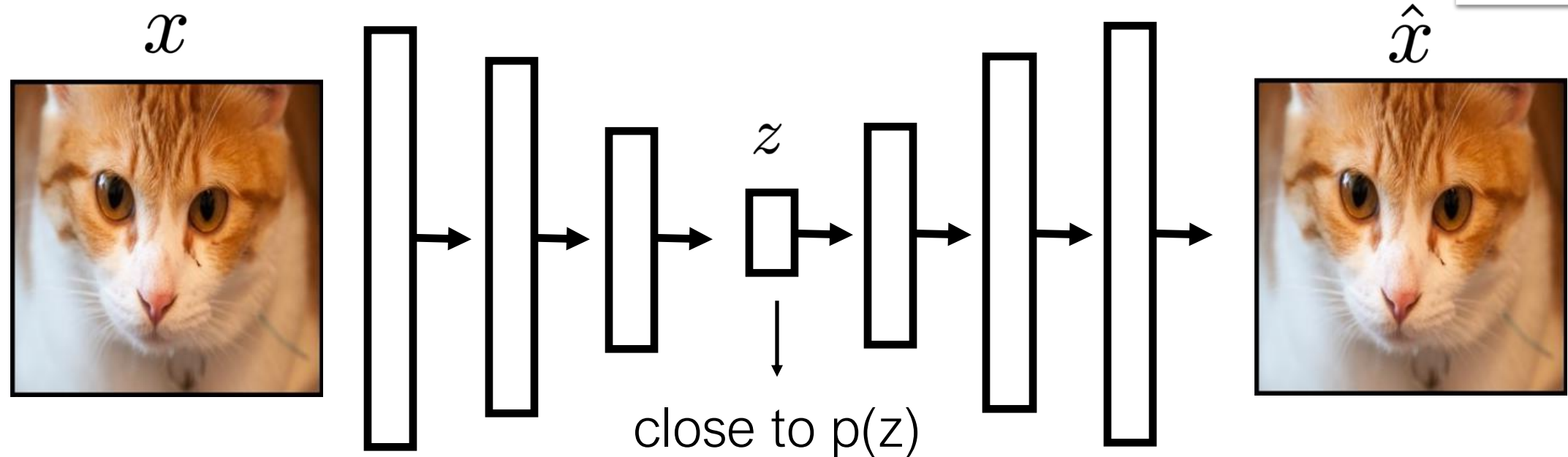
Train “inference network” $q_{\psi}(z|x)$

to give distribution over the z 's that are likely to produce x

Approximate $p_{\theta}(x)$ with $\mathbb{E}_{q_{\psi}(z|x)}[p_{\theta}(x|z)]$

Variational Autoencoders (VAEs)

encoder $q_{\psi}(z|x)$ $z = E_{\psi}^{\mu}(x) + E_{\psi}^{\sigma}(x) \cdot \epsilon_z$ generator $p_{\theta}(x|z)$ $\hat{x} = G_{\theta}^{\mu}(z)$

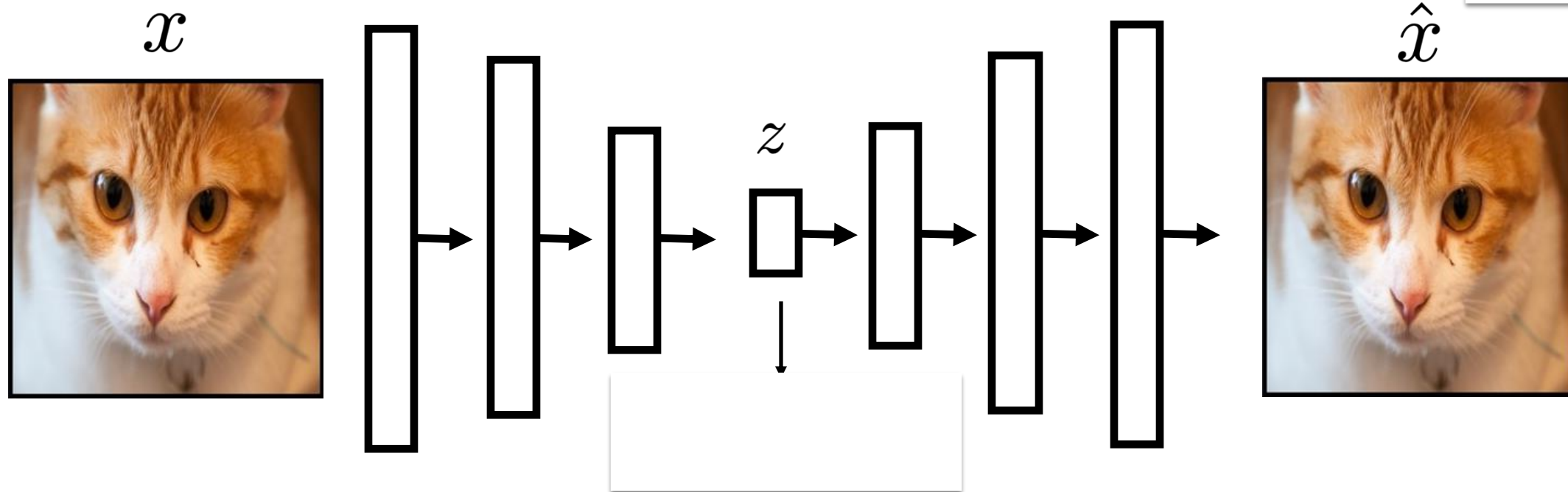


$$\begin{aligned} & \max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} [\log p_{\theta}(x)] && \text{Multi-variate Gaussian} \\ & \geq \max_{\theta, \psi} \mathbb{E}_{x_i \sim p_{\text{data}}} [\mathbb{E}_{q_{\psi}(z|x_i)} [p_{\theta}(x|z)] - \text{KL}(q_{\psi}(z|x_i) || p(z))] \\ & \quad \uparrow && \downarrow \\ & \text{reconstruction loss} && \text{KLD loss} \\ & ||x - \hat{x}||_2 && \text{KLD}(\mathcal{N}(E_{\psi}^{\mu}(x), E_{\psi}^{\sigma}(x)) | \mathcal{N}(0, I)) \end{aligned}$$

Autoencoders (AEs)

encoder $z = E_{\psi}^{\mu}(x)$

generator $\hat{x} = G_{\theta}^{\mu}(z)$



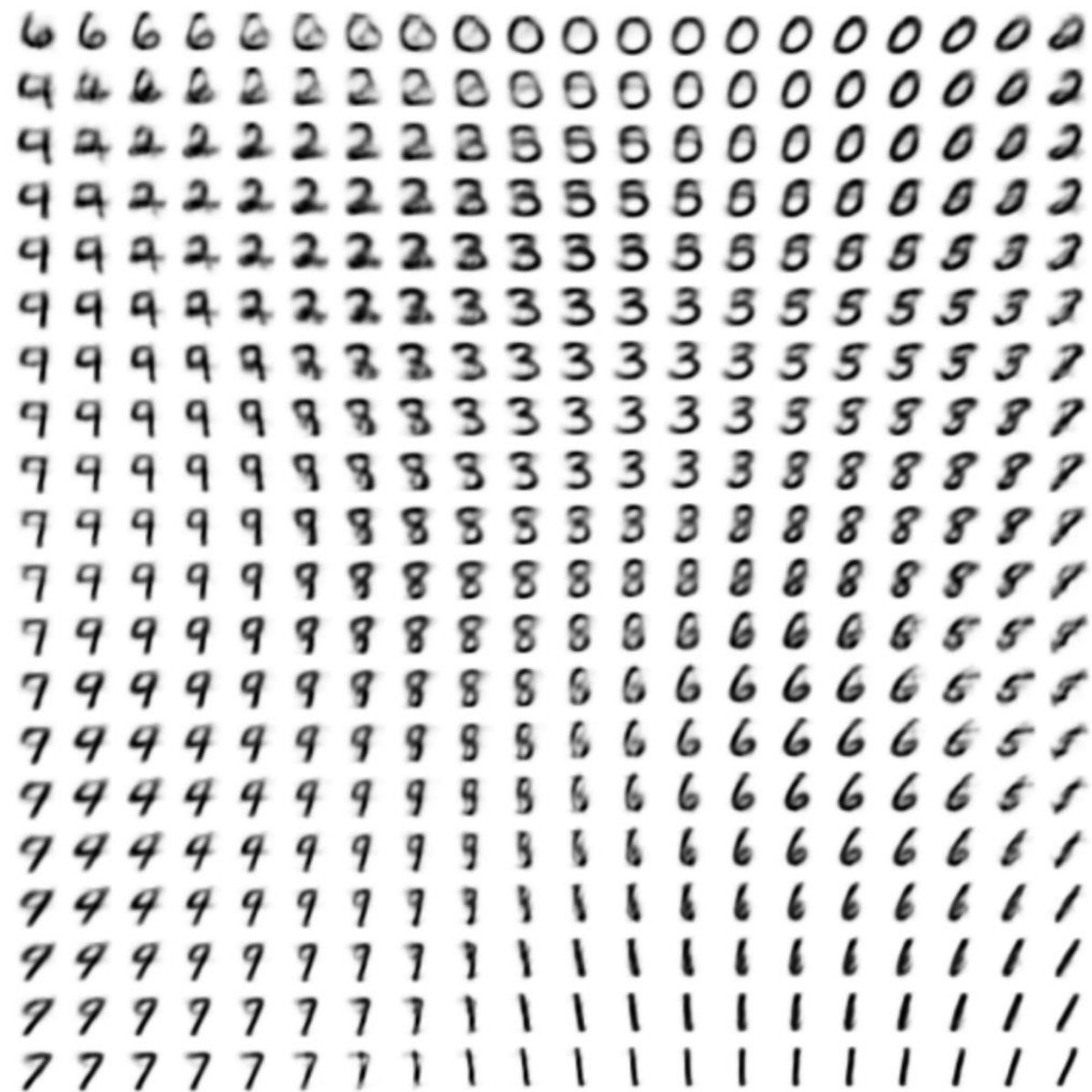
$$\max_{\theta, \psi} \mathbb{E}_{x_i \sim p_{\text{data}}} [\mathbb{E}_{q_{\psi}(z|x_i)} [p_{\theta}(x|z)]]$$

↑
reconstruction loss

$$\|x - \hat{x}\|_2$$



Variational Autoencoders (VAEs)



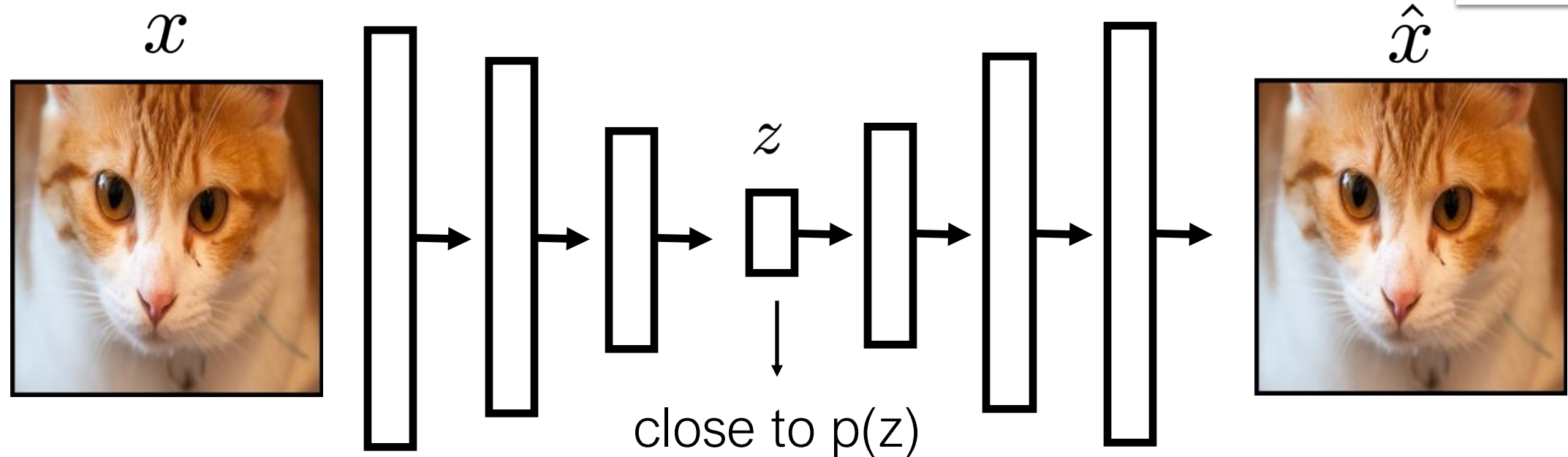
VAE with two-dimensional latent space

How to improve VAE?

- Why are the results blurry?
 - L2 reconstruction loss?
 - Lower bound might not be tight?
- How can we further improve results?

