



10-423/10-623 Generative AI

Machine Learning Department
School of Computer Science
Carnegie Mellon University

Variational Inference + Variational Autoencoders (VAEs)

Matt Gormley & Henry Chai

Lecture 7

Sep. 18, 2024

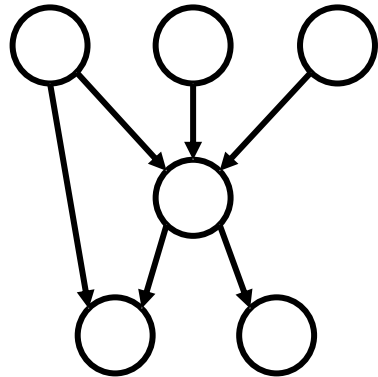
Reminders

- **Homework 1: Generative Models of Text**
 - **Out: Mon, Sep 9**
 - **Due: Mon, Sep 23 at 11:59pm**

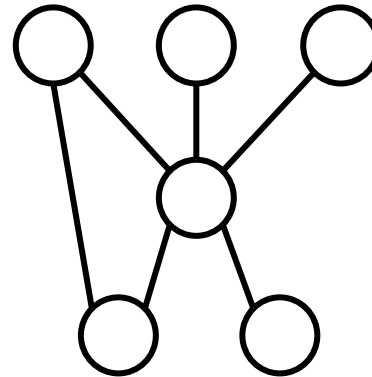
DIRECTED GRAPHICAL MODEL

Three Types of Graphical Models

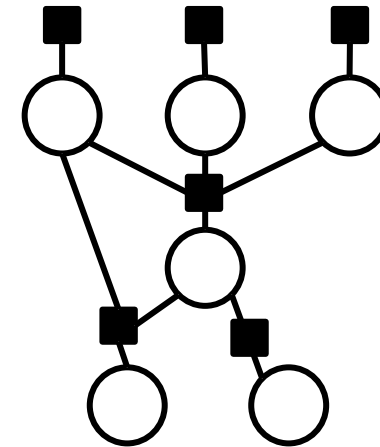
Directed Graphical Model



Undirected Graphical Model



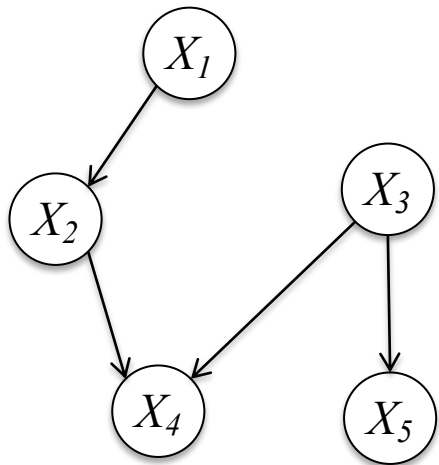
Factor Graph



Directed Graphical Model

$$P(X_1, X_2, X_3) = P(X_1) P(X_2 | X_1) P(X_3 | X_2, X_1)$$

Example



$$P(X_1, X_2, X_3, X_4, X_5) = \\ P(X_5 | X_3) P(X_4 | X_2, X_3) \\ P(X_3) P(X_2 | X_1) P(X_1)$$

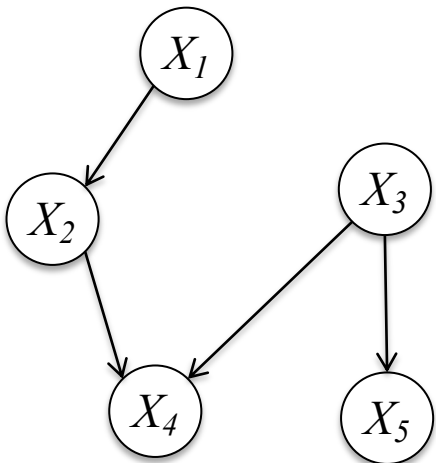
Definition

- A directed graphical model (aka. Bayesian network) is a **directed acyclic graph** that represents the conditional independencies of a set of variables X_1, \dots, X_T
- Each **node** is variable X_t and each **edge** implies a directional influence between a pair of variables
- The DGM factorizes the joint distribution over the variables as a product of conditional probabilities:

$$P(X_1, \dots, X_T) = \prod_{t=1}^T P(X_t | \text{parents}(X_t))$$

Directed Graphical Model

Example



the graph (*qualitative specification*) could be:

- specified using domain expertise about causal relationships
- learned from data
- chosen because of nice computational properties