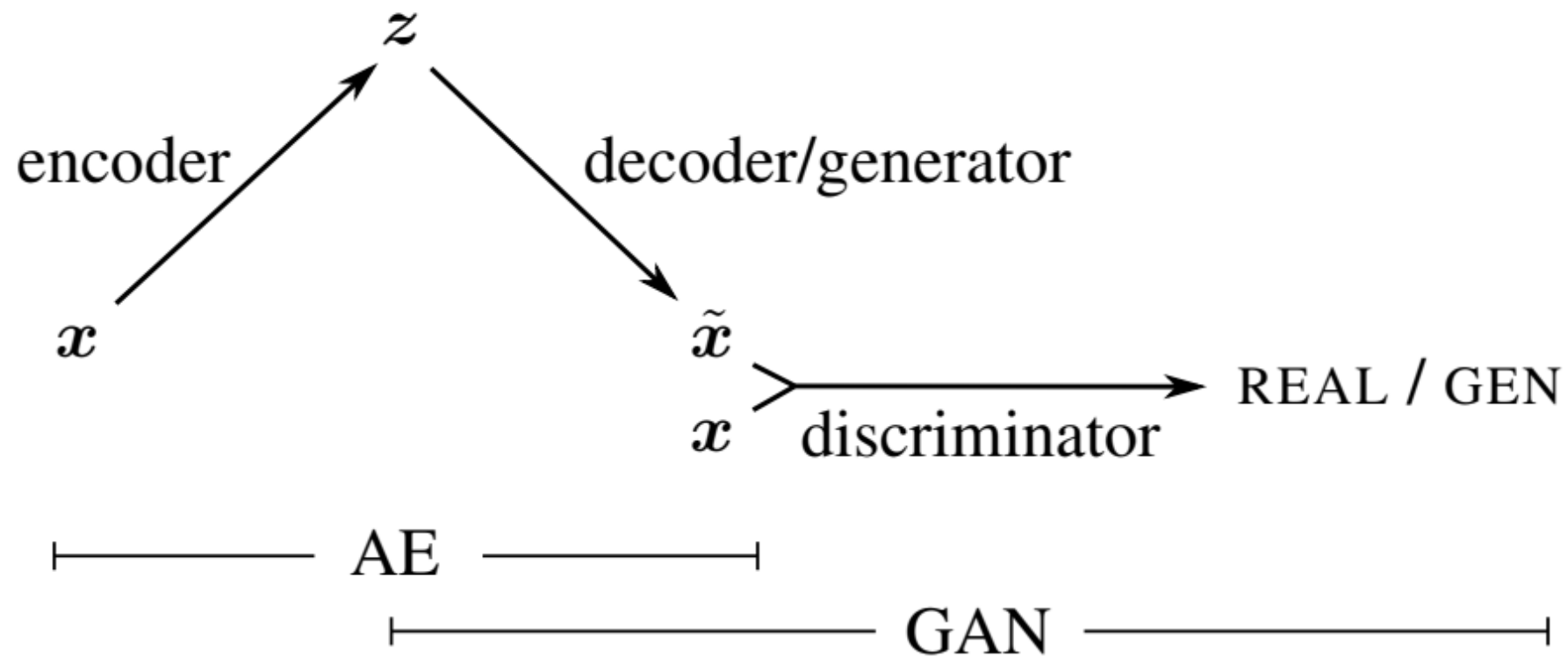
VAE + Perceptual Loss

encoder $q_\psi(z|x)$ $z = E_\psi^\mu(x) + E_\psi^\sigma(x) \cdot \epsilon_z$ generator $p_\theta(x|z)$ $\hat{x} = G_\theta^\mu(z)$



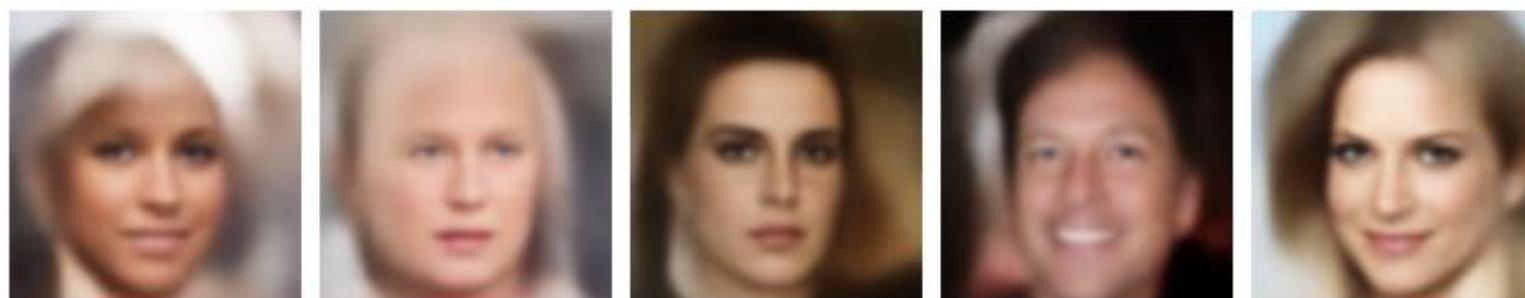
$$\begin{aligned} & \max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} [\log p_{\theta}(x)] && \text{Multi-variate Gaussian} \\ & \geq \max_{\theta, \psi} \mathbb{E}_{x_i \sim p_{\text{data}}} [\mathbb{E}_{q_{\psi}(z|x_i)} [p_{\theta}(x|z)] - \text{KL}(q_{\psi}(z|x_i) || p(z))] \\ & \quad \uparrow && \uparrow \\ & \text{Perceptual loss} && \text{KLD loss} \\ & ||F(x) - F(\hat{x})||_2 \quad \text{KLD}(\mathcal{N}(E_{\psi}^{\mu}(x), E_{\psi}^{\sigma}(x)) | \mathcal{N}(0, I)) \end{aligned}$$

VAE + GANs



VAE + GANs

VAE



VAE_{DisI}



VAE/GAN

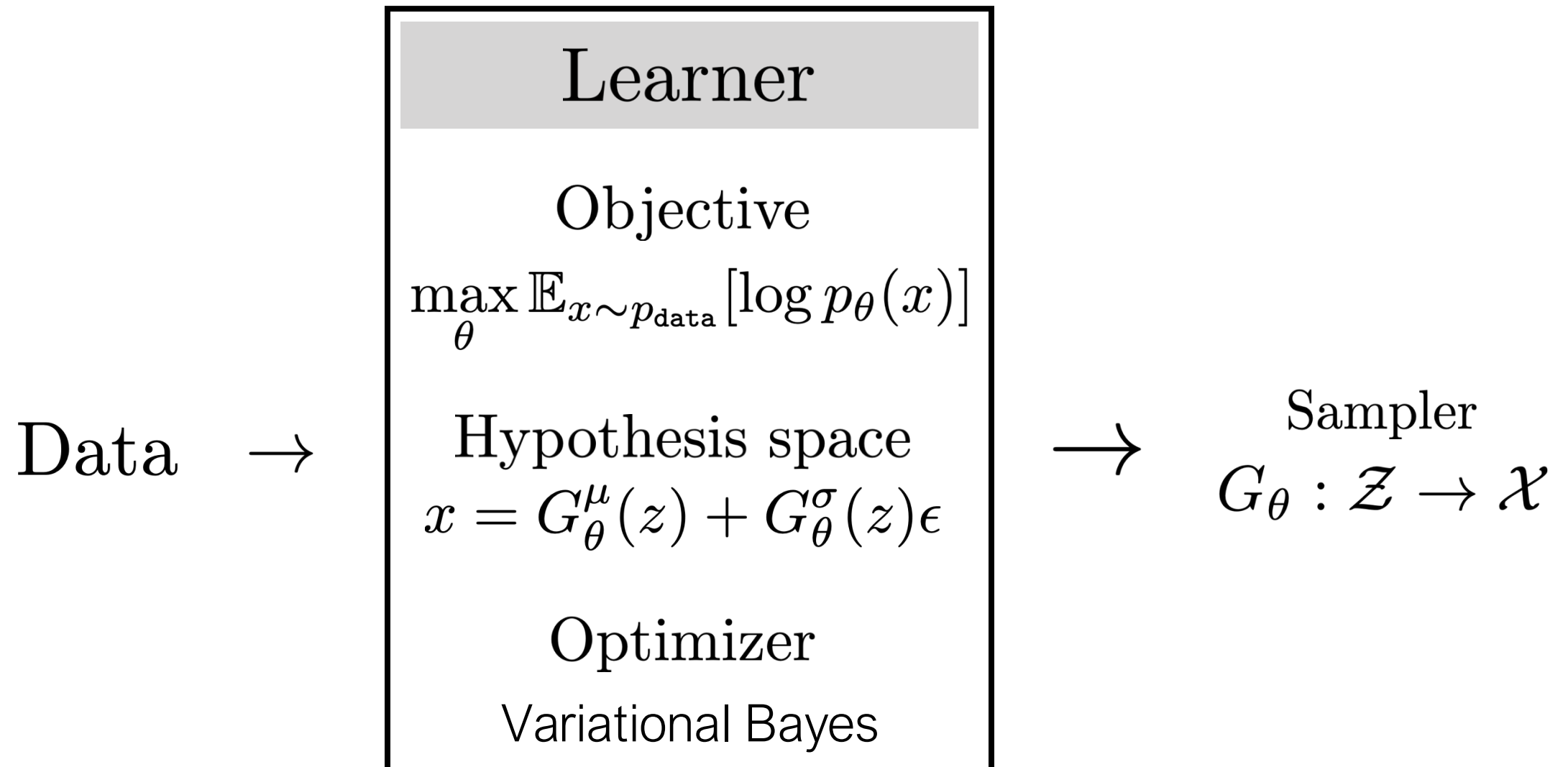


GAN



$$\text{VAE(DisI)} = \text{VAE} + \text{feature matching loss}$$

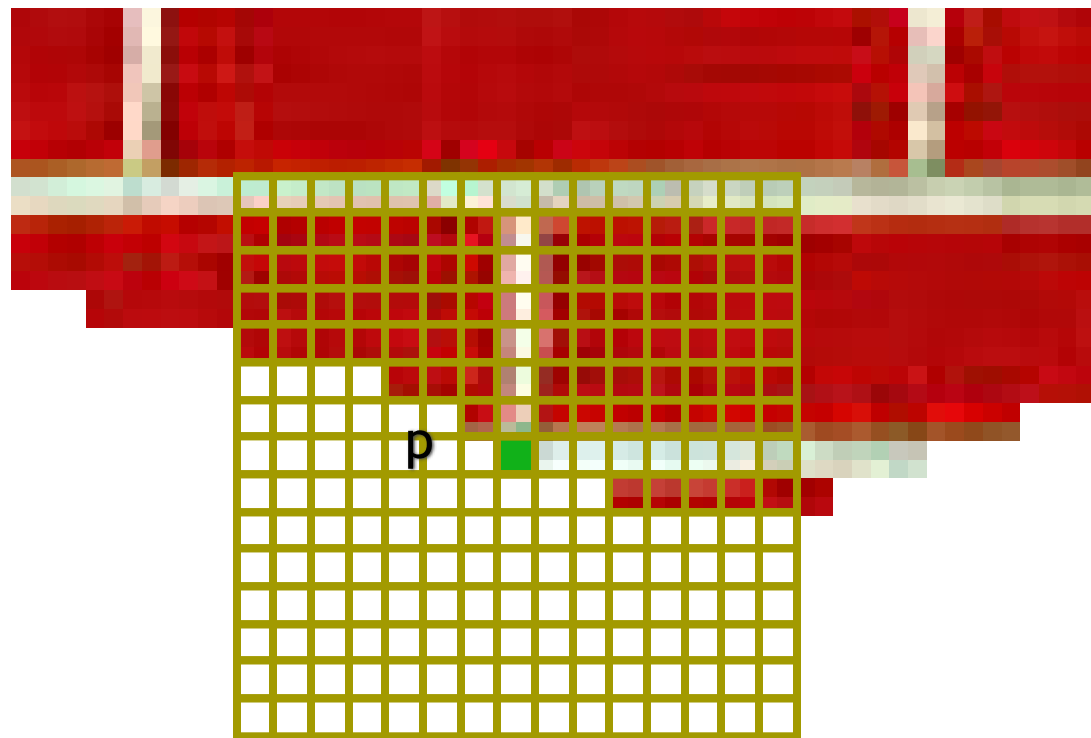
Variational Autoencoder (VAE)



Autoregressive Model

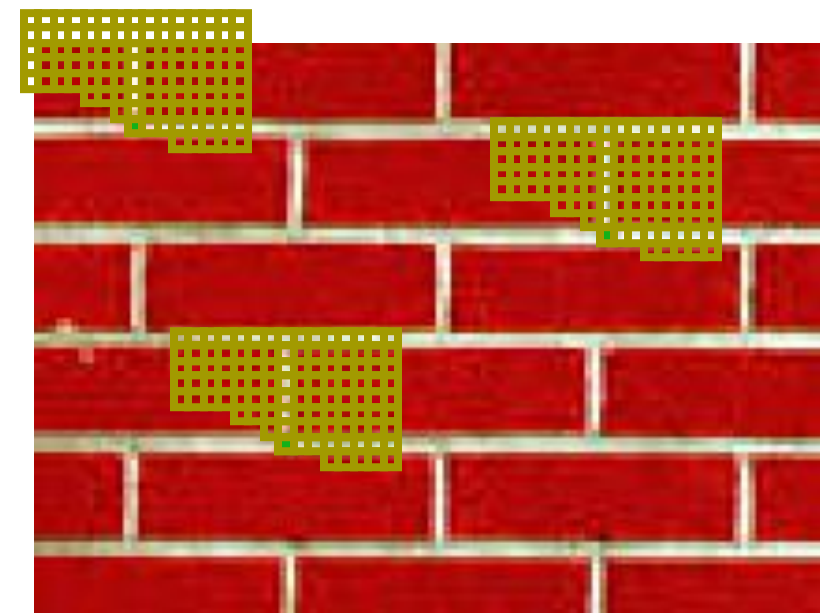

Texture synthesis by non-parametric sampling

[Efros & Leung 1999]



Synthesizing a pixel

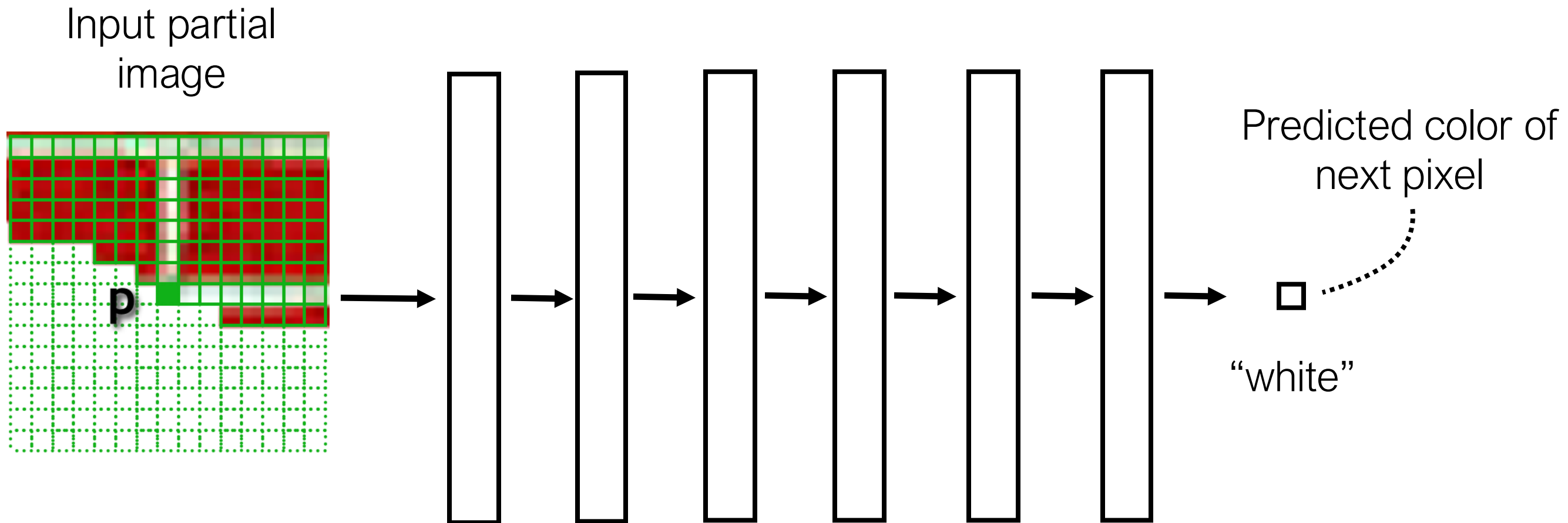
non-parametric
sampling



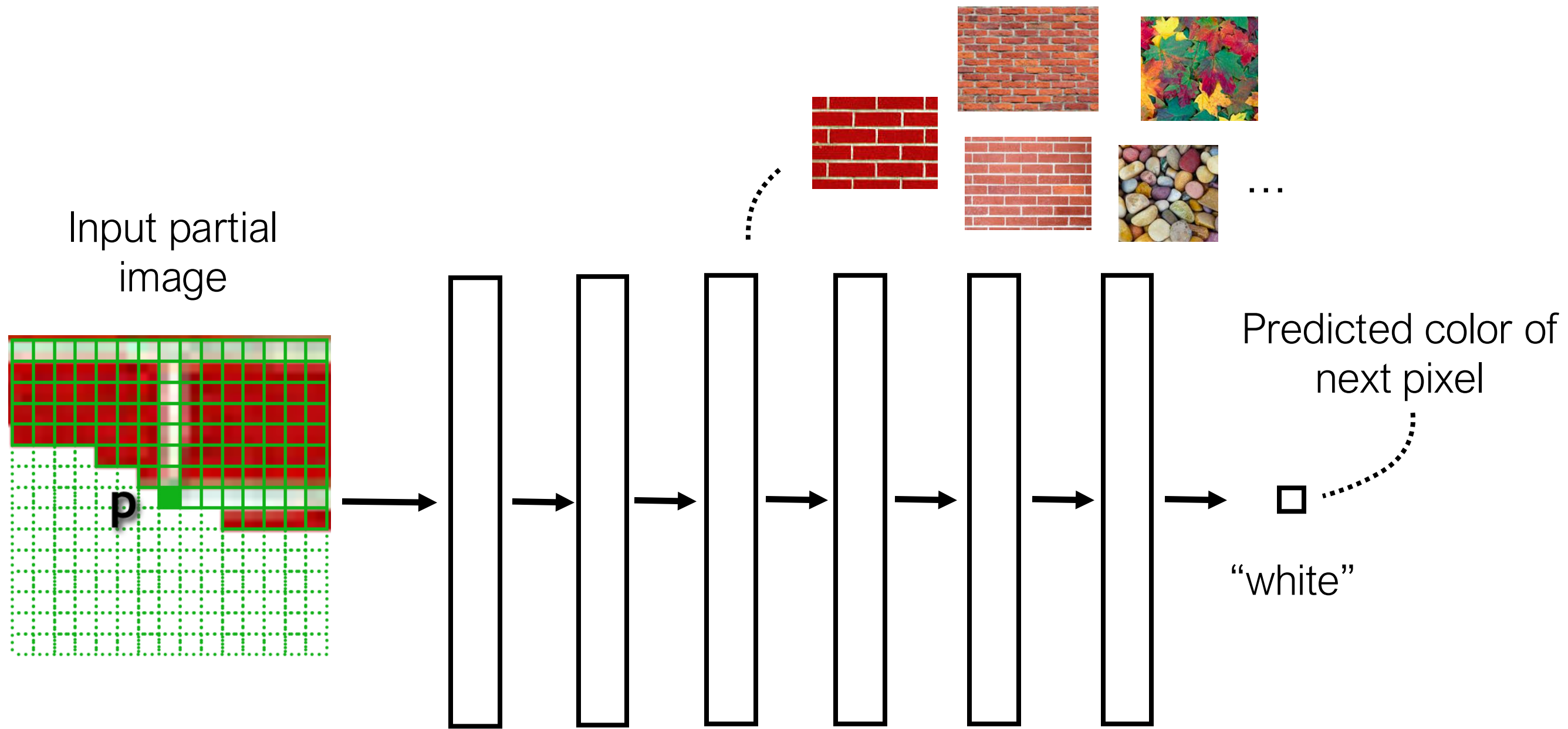
Input image

Models $P(p|N(p))$

Autoregressive image synthesis



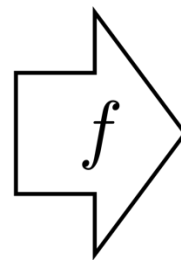
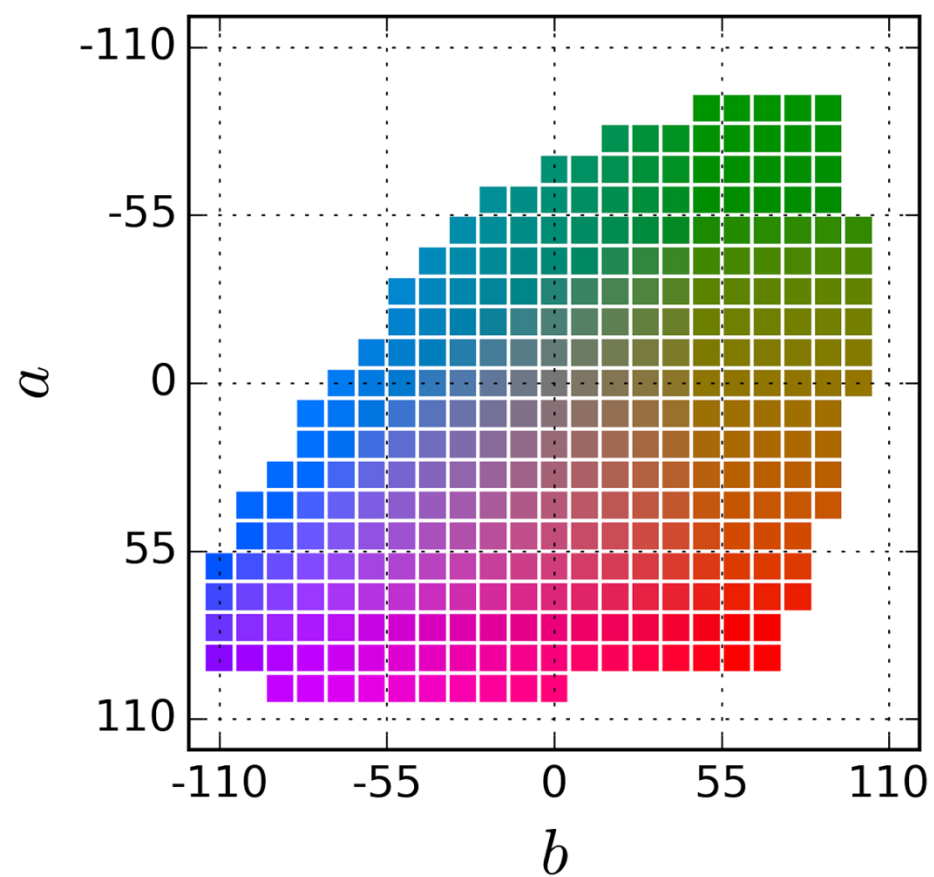
[PixelRNN, PixelCNN, van der Oord et al. 2016]



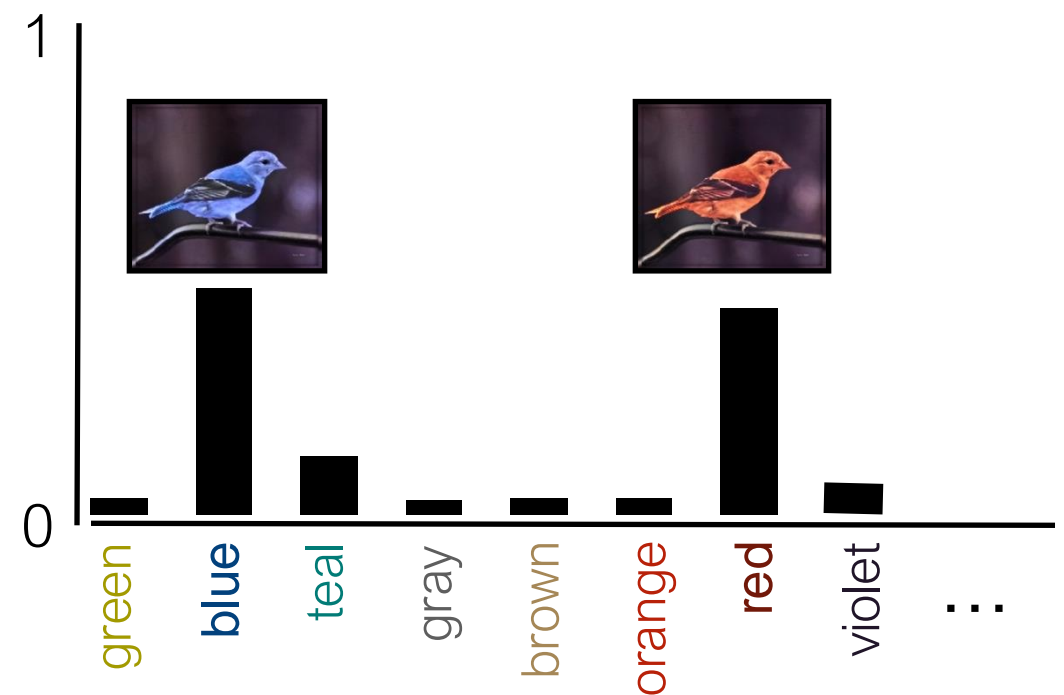
[PixelRNN, PixelCNN, van der Oord et al. 2016]