$$P(X_1, X_2, X_3, X_4, X_5) =$$

$$P(X_5 | X_3)P(X_4 | X_2, X_3)$$

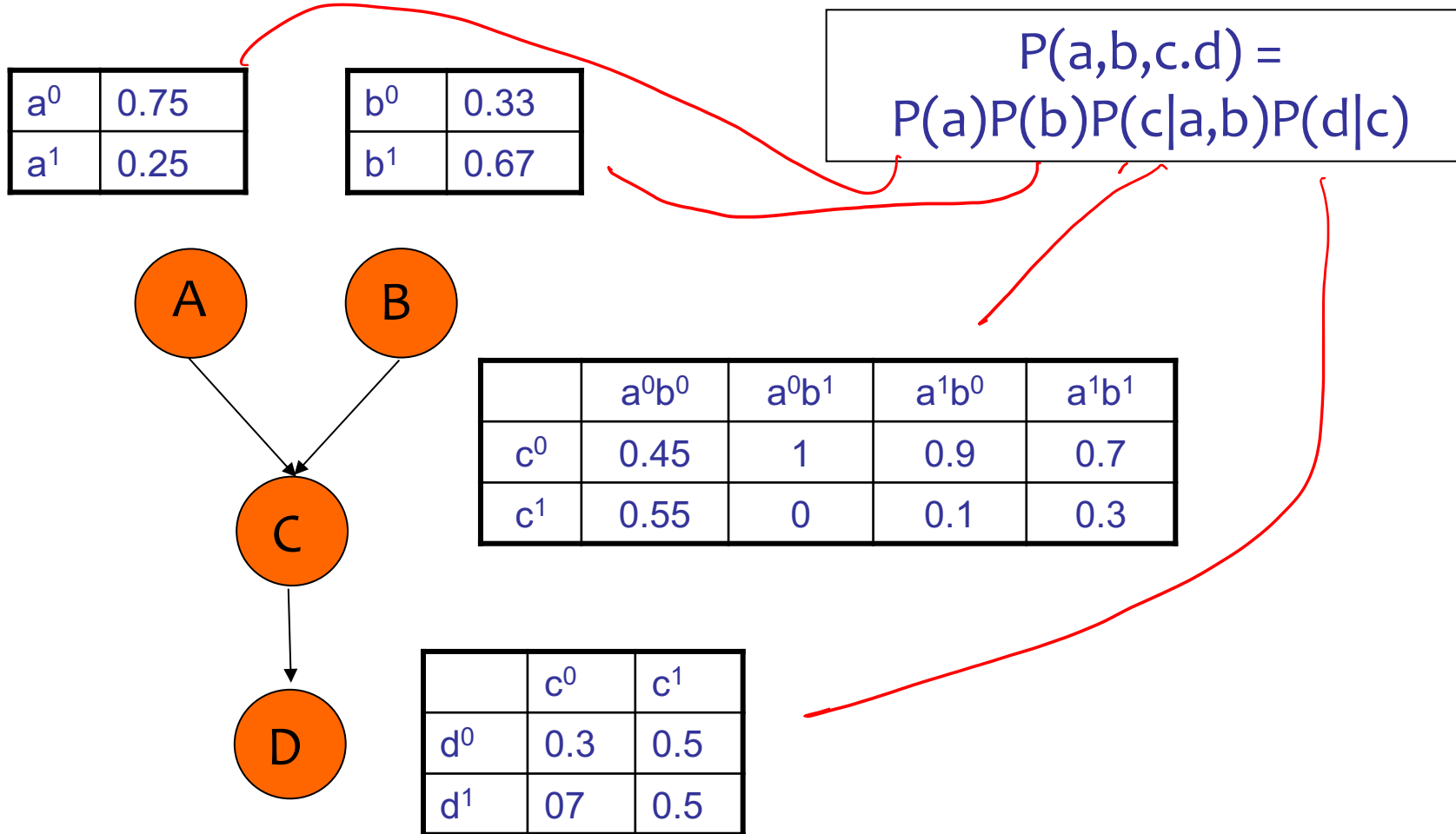
$$P(X_3)P(X_2 | X_1) \mathcal{P}(x_1)$$

the conditional probabilities (*quantitative specification*) is:

- depends on the types of variables involved
- typically learned from data

Quantitative Specification

Example: Conditional probability tables (CPTs)
for discrete random variables

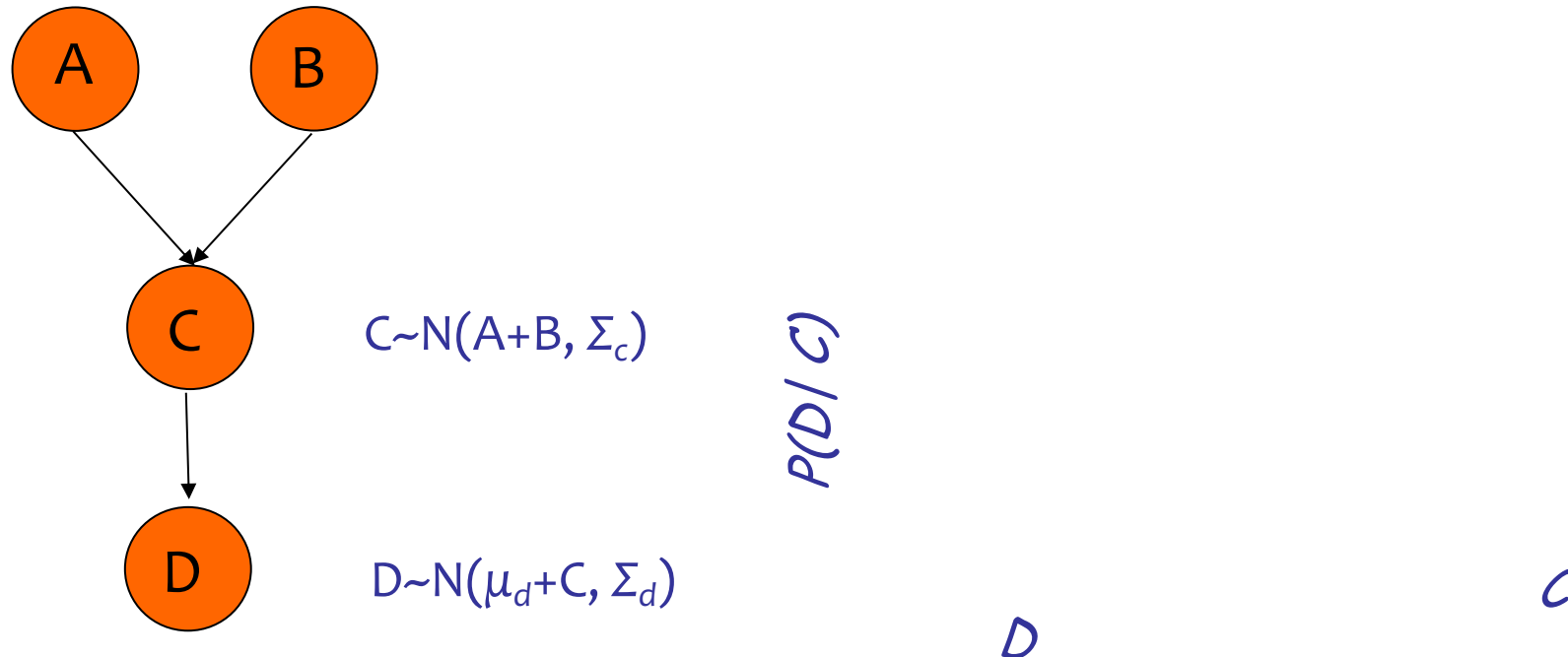


Quantitative Specification

Example: Conditional probability density functions (CPDs)
for continuous random variables

$$A \sim N(\mu_a, \Sigma_a) \quad B \sim N(\mu_b, \Sigma_b)$$

$$P(a,b,c,d) = P(a)P(b)P(c|a,b)P(d|c)$$



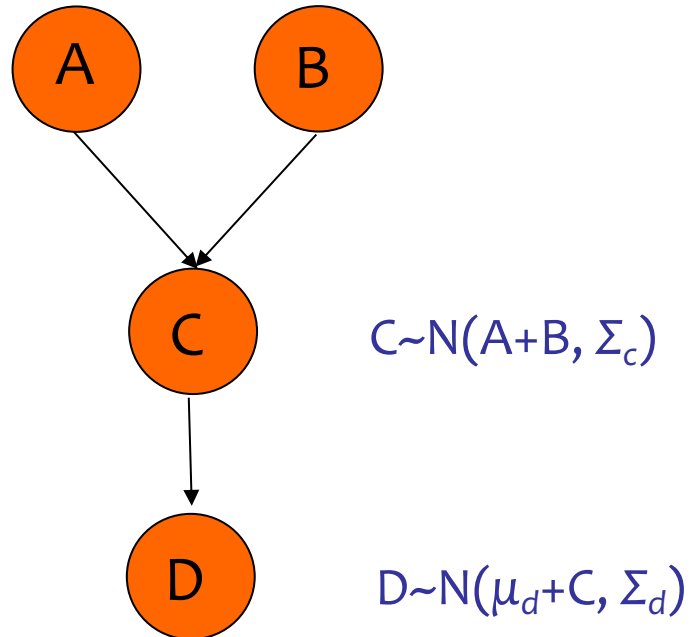
Quantitative Specification

Example: Combination of CPTs and CPDs
for a mix of discrete and continuous variables

a^0	0.75
a^1	0.25

b^0	0.33
b^1	0.67

$$P(a,b,c,d) = P(a)P(b)P(c|a,b)P(d|c)$$

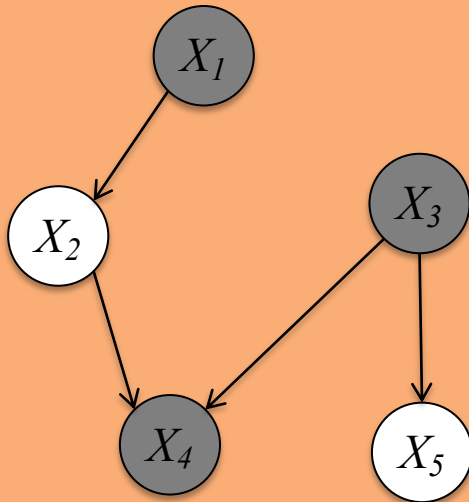


Observed Variables

- In a graphical model, **shaded nodes** are “**observed**”, i.e. their values are given