Recall: we can represent colors as discrete classes

$$\mathbf{y} \in \mathbb{R}^{H \times W \times K}$$



Prediction for a single pixel i, j



$$\mathcal{L}(\mathbf{y}, f_{\theta}(\mathbf{x})) = H(\mathbf{y}, \text{softmax}(f_{\theta}(\mathbf{x})))$$

And we can interpret the learner as modeling $P(\text{next pixel} \mid \text{previous pixels})$:

Softmax regression (a.k.a. multinomial logistic regression)

$$\hat{\mathbf{y}} \equiv [P_\theta(Y = 1 \mid X = \mathbf{x}), \dots, P_\theta(Y = K \mid X = \mathbf{x})] \quad \leftarrow \text{predicted probability of each class given input } \mathbf{x}$$

$$H(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{k=1}^K y_k \log \hat{y}_k \quad \leftarrow \text{picks out the } -\log \text{ likelihood of the ground truth class under the model prediction}$$

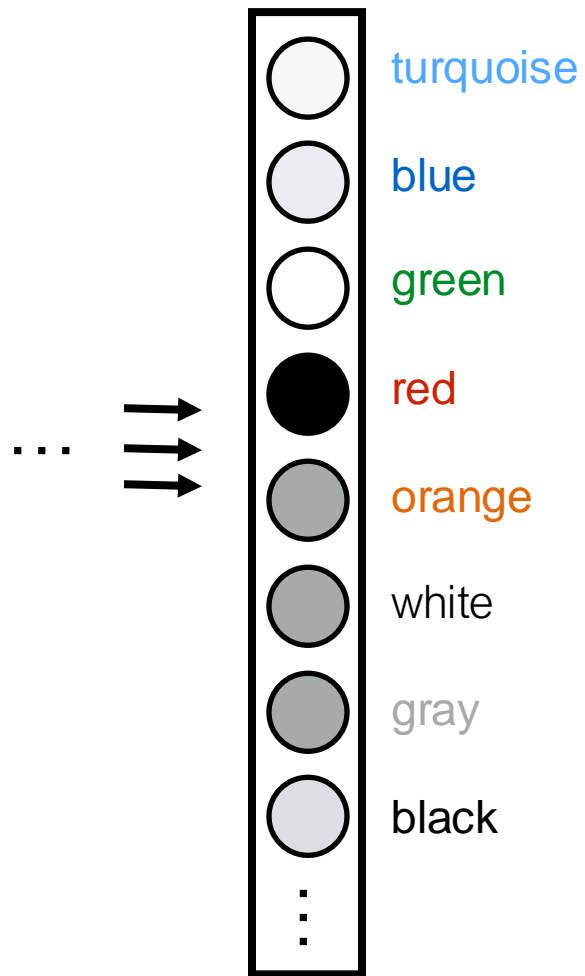
One-hot vector \nearrow

$\hat{\mathbf{y}}$

$$f^* = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^N H(\mathbf{y}_i, \hat{\mathbf{y}}_i) \quad \leftarrow \text{max likelihood learner!}$$

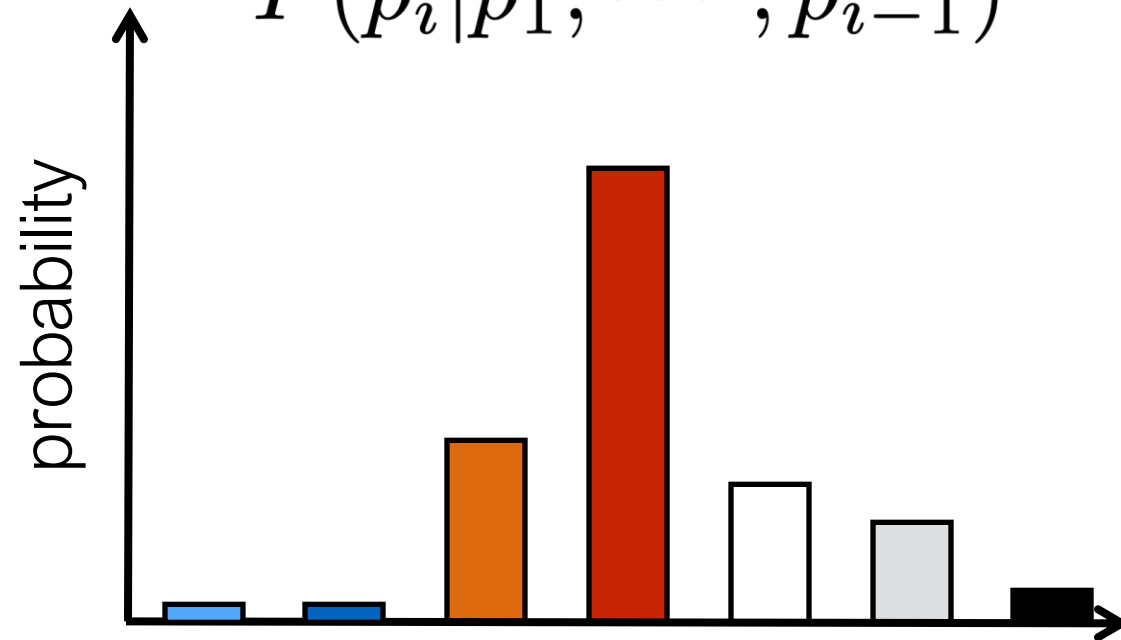
Cross-entropy loss \nearrow

Network output

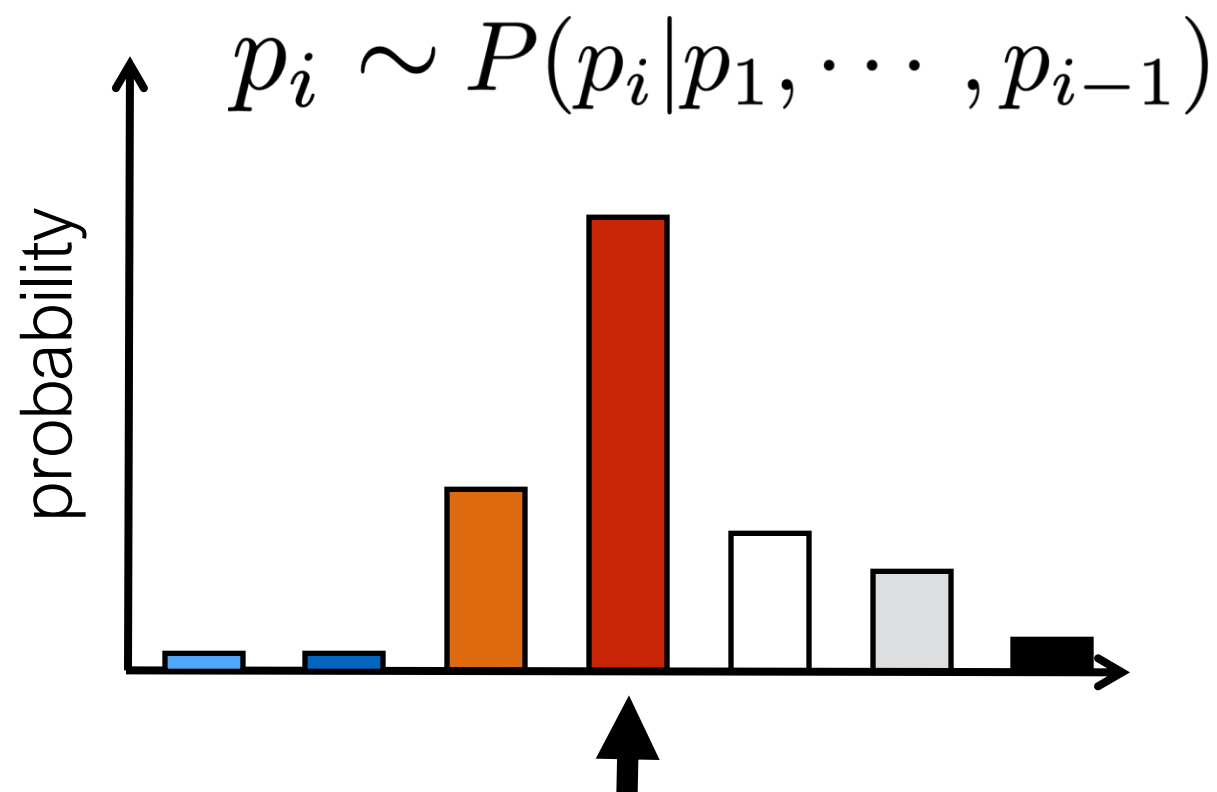
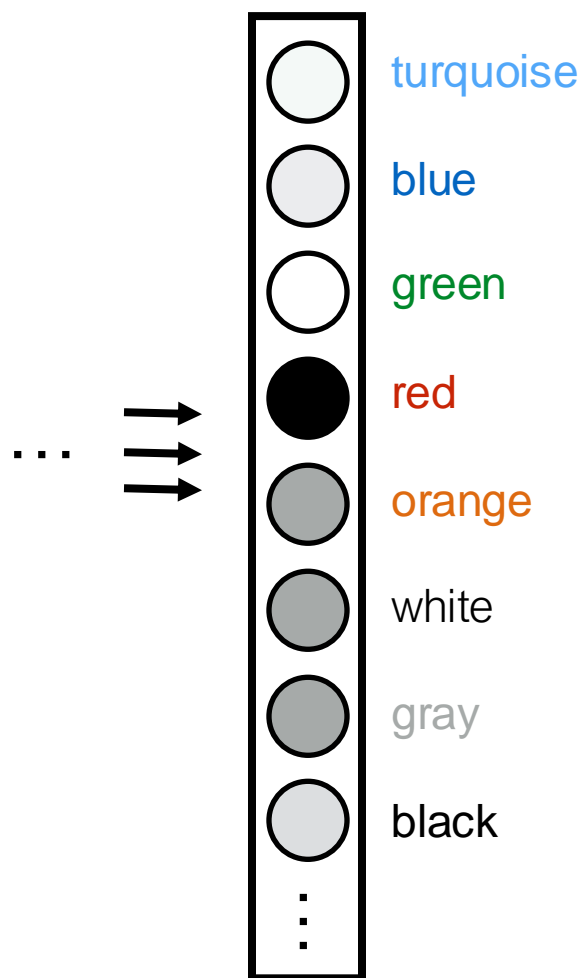


P(next pixel | previous pixels)

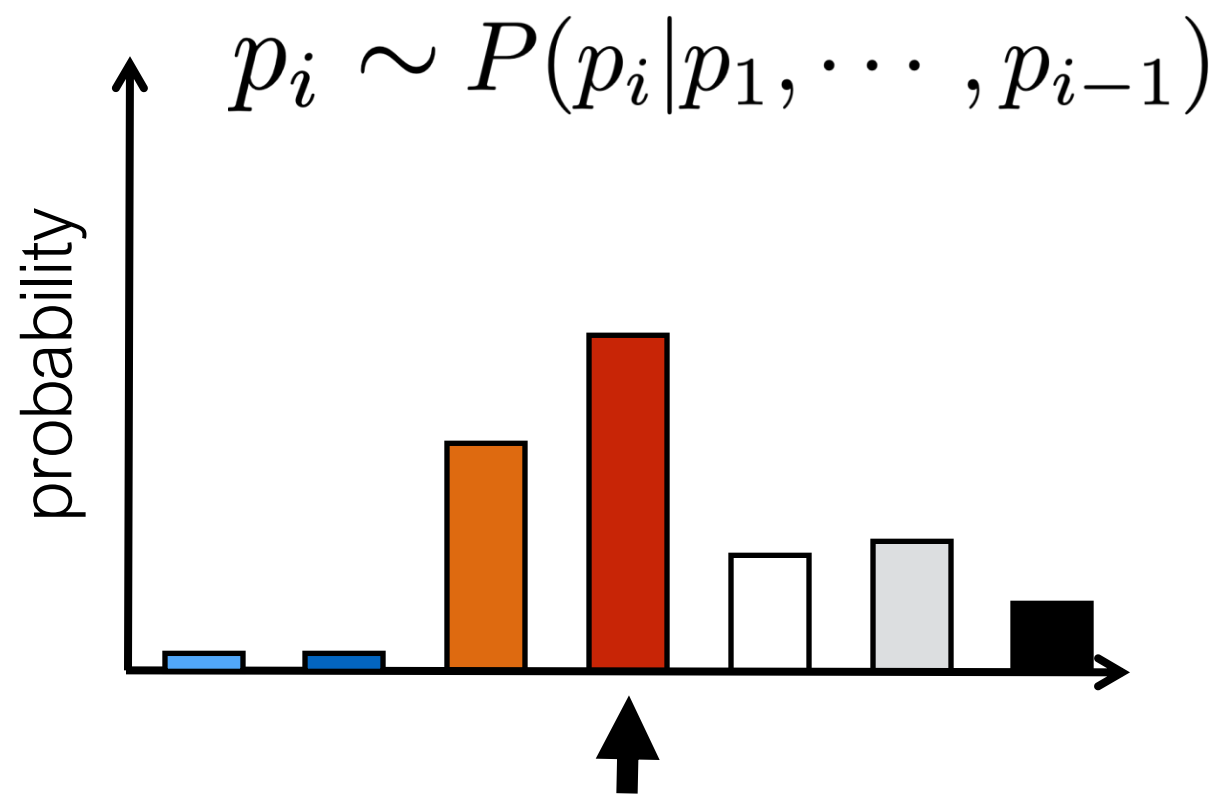
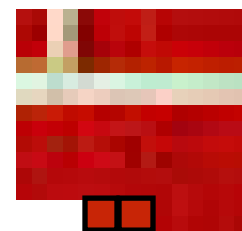
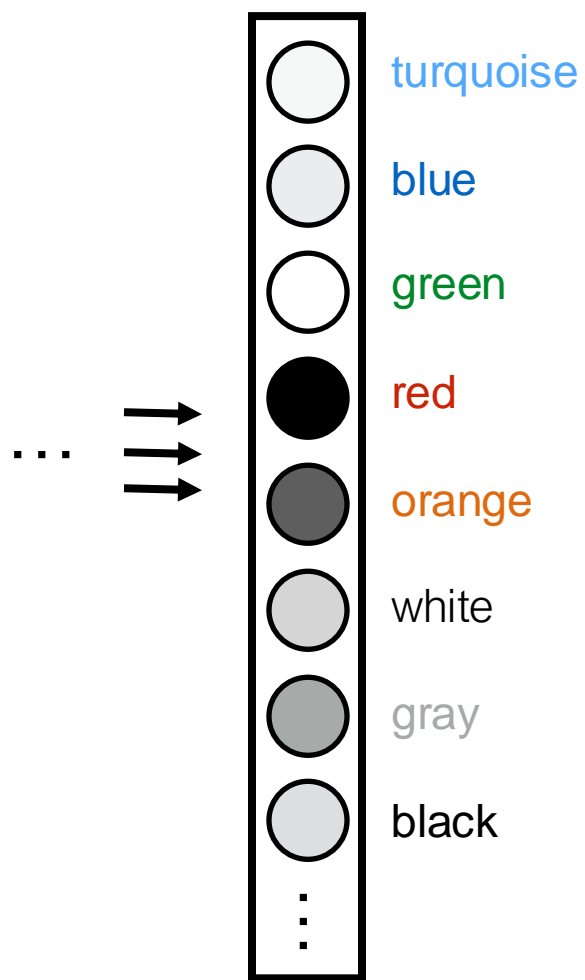
$$P(p_i | p_1, \dots, p_{i-1})$$



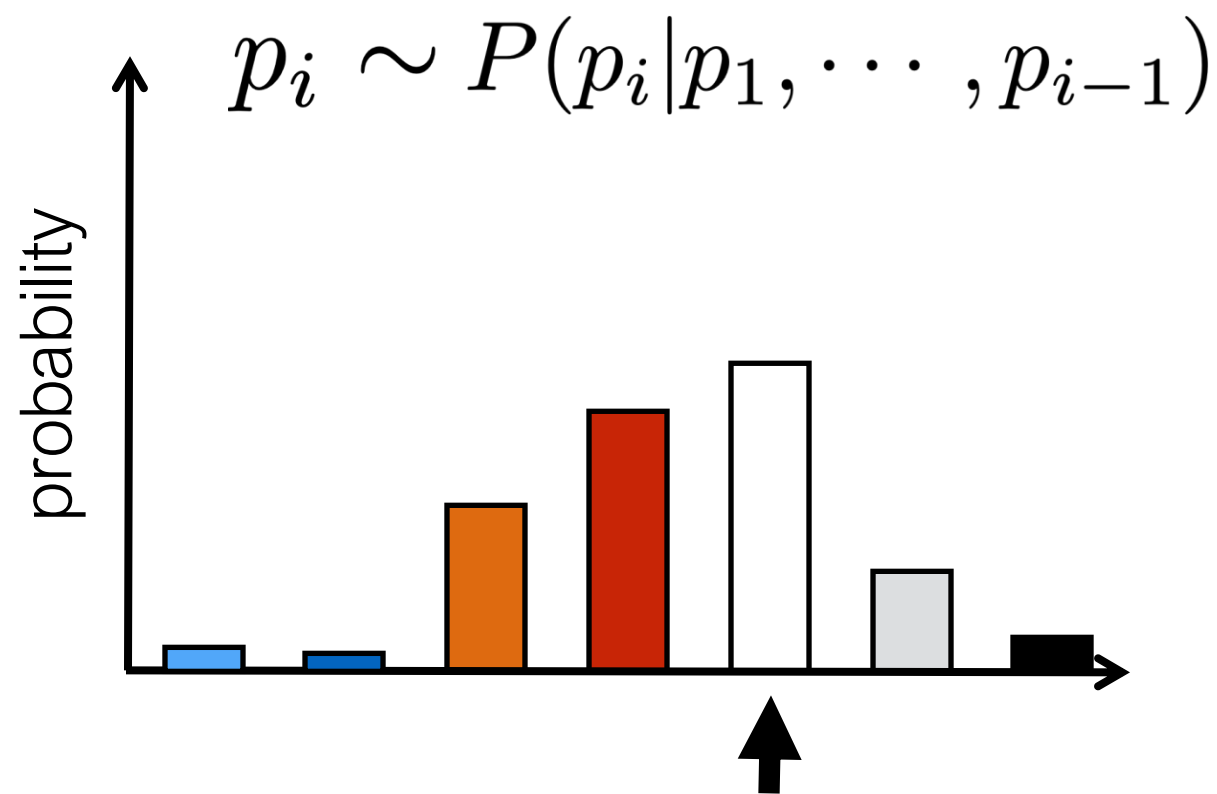
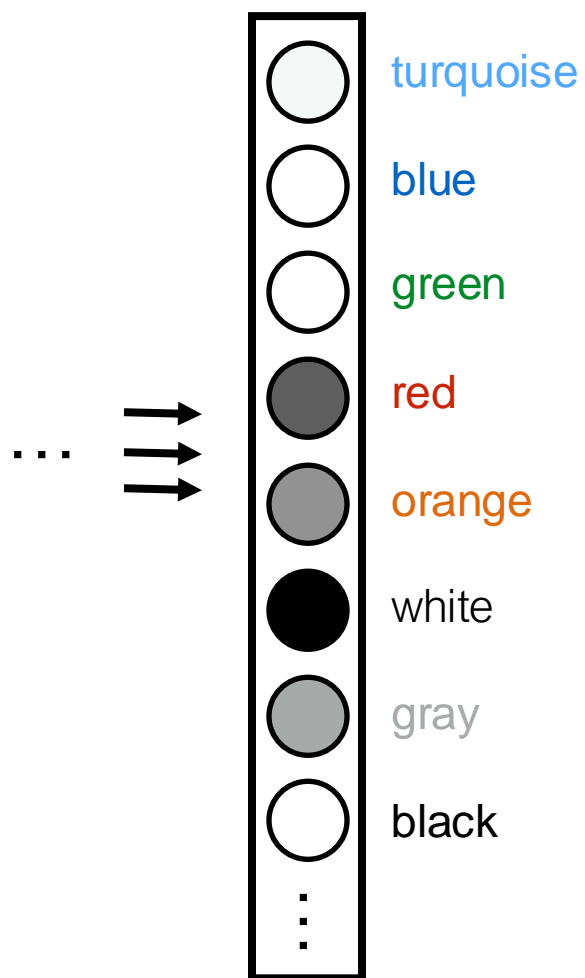
Network output



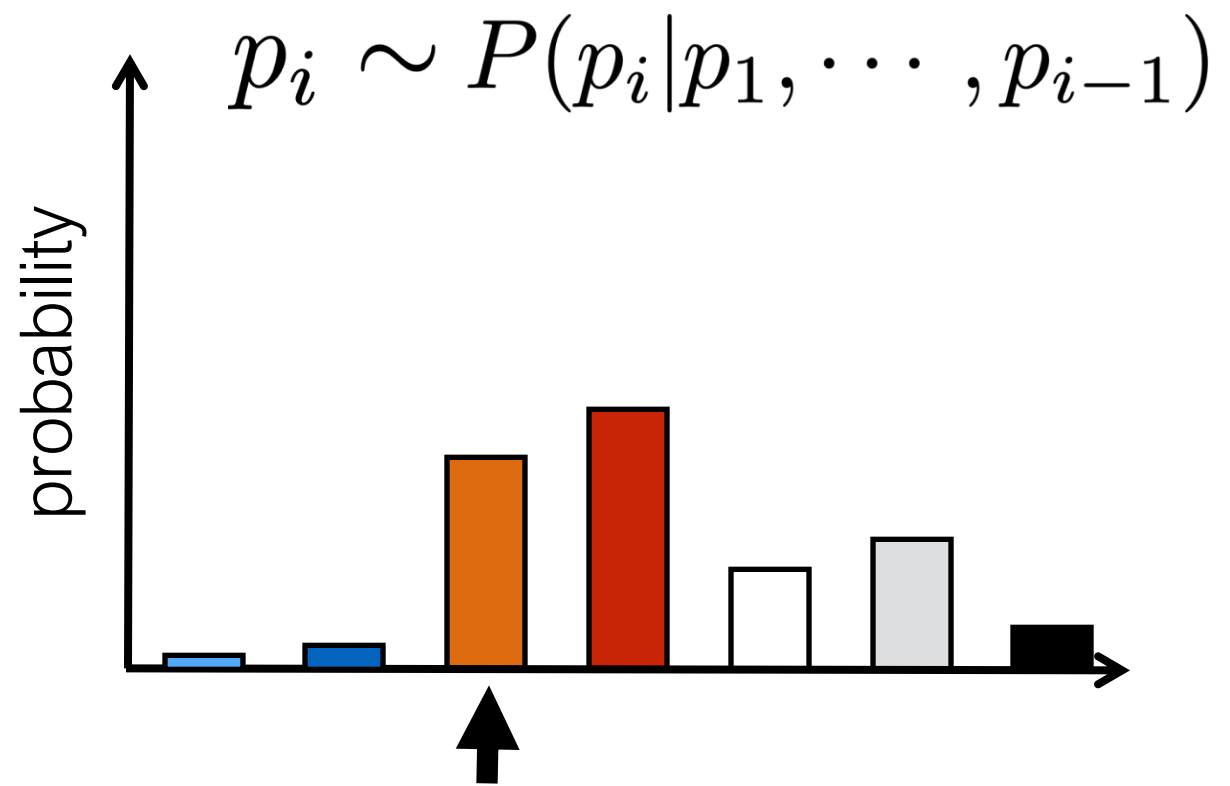
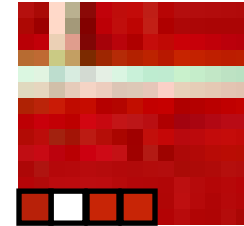
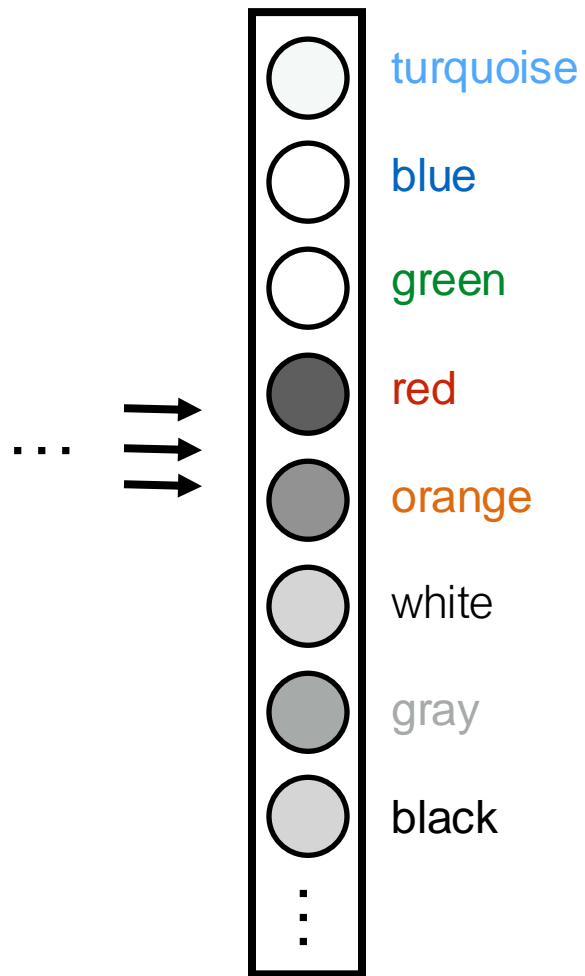
Network output



Network output



Network output



$$p_1 \sim P(p_1)$$

$$p_2 \sim P(p_2|p_1)$$

$$p_3 \sim P(p_3|p_1, p_2)$$

$$p_4 \sim P(p_4|p_1, p_2, p_3)$$

$$\{p_1, p_2, p_3, p_4\} \sim P(p_4|p_1, p_2, p_3)P(p_3|p_1, p_2)P(p_2|p_1)P(p_1)$$