Example:

$$P(X_2, X_5 \mid X_1 = 0, X_3 = 1, X_4 = 1)$$



MARKOV MODEL

Markov Model

1. **1st-order Markov assumption:**

for a sequence of random variables, the probability distribution over x_t random variables is conditionally independent of x_1, \dots, x_{t-2} given x_{t-1}

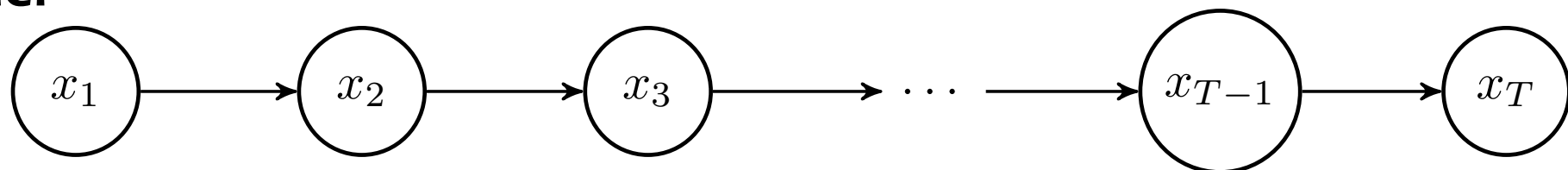
$$p(x_t \mid x_1, \dots, x_{t-1}) = p(x_t \mid x_{t-1})$$

2. **1st-order Markov model:**

defines a joint distribution over a sequence of variables using a Markov assumption

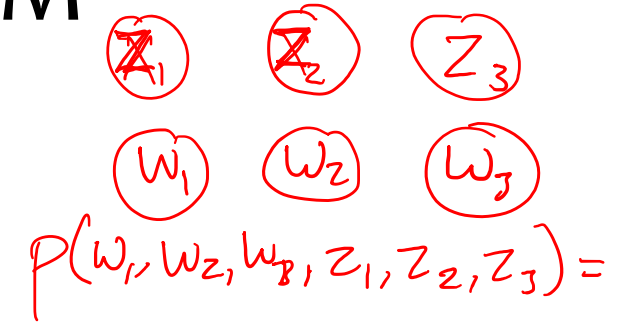
$$p(x_1, \dots, x_T) = p(x_1) \prod_{t=2}^T p(x_t \mid x_{t-1})$$

3. We can represent the Markov model as a **directed graphical model**

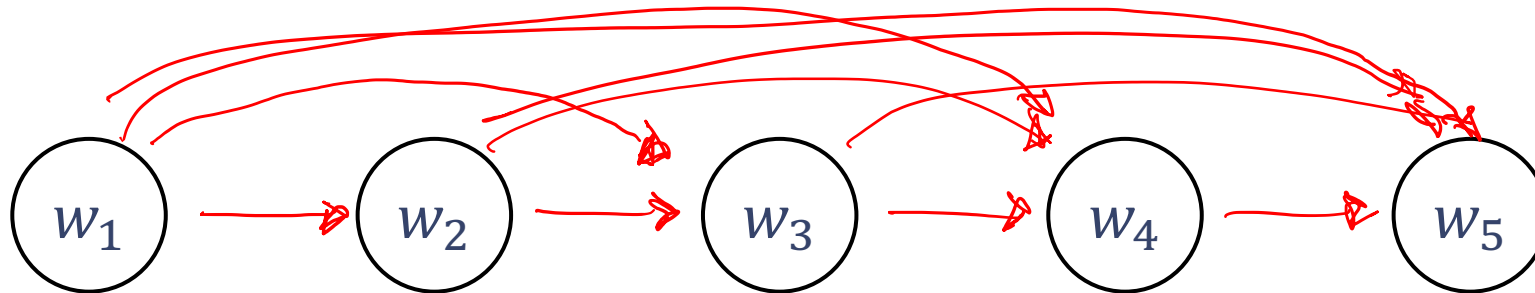


In-class Exercise: RNN as a DGM

Given a five-word sequence, $[w_1, w_2, w_3, w_4, w_5]$, how could we represent the implied probability distributions of an RNN as a directed graphical model?



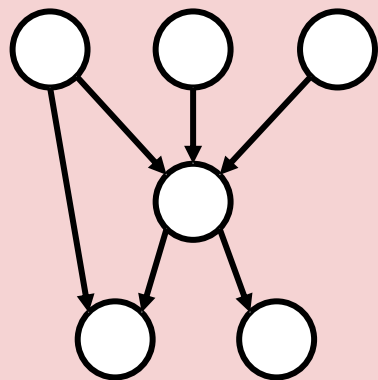
$$p(w_1, w_2, w_3, w_4, w_5) = p(w_1) p(w_2|w_1) p(w_3|w_2, w_1) \dots p(w_5|w_1, \dots, w_4)$$



Locally Normalized

autoregressive LM

Directed Graphical Model



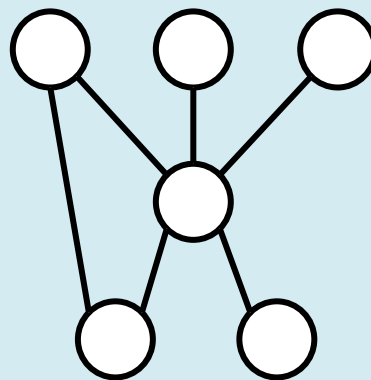
$$P(X_1, \dots, X_T) = \prod_{t=1}^T P(X_t \mid \text{parents}(X_t))$$

vs.

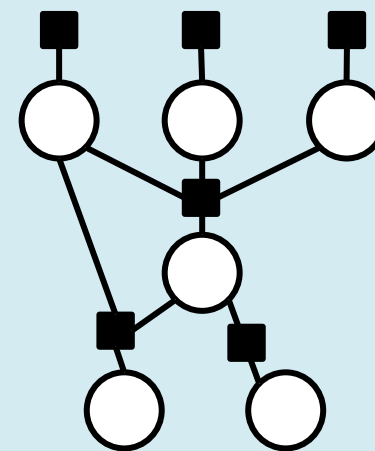
Globally Normalized

encoder-only Transformer

Undirected Graphical Model



Factor Graph



$$p(\mathbf{x}) = \frac{1}{Z} \prod_{\alpha} \psi_{\alpha}(\mathbf{x}_{\alpha})$$

UNSUPERVISED LEARNING

Unsupervised Learning

Assumptions:

1. our data comes from some distribution $p^*(\mathbf{x}_o)$
2. we choose a distribution $p_\theta(\mathbf{x}_o)$ for which sampling $\mathbf{x}_o \sim p_\theta(\mathbf{x}_o)$ is tractable

Goal: learn θ s.t. $p_\theta(\mathbf{x}_o) \approx p^*(\mathbf{x}_o)$

Unsupervised Learning

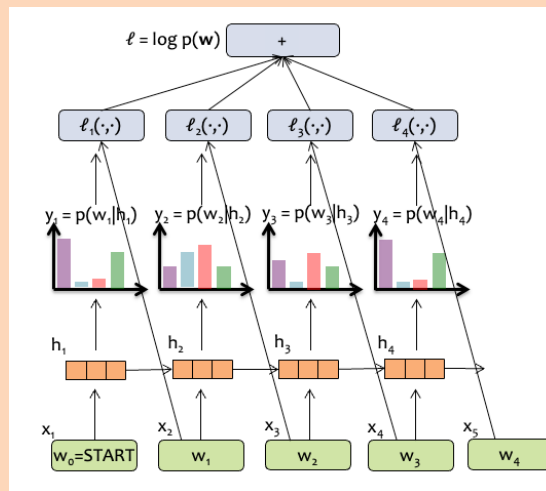
Assumptions:

1. our data comes from some distribution $p^*(\mathbf{x}_0)$
2. we choose a distribution $p_\theta(\mathbf{x}_0)$ for which sampling $x_0 \sim p_\theta(\mathbf{x}_0)$ is tractable

Goal: learn θ s.t. $p_\theta(\mathbf{x}_0) \approx p^*(\mathbf{x}_0)$

Example: autoregressive LMs

- true $p^*(\mathbf{x}_0)$ is the (human) process that produced text on the web
- choose $p_\theta(\mathbf{x}_0)$ to be an autoregressive language model
 - autoregressive structure means that $p(\mathbf{x}_t | \mathbf{x}_1, \dots, \mathbf{x}_{t-1}) \sim \text{Categorical}(\cdot)$ and ancestral sampling is exact/efficient
- learn by finding $\theta \approx \operatorname{argmax}_\theta \log(p_\theta(\mathbf{x}_0))$ using gradient based updates on $\nabla_\theta \log(p_\theta(\mathbf{x}_0))$



Unsupervised Learning

Assumptions:

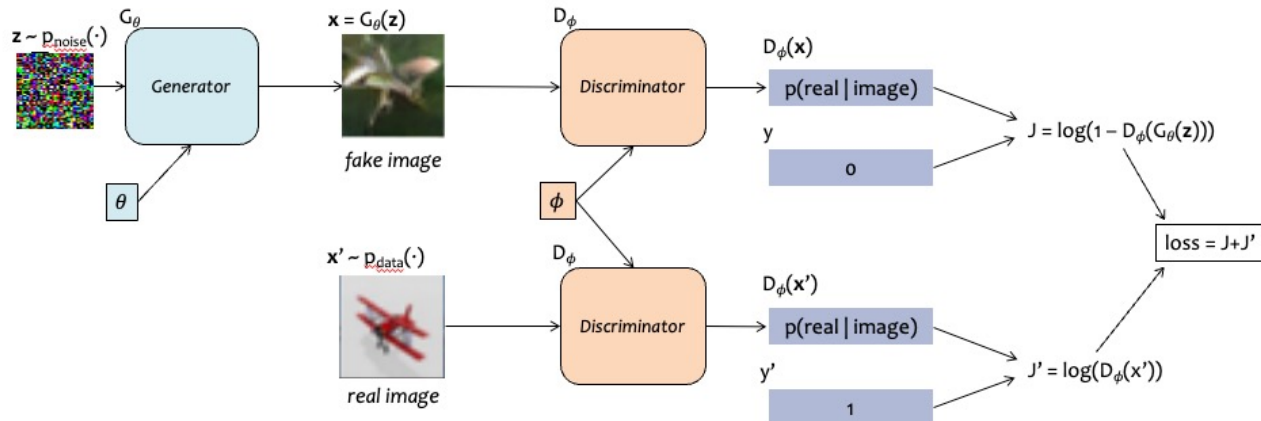
1. our data comes from some distribution $p^*(\mathbf{x}_0)$
2. we choose a distribution $p_\theta(\mathbf{x}_0)$ for which sampling $x_0 \sim p_\theta(\mathbf{x}_0)$ is tractable

Goal: learn θ s.t. $p_\theta(\mathbf{x}_0) \approx p^*(\mathbf{x}_0)$

Example: GANs

- true $p^*(\mathbf{x}_0)$ is distribution over photos taken and posted to Flickr
- choose $p_\theta(\mathbf{x}_0)$ to be an expressive model (e.g. noise fed into inverted CNN) that can generate images
 - sampling is typically easy:
 $\mathbf{z} \sim N(\mathbf{0}, \mathbf{I})$ and $\mathbf{x}_0 = f_\theta(\mathbf{z})$
- learn by finding $\theta \approx \operatorname{argmax}_\theta \log(p_\theta(\mathbf{x}_0))$
 - No! Because we can't even compute $\log(p_\theta(\mathbf{x}_0))$ or its gradient
 - Why not? Because the integral is intractable even for a simple 1-hidden layer neural network with nonlinear activation

$$p(\mathbf{x}_0) = \int_{\mathbf{z}} p(\mathbf{x}_0 | \mathbf{z}) p(\mathbf{z}) d\mathbf{z}$$



so optimize a minimax loss instead

Unsupervised Learning

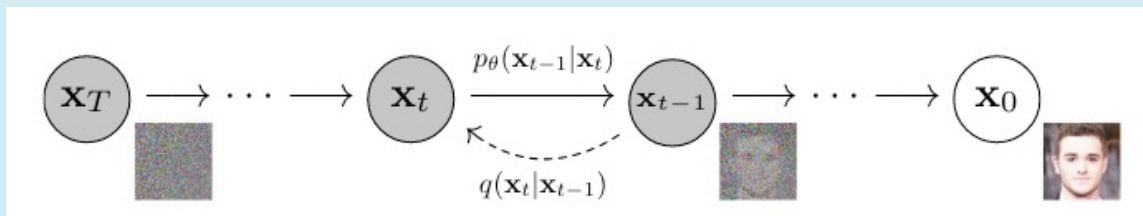
Assumptions:

1. our data comes from some distribution $p^*(\mathbf{x}_0)$
2. we choose a distribution $p_\theta(\mathbf{x}_0)$ for which sampling $\mathbf{x}_0 \sim p_\theta(\mathbf{x}_0)$ is tractable

Goal: learn θ s.t. $p_\theta(\mathbf{x}_0) \approx p^*(\mathbf{x}_0)$

Example: VAEs / Diffusion Models

- true $p^*(\mathbf{x}_0)$ is distribution over photos taken and posted to Flickr
- choose $p_\theta(\mathbf{x}_0)$ to be an expressive model (e.g. noise fed into inverted CNN) that can generate images
 - sampling is will be easy
- learn by finding $\theta \approx \operatorname{argmax}_\theta \log(p_\theta(\mathbf{x}_0))$?
 - Sort of! We can't compute the gradient $\nabla_\theta \log(p_\theta(\mathbf{x}_0))$
 - So we instead optimize a variational lower bound (more on that later)



Latent Variable Models

- For GANs and VAEs, we assume that there are (unknown) **latent variables** which give rise to our observations
- The **vector z** are those latent variables
- After learning a GAN or VAE, we can **interpolate** between images in latent z space