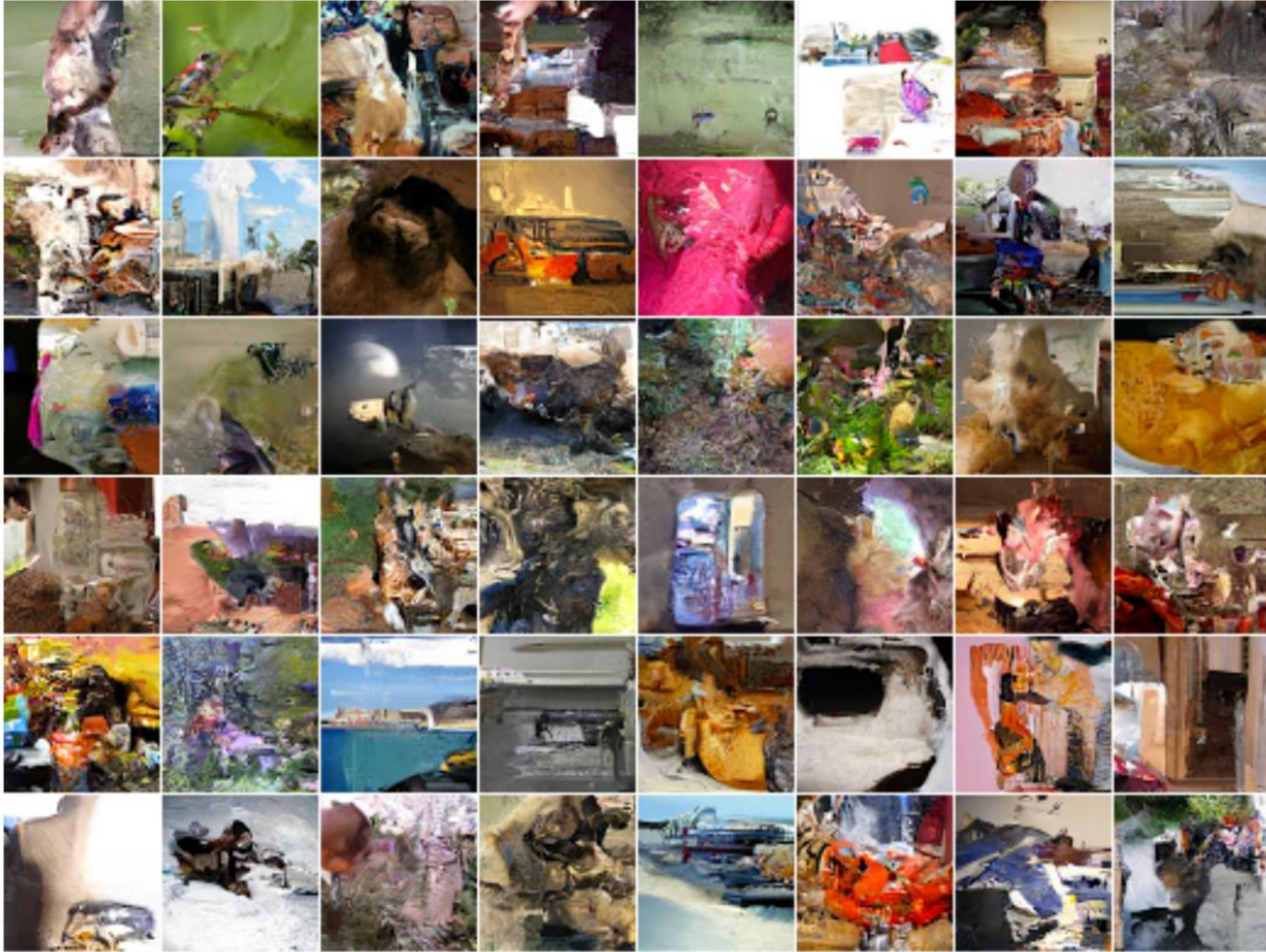
$$p_i \sim P(p_i|p_1, \dots, p_{i-1})$$

p_3 p_4 p_2 p_1



$$\mathbf{p} \sim \prod_{i=1}^N P(p_i|p_1, \dots, p_{i-1})$$

Samples from PixelRNN



[PixelRNN, van der Oord et al. 2016]

Image completions (conditional samples) from PixelRNN

occluded

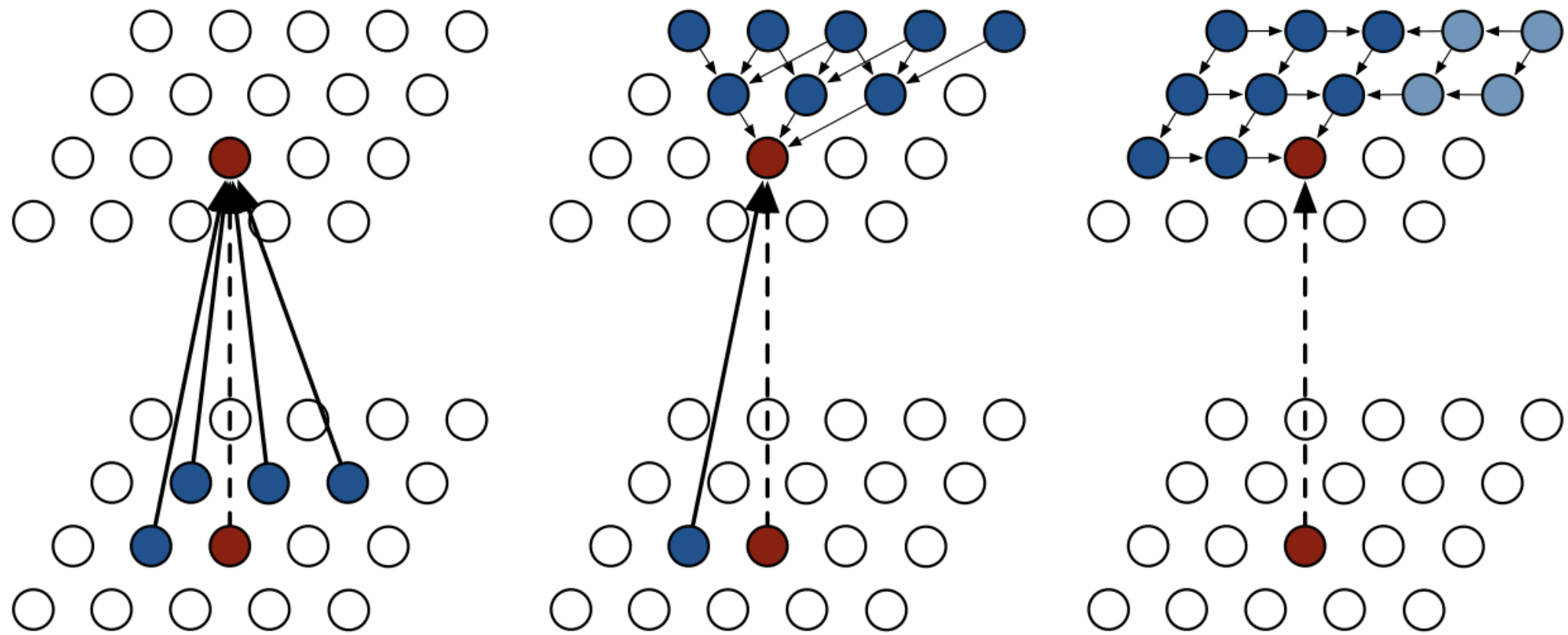
completions

original



[PixelRNN, van der Oord et al. 2016]

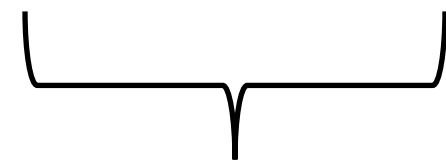
PixelCNN vs. PixelRNN



PixelCNN

Row LSTM

Diagonal BiLSTM



PixelRNN

Checkout PixelCNN++ [Salimans et al., 2017] (+ coarse-to-fine, ResNet, whole pixels, etc.)

How to improve PixelCNN?

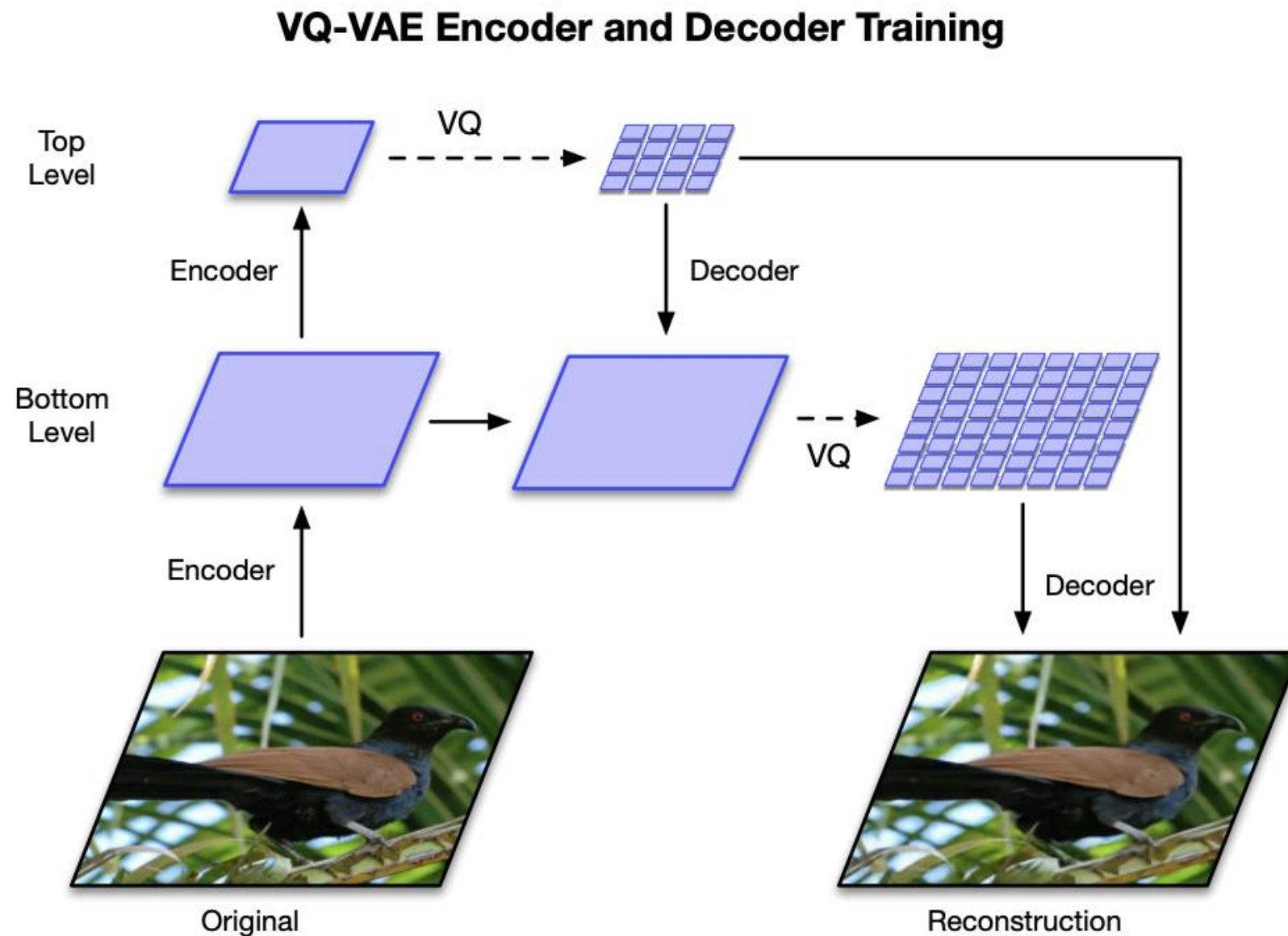
- What are the limitations of PixelCNN/RNNs?
 - Slow sampling time.
 - May accumulate errors over multiple steps. (might not be a big issue for image completion)
- How can we make it faster?
- How can we further improve results?

VQ-VAE-2: VAE+PixelCNN



VQ (Vector quantization) maps continuous vectors into discrete codes
Common methods: clustering (e.g., k-means)

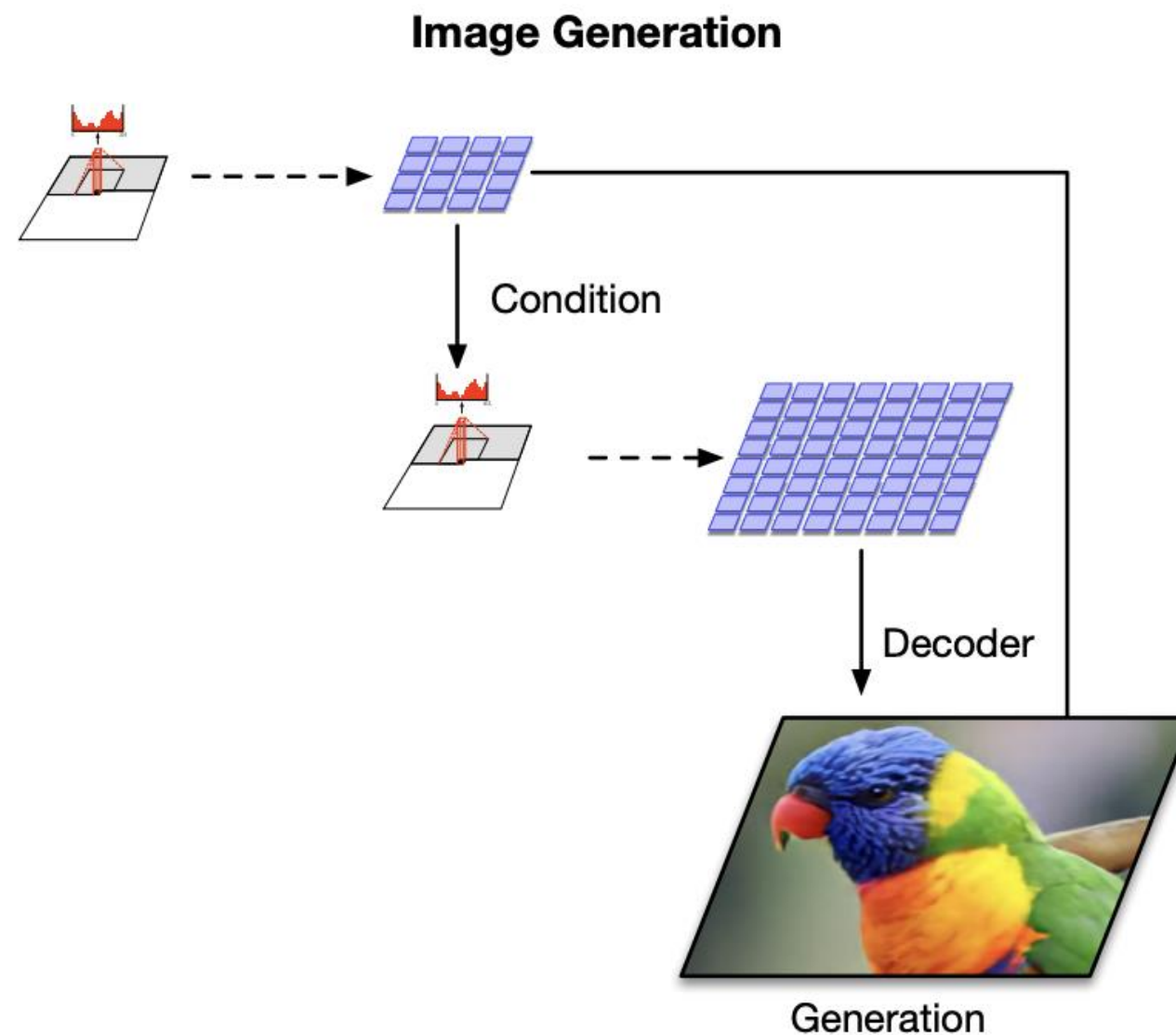
VQ-VAE-2: VAE+PixelCNN



VAE+VQ: learn a more compact codebook for PixelCNN (instead of pixels)

PixelCNN: use a more expressive bottleneck for VAE (instead of Gaussian)

VQ-VAE-2: VAE+PixelCNN



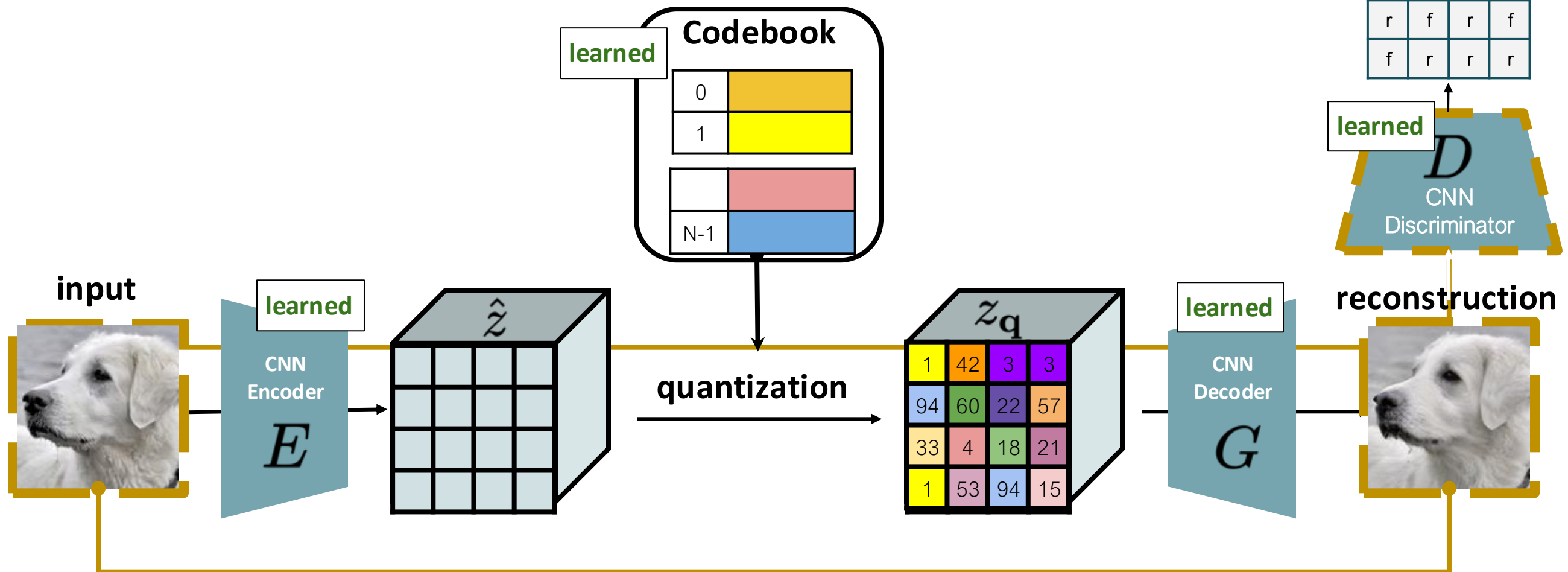
VAE+VQ: learn a more compact codebook for PixelCNN (instead of pixel colors)

PixelCNN: use a more expressive bottleneck for VAE (instead of Gaussian prior)

How to Improve further?

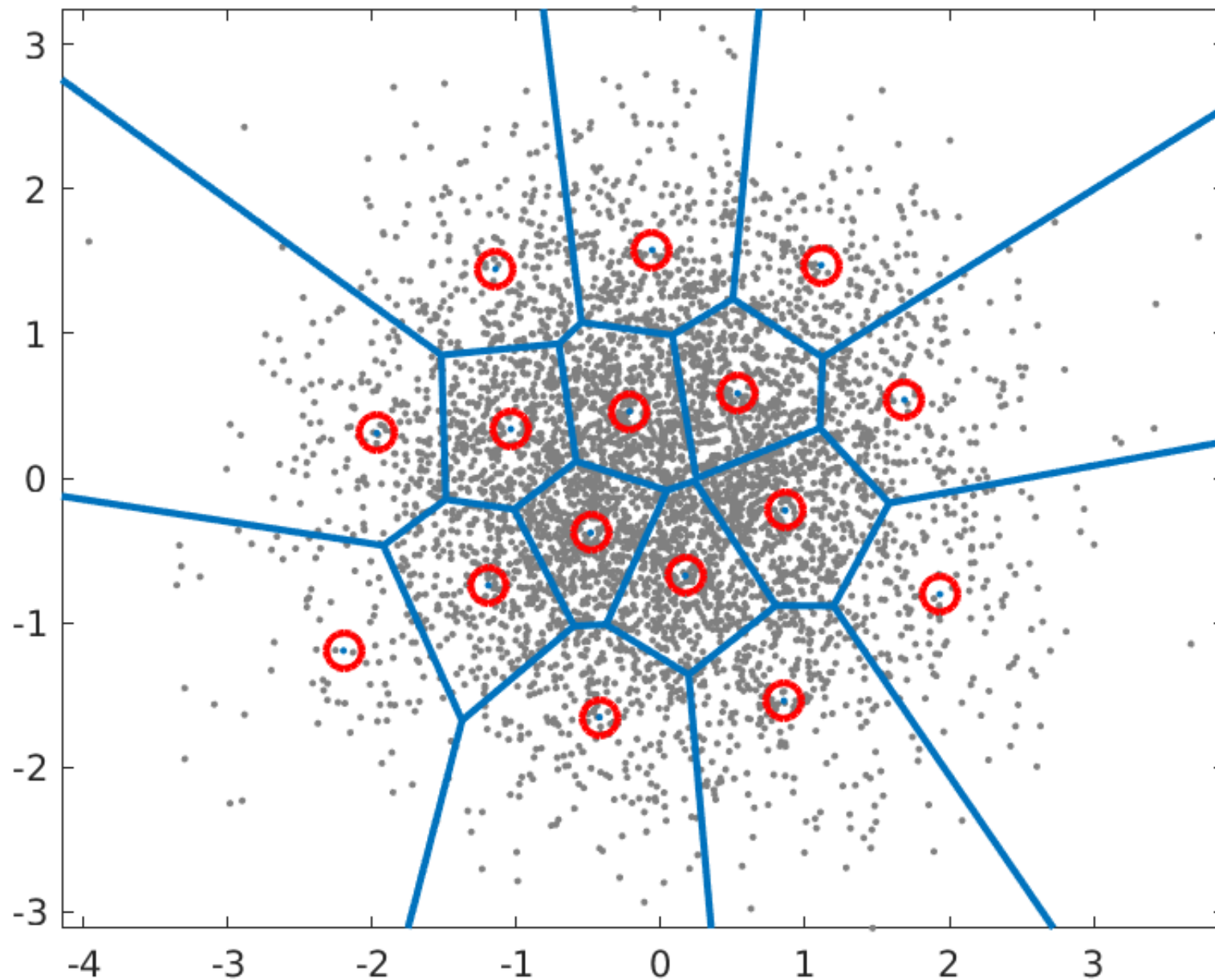
- Better architectures
- Better loss functions for encoder-decoder

From VQ-VAE to VQGAN



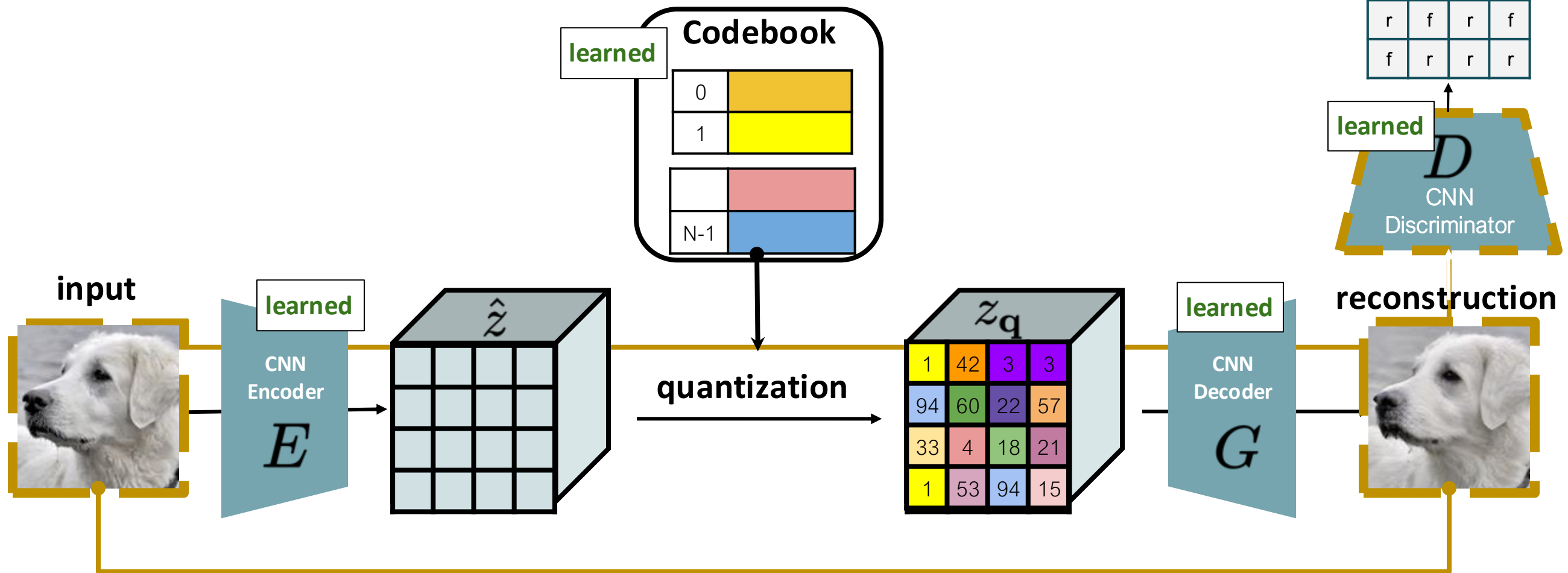
- 1) replace L2/L1 rec. loss with Perceptual loss (includes pixel-level)
- 2) add Discriminator to favor realism over per-pixel reconstruction
- 3) use transformer rather than CNN