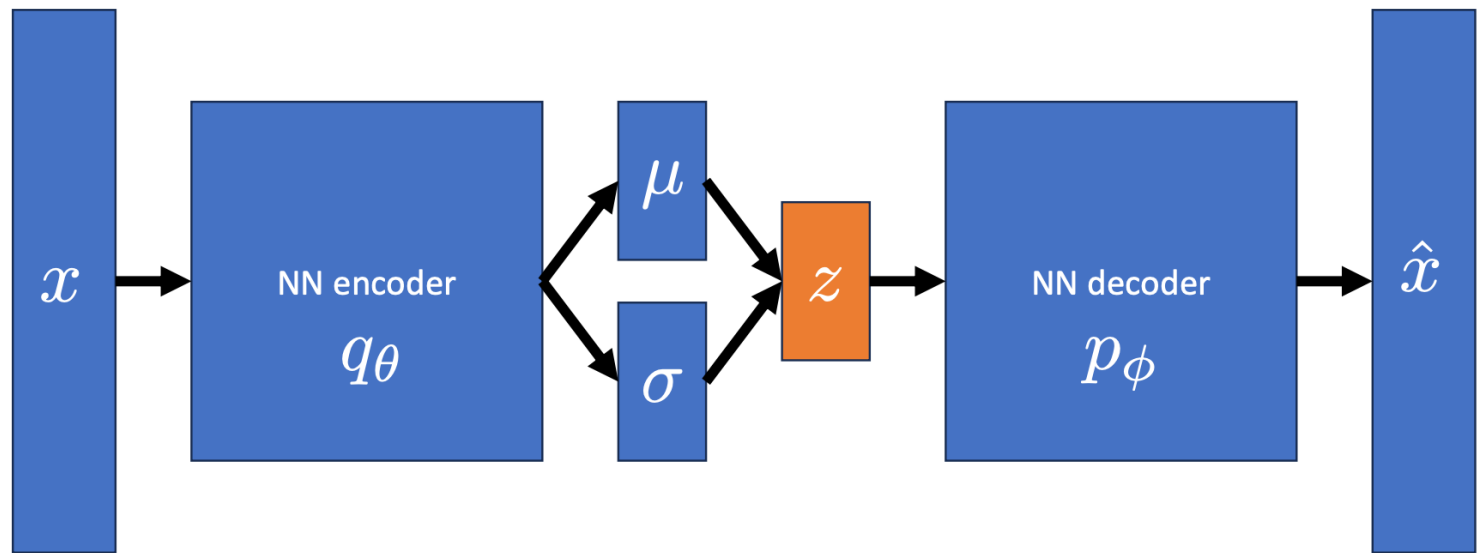
Figure 4: Top rows: Interpolation between a series of 9 random points in Z show that the space learned has smooth transitions, with every image in the space plausibly looking like a bedroom. In the 6th row, you see a room without a window slowly transforming into a room with a giant window. In the 10th row, you see what appears to be a TV slowly being transformed into a window.

From VAEs to Diffusion Models

- Next we will consider (1) **diffusion models** and (2) **variational autoencoders (VAEs)**
- The steps in defining these models is roughly:
 - Define a probability distribution involving Gaussian noise
 - Use a variational lower bound as an objective function
 - Learn the parameters of the probability distribution by optimizing the objective function
- So what is a variational lower bound?

HIGH-LEVEL INTRO TO VARIATIONAL INFERENCE

Variational Autoencoder: Network Perspective

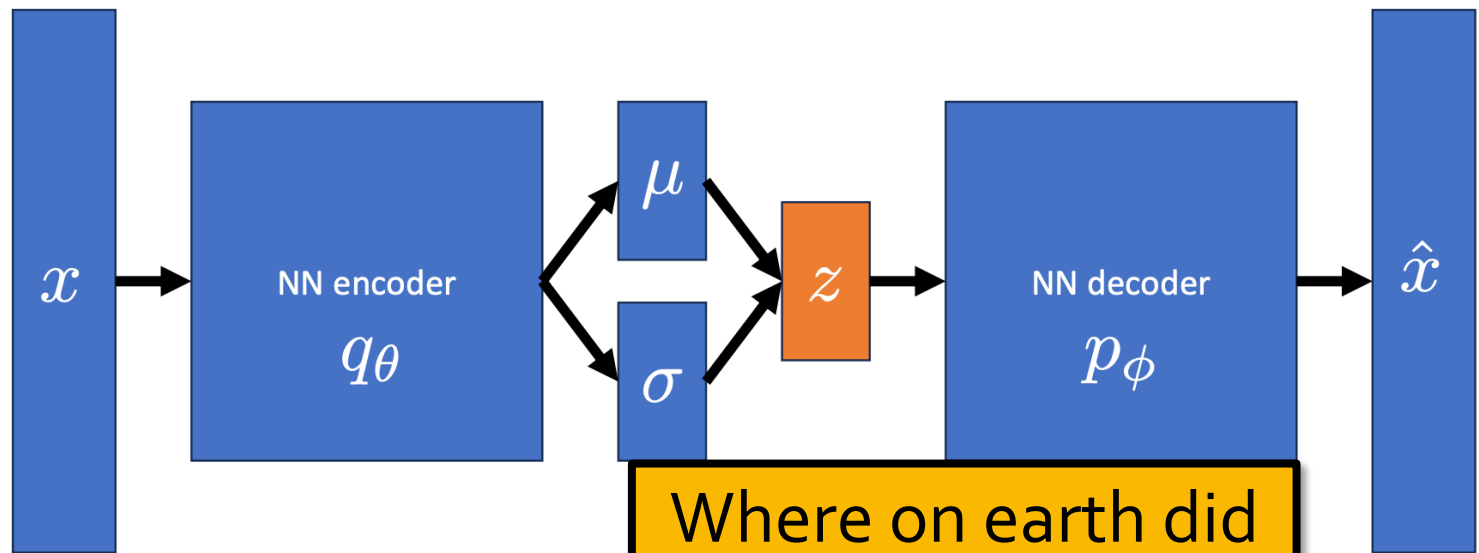


- Objective: minimize the negative log-likelihood of the dataset plus a *regularization term* that encourages a dense latent space

$$J(\theta, \phi) = \sum_{i=1}^N \ell_i(\theta, \phi)$$

$$\ell_i(\theta, \phi) = -\mathbb{E}_{q_{\theta}(z|x^{(i)})}[\log p_{\phi}(x^{(i)}|z)] + KL(q_{\theta}(z|x^{(i)}) \parallel p(z))$$

Variational Autoencoder: Network Perspective



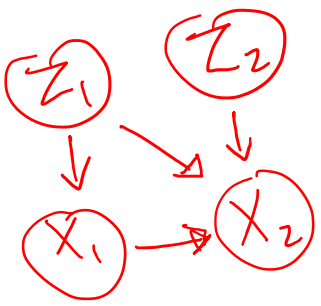
- Objective: minimize the negative log-likelihood of the dataset plus a *regularization term* that

$$J(\theta, \phi) = \sum_{i=1}^N \ell_i(\theta, \phi)$$

$$\ell_i(\theta, \phi) = -\mathbb{E}_{q_\theta(z|x^{(i)})} [\log p_\phi(x^{(i)}|z)] + KL(q_\theta(z|x^{(i)}) \parallel p(z))$$

Where on earth did
this objective
function come
from?

dataset
latent space

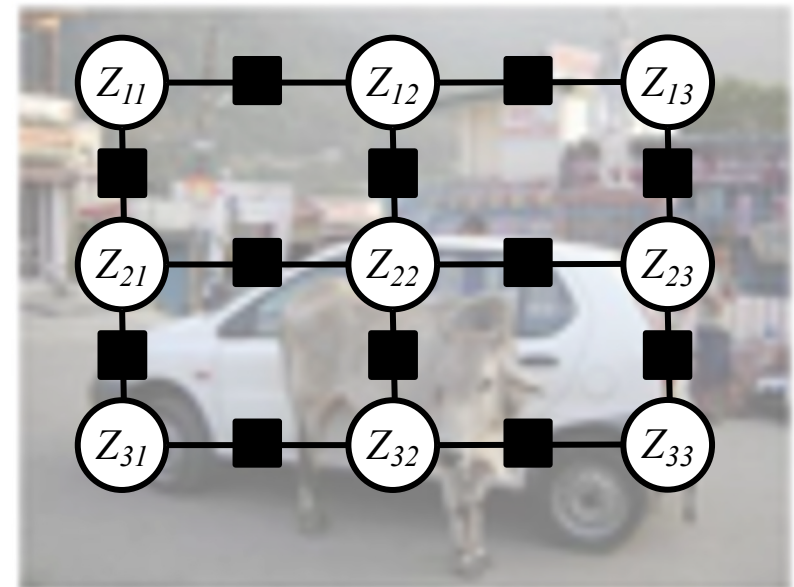
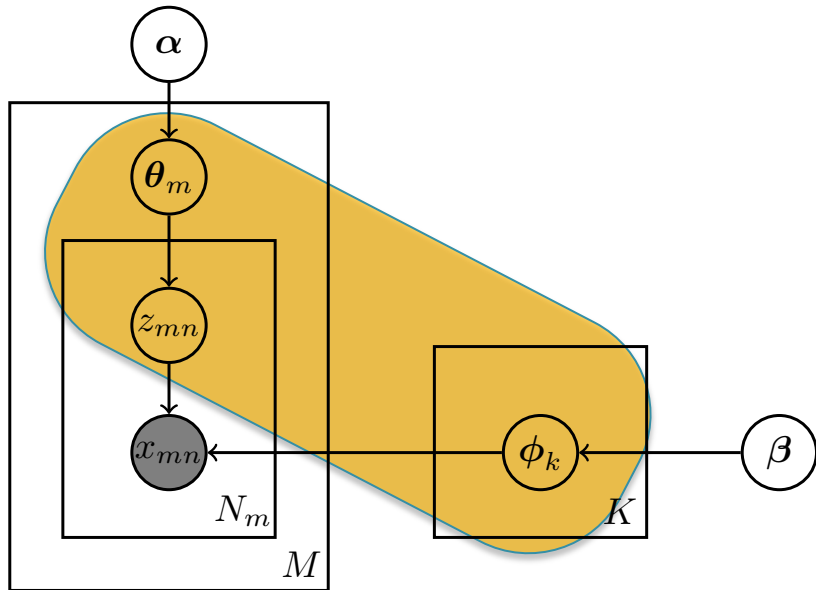


Variational Inference

$$p(x_1, x_2, z_1, z_2) = p(z_1) p(z_2) p(x_1 | z_1) p(x_2 | x_1, z_1, z_2)$$

A Common Problem:

- Suppose we have an interesting distribution $p(\mathbf{x}, \mathbf{z})$ and we wish to work with its posterior $p(\mathbf{z} | \mathbf{x})$
- For training data \mathbf{x} and latent variables \mathbf{z} , estimating the posterior $p(\mathbf{z} | \mathbf{x})$ is usually intractable!



Variational Inference

A Common Problem:

- Suppose we have an interesting distribution $p(\mathbf{x}, \mathbf{z})$ and we wish to work with its posterior $p(\mathbf{z} | \mathbf{x})$ or the marginal $p(\mathbf{x})$
- For training data \mathbf{x} and latent variables \mathbf{z} , estimating the posterior $p(\mathbf{z} | \mathbf{x})$ or the marginal $p(\mathbf{x})$ is usually intractable!

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

we assume
 $p(\mathbf{x})$ is
intractable

Question: Why is $p(\mathbf{x})$ often intractable to compute?

Answer:

$$p(\vec{x}) = \int_{\vec{z}} p(\vec{x}, \vec{z}) d\vec{z}$$

$$p(\vec{x}) = \sum_{\vec{z}} p(\vec{x}, \vec{z})$$

Variational Inference

A Common Problem:

- Suppose we have an interesting distribution $p(\mathbf{x}, \mathbf{z})$ and we wish to work with its posterior $p(\mathbf{z} | \mathbf{x})$
- For training data \mathbf{x} and latent variables \mathbf{z} , estimating the posterior $p(\mathbf{z} | \mathbf{x})$ is usually intractable!

Solution:

- Approximate $p(\mathbf{z} | \mathbf{x})$ with a simpler $q(\mathbf{z} | \mathbf{x})$
- Typically $q(\mathbf{z} | \mathbf{x})$ has more independence assumptions than $p(\mathbf{z} | \mathbf{x})$, which is fine b/c $q(\mathbf{z} | \mathbf{x})$ is tuned for a specific \mathbf{x}
- **Key idea:** pick a single $q(\mathbf{z} | \mathbf{x})$ from some family Q that best approximates $p(\mathbf{z} | \mathbf{x})$

Variational Inference

Terminology:

- $q(\mathbf{z} \mid \mathbf{x})$: the **variational approximation**
- Q : the **variational family**
- Usually $q_{\theta}(\mathbf{z} \mid \mathbf{x})$ is parameterized by some θ called **variational parameters**
- Usually $p_{\alpha}(\mathbf{z} \mid \mathbf{x})$ is parameterized by some fixed α – we'll call them the parameters

Example Algorithms:

- mean-field variational inference
- loopy belief propagation
- tree-reweighted belief propagation
- expectation propagation

Variational Inference

- **Question:** Do we learn a single distribution $q_{\theta}(z | x)$ for all x 's?
- **Answer:** Not necessarily, it's quite common to infer a separate q_{θ} for each x !
 - Consider the sampling equivalent of this:
 - you could draw samples $z^{(i)} \sim p(z | x')$
 - then train some simple $q_{\theta}(z | x')$ on $z^{(1)}, z^{(2)}, \dots, z^{(N)}$
 - hope that the sample adequately represents the posterior for the given x'
 - How is VI different from this?
 - VI doesn't require sampling
 - VI is fast and deterministic
 - Why? b/c we choose an objective function (KL divergence) that defines which q_{θ} best approximates p_{α} , and exploit the special structure of q_{θ} to optimize it

Variational Inference

V.I. offers a new design decision

- Choose the distribution $p_\alpha(z | x)$ that you really want, i.e. don't just simplify it to make it computationally convenient
- Then design the structure of another distribution $q_\theta(z | x)$ such that V.I. is efficient

VAE

$$p(x, z) = p(x|z) p(z)$$

$$p(z) = \text{Gaussian}(0, I)$$

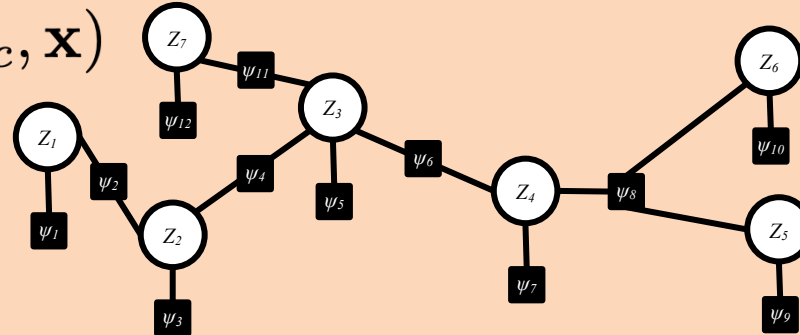
$$p(x|z) = \begin{cases} 1 & \text{if } x = \text{MLP}(z) \\ 0 & \text{otherwise} \end{cases}$$

THE MEAN FIELD APPROXIMATION

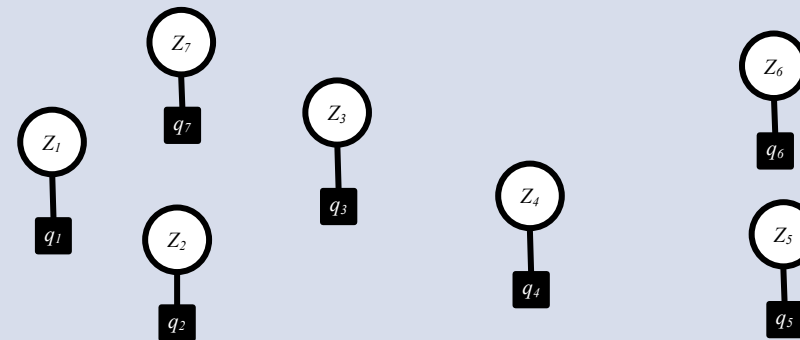
Mean Field Approximation

The **mean field approximation** assumes our variational approximation $q_{\theta}(\mathbf{z})$ treats each variable as independent

$$p_{\alpha}(\mathbf{z} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{c \in \mathcal{C}} \psi_c(\mathbf{z}_c, \mathbf{x})$$



$$q_{\theta}(\mathbf{z}) = \prod_{t=1}^T q_t(z_t)$$

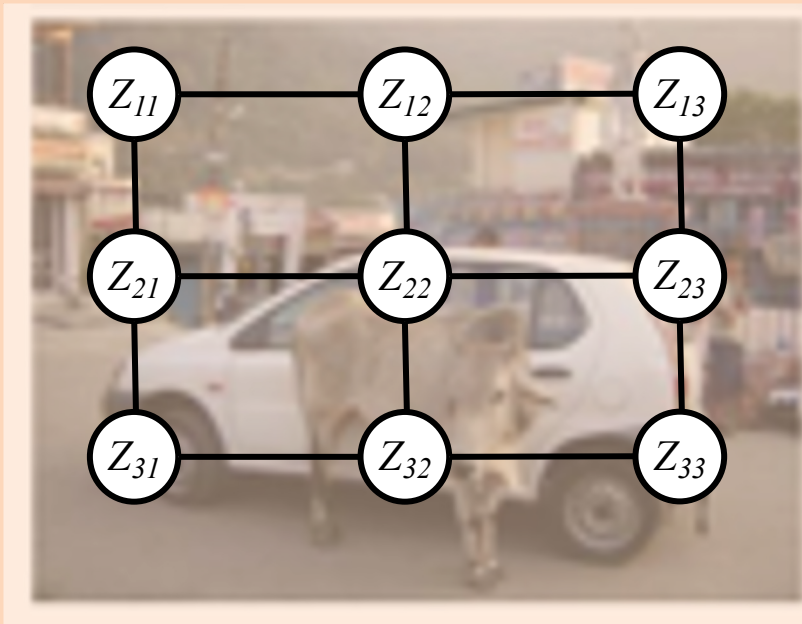


Mean Field Approximation

The **mean field approximation** assumes our variational approximation $q_{\theta}(\mathbf{z})$ treats each variable as independent

Ising Model

$$p_{\alpha}(\mathbf{z} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{c \in \mathcal{C}} \psi_c(\mathbf{z}_c, \mathbf{x})$$



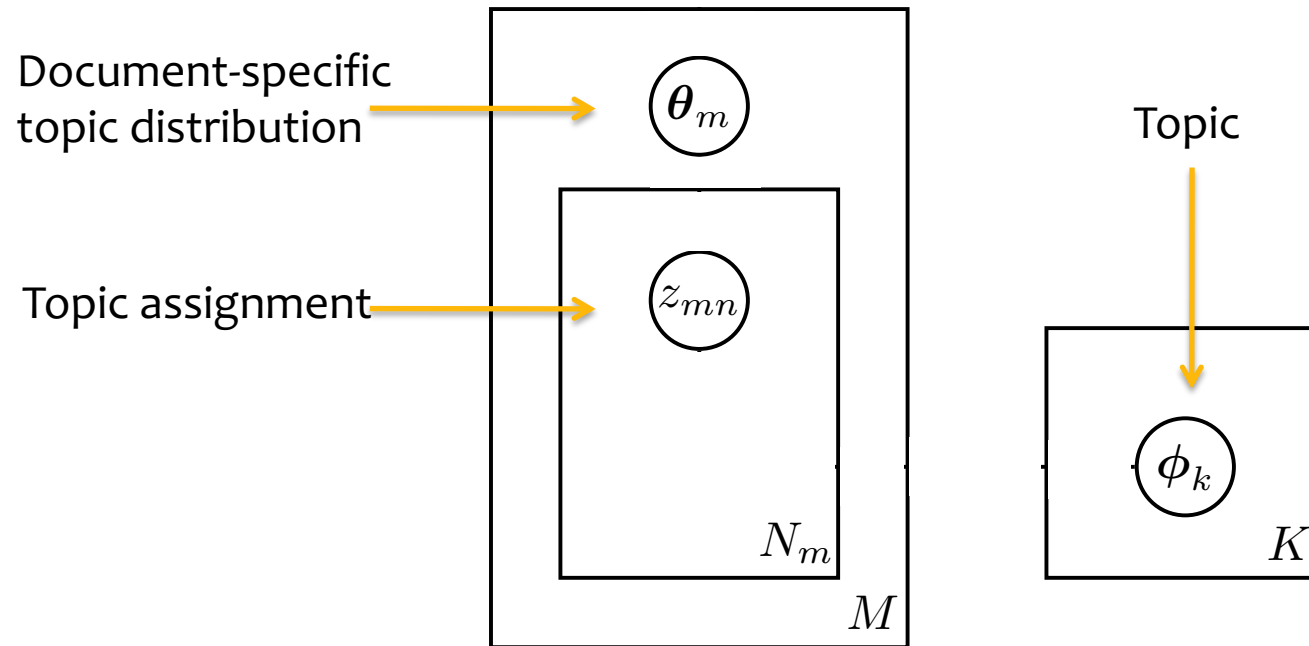
$$q_{\theta}(\mathbf{z}) = \prod_{t=1}^T q_t(z_t)$$



Mean Field Approximation

Latent Dirichlet Allocation (LDA)

- Uncollapsed Variational Inference, aka. Explicit V.I. (mean field variational approximation)



MEAN FIELD VARIATIONAL INFERENCE

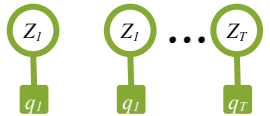
KL Divergence

- Definition: for two distributions $q(x)$ and $p(x)$ over $x \in \mathcal{X}$, the **KL Divergence** is:

$$\text{KL}(q||p) = E_{q(x)} \left[\log \frac{q(x)}{p(x)} \right] = \begin{cases} \sum_x q(x) \log \frac{q(x)}{p(x)} \\ \int_x q(x) \log \frac{q(x)}{p(x)} dx \end{cases}$$

- Properties:
 - $\text{KL}(q || p)$ measures the **proximity** of two distributions q and p
 - KL is **not** symmetric: $\text{KL}(q || p) \neq \text{KL}(p || q)$
 - KL is minimized when $q(x) = p(x)$ for all $x \in \mathcal{X}$
 - $\text{KL}(q || p) \geq 0$

Mean Field V.I. Overview

1. Goal: estimate $p_\alpha(\mathbf{z} \mid \mathbf{x})$
we assume this is intractable to compute exactly
2. Idea: approximate with another distribution $q_\theta(\mathbf{z} \mid \mathbf{x}) \approx p_\alpha(\mathbf{z} \mid \mathbf{x})$ for each \mathbf{x}
3. Mean Field: assume $q_\theta(\mathbf{z} \mid \mathbf{x}) = \prod_t q_t(z_t \mid \mathbf{x}; \theta)$
i.e., we decompose over variables
other choices for the decomposition of $q_\theta(\mathbf{z})$ give rise to “structured mean field”

4. Optimization Problem: pick the q that minimizes $\text{KL}(q \parallel p)$

$$\hat{q}(\mathbf{z} \mid \mathbf{x}) = \underset{q(\mathbf{z} \mid \mathbf{x}) \in \mathcal{Q}}{\text{argmin}} \text{KL}(q(\mathbf{z} \mid \mathbf{x}) \parallel p(\mathbf{z} \mid \mathbf{x}))$$

$$\hat{\theta} = \underset{\theta \in \Theta}{\text{argmin}} \text{KL}(q_\theta(\mathbf{z} \mid \mathbf{x}) \parallel p_\alpha(\mathbf{z} \mid \mathbf{x}))$$

equivalent

5. Optimization Algorithm: various options
 - e.g. coordinate descent repeatedly picks the best $q_t(z_t \mid \mathbf{x})$ based on the other $\{q_s(z_s \mid \mathbf{x})\}_{s \neq t}$ being fixed
 - e.g. gradient descent optimizes a surrogate objective $\text{ELBO}(q_\theta)$ to find θ

Optimizing KL Divergence

- Question: How do we minimize KL?

$$\hat{\theta} = \operatorname{argmin}_{\theta \in \Theta} \underbrace{\text{KL}(q_{\theta}(\mathbf{z} | \mathbf{x}) \parallel p_{\alpha}(\mathbf{z} | \mathbf{x}))}$$

- Answer #1: Oh no! We can't even compute this KL.

Why we can't compute KL...

$$\begin{aligned} \text{KL}(q(\mathbf{z} | \mathbf{x}) \parallel p(\mathbf{z} | \mathbf{x})) &= E_{q(\mathbf{z}|\mathbf{x})} \left[\log \left(\frac{q(\mathbf{z} | \mathbf{x})}{p(\mathbf{z} | \mathbf{x})} \right) \right] \\ &= E_{q(\mathbf{z}|\mathbf{x})} [\log q(\mathbf{z} | \mathbf{x})] - E_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{z} | \mathbf{x})] \\ &= E_{q(\mathbf{z}|\mathbf{x})} [\log q(\mathbf{z} | \mathbf{x})] - E_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}, \mathbf{z})] + \boxed{E_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x})]} \\ &= E_{q(\mathbf{z}|\mathbf{x})} [\log q(\mathbf{z} | \mathbf{x})] - E_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}, \mathbf{z})] + \underbrace{\log p(\mathbf{x})} \end{aligned}$$

this
expectation
does not
depend on q

we assumed this
is intractable to
compute!

$= p(\mathbf{x}, \mathbf{z}) / p(\mathbf{x})$

Optimizing KL Divergence

- Question: How do we minimize KL?

$$\hat{\theta} = \operatorname{argmin}_{\theta \in \Theta} \underbrace{\text{KL}(q_{\theta}(\mathbf{z} | \mathbf{x}) \parallel p_{\alpha}(\mathbf{z} | \mathbf{x}))}$$

- Answer #2: We don't need to compute this KL
We can instead maximize the ELBO (i.e. **E**vidence **L**ower **B**ound)

$$\text{ELBO}(q_{\theta}) = E_{q_{\theta}(\mathbf{z}|\mathbf{x})} [\log p_{\alpha}(\mathbf{x}, \mathbf{z})] - E_{q_{\theta}(\mathbf{z}|\mathbf{x})} [\log q_{\theta}(\mathbf{z} | \mathbf{x})]$$

The ELBO for a DGM

Here is why...

$$\begin{aligned} \theta &= \operatorname{argmin}_{\theta} \text{KL}(q_{\theta}(\mathbf{z} | \mathbf{x}) \parallel p_{\alpha}(\mathbf{z} | \mathbf{x})) \\ &= \operatorname{argmin}_{\theta} E_{q_{\theta}(\mathbf{z}|\mathbf{x})} [\log q_{\theta}(\mathbf{z} | \mathbf{x})] - E_{q_{\theta}(\mathbf{z}|\mathbf{x})} [\log p_{\alpha}(\mathbf{x}, \mathbf{z})] + \log p_{\alpha}(\mathbf{x}) \\ &= \operatorname{argmin}_{\theta} E_{q_{\theta}(\mathbf{z}|\mathbf{x})} [\log q_{\theta}(\mathbf{z} | \mathbf{x})] - E_{q_{\theta}(\mathbf{z}|\mathbf{x})} [\log p_{\alpha}(\mathbf{x}, \mathbf{z})] \\ &= \operatorname{argmax}_{\theta} \text{ELBO}(q_{\theta}) \end{aligned}$$

dropping the
intractable term
gives the ELBO

ELBO as Objective Function

What does maximizing $\text{ELBO}(q_\theta)$ accomplish?

$$\text{ELBO}(q_\theta) = E_{q_\theta(\mathbf{z}|\mathbf{x})} [\log p_\alpha(\mathbf{x}, \mathbf{z})] - E_{q_\theta(\mathbf{z}|\mathbf{x})} [\log q_\theta(\mathbf{z} | \mathbf{x})]$$

1. The first expectation is high if q_θ puts probability mass on the same values of \mathbf{z} that p_α puts probability mass

2. The second term is the entropy of q_θ and the entropy will be high if q_θ spreads its probability mass evenly

ELBO as lower bound

Theorem: For any q , $\log p(\mathbf{x}) \geq \text{ELBO}(q)$
i.e. $\text{ELBO}(q)$ is a lower bound for $\log p(\mathbf{x})$

Note:

$$\text{ELBO}(q_\theta) = E_{q_\theta(\mathbf{z}|\mathbf{x})} [\log p_\alpha(\mathbf{x}, \mathbf{z})] - E_{q_\theta(\mathbf{z}|\mathbf{x})} [\log q_\theta(\mathbf{z} | \mathbf{x})]$$

$$\text{KL}(q(\mathbf{z} | \mathbf{x}) || p(\mathbf{z} | \mathbf{x})) = E_{q(\mathbf{z}|\mathbf{x})} [\log q(\mathbf{z} | \mathbf{x})] - E_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}, \mathbf{z})] + E_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x})]$$

Proof #1:

1. $\log p(x) = \text{KL}(q||p) + \text{ELBO}(q)$
2. $\text{KL}(q||p) \geq 0$
3. $\log(p(x)) \geq \text{ELBO}(q)$

Takeaways:

1. in variational inference, we find the q that gives the **tightest bound** on the normalization constant for $p(\mathbf{z} | \mathbf{x})$
2. maximizing the ELBO is equivalent to minimizing KL
3. maximizing the ELBO is maximizing a lower bound on the likelihood $p(\mathbf{x})$

ELBO's relation