Vector Quantization (VQ)



K-means, EM (GMM), end-to-end learning

From VQ-VAE to VQGAN



$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{VQ}} + \lambda \mathcal{L}_{\text{GAN}} \quad \text{where} \quad \lambda = \frac{\nabla_{G_L} [\mathcal{L}_{\text{rec}}]}{\nabla_{G_L} [\mathcal{L}_{\text{GAN}}] + \delta}$$

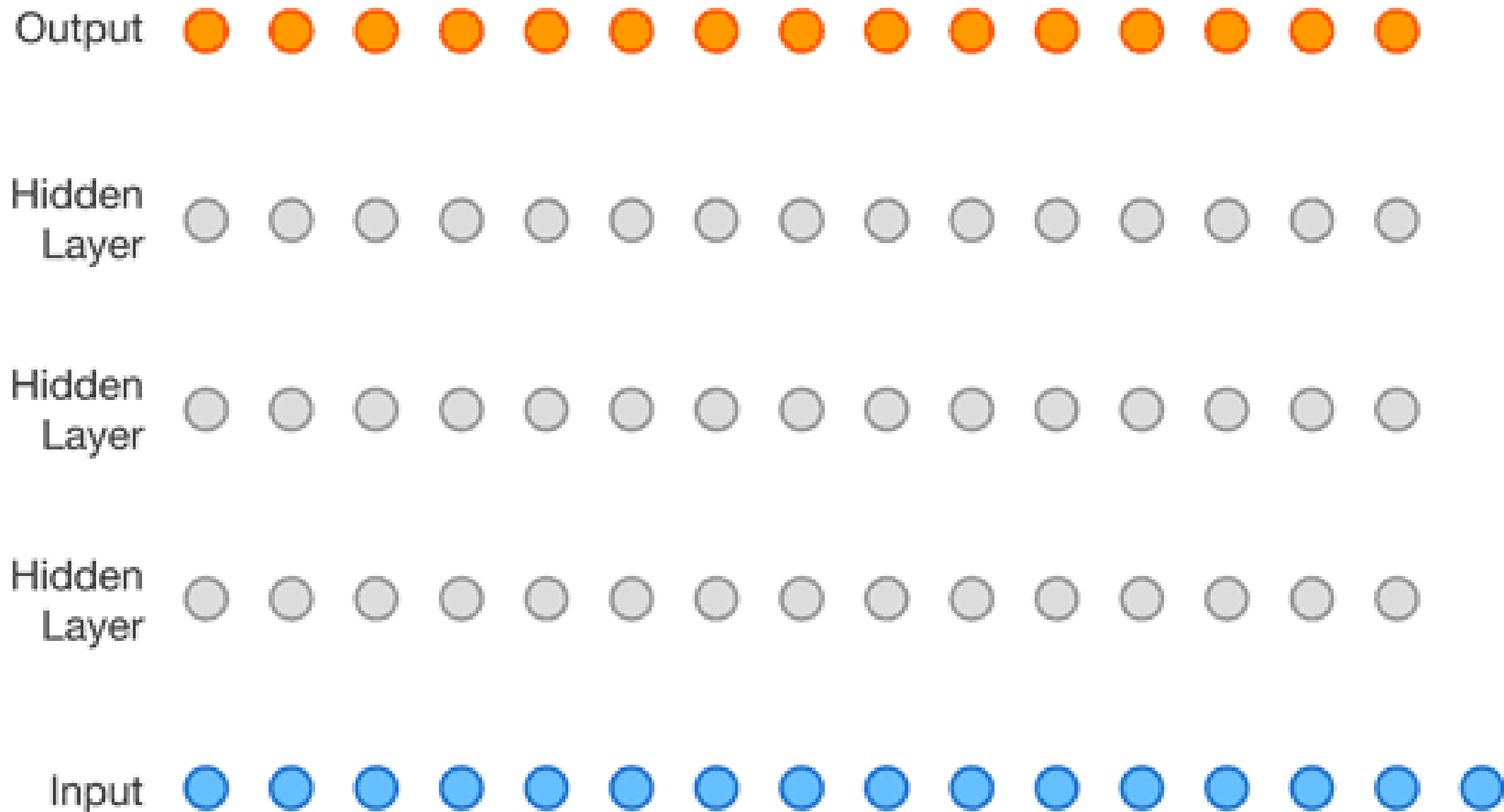


Slide credit: Robin Rombach



Slide credit: Robin Rombach

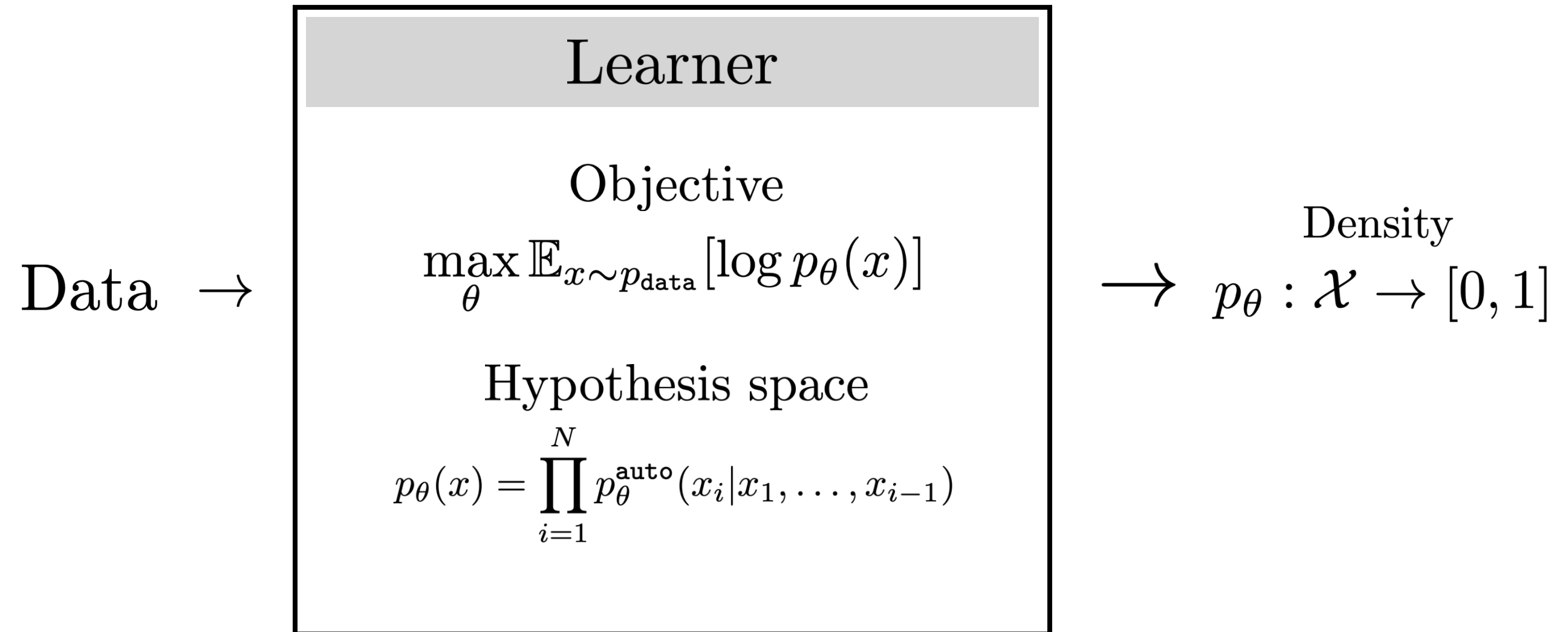
WaveNet



[Wavenet, <https://deepmind.com/blog/wavenet-generative-model-raw-audio/>]

Auto-regressive models works extremely well for audio/music data.

Autoregressive Model



Autoregressive probability model

$$\mathbf{p} \sim \prod_{i=1}^N P(p_i | p_1, \dots, p_{i-1})$$

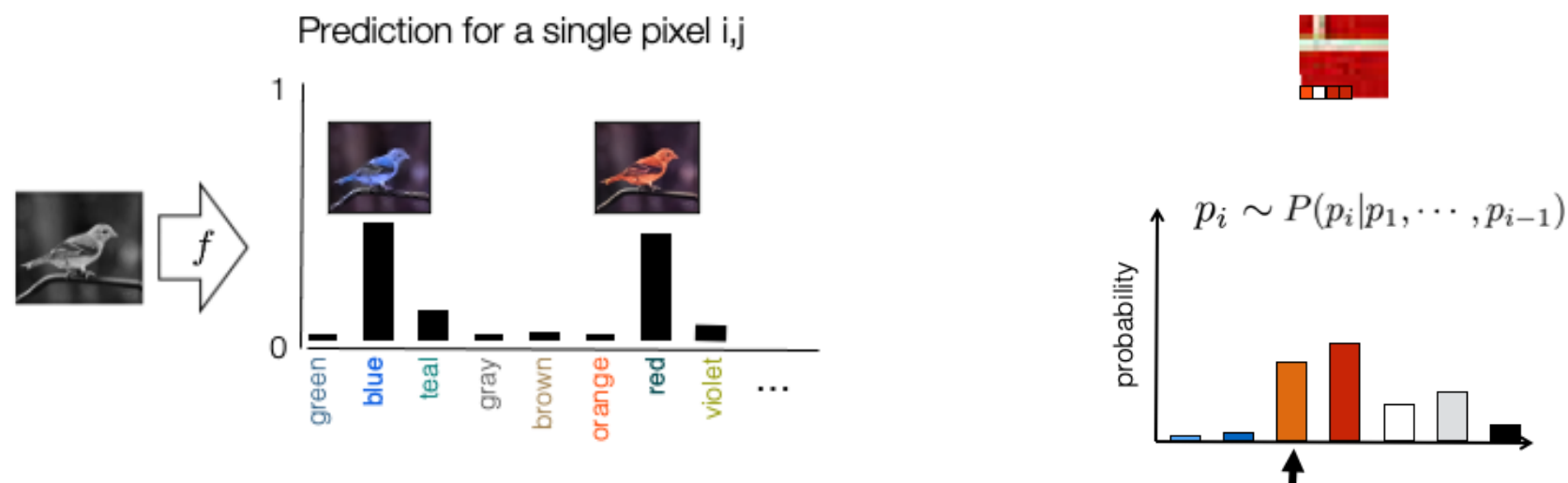
$$P(\mathbf{p}) = \prod_{i=1}^N P(p_i | p_1, \dots, p_{i-1}) \quad \leftarrow \text{General product rule}$$

The sampling procedure we defined above takes exact samples from the learned probability distribution (pmf).

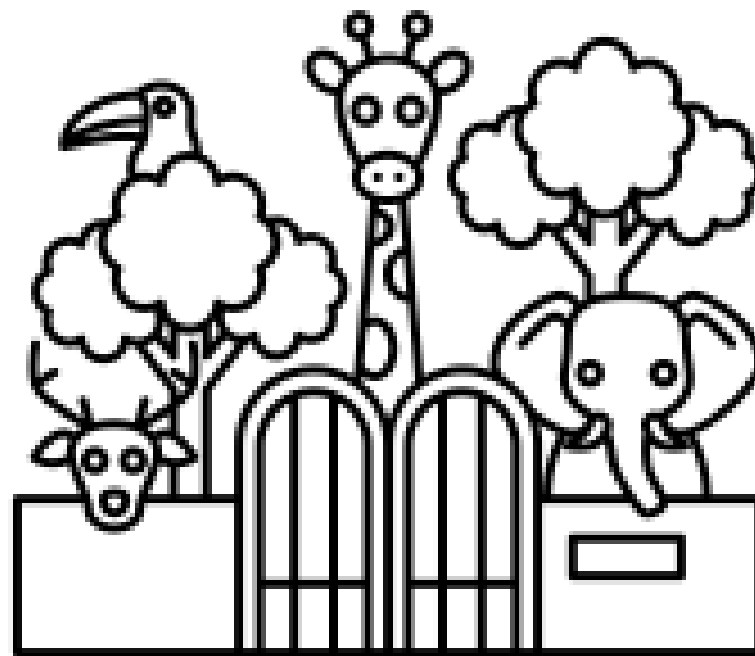
Multiplying all conditionals evaluates the probability of a full joint configuration of pixels.

Per-pixel classification vs. Autoregressive

- Image colorization: per-pixel classification loss
- PixelCNN, VQ-VAE2: autoregressive model
- Key idea: only produce discrete representations (e.g., VQ codebook, 313 ab color bins)
- Differences: Independent vs. sequential prediction



Thank You!



16-726 Learning-based Image Synthesis

<https://learning-image-synthesis.github.io/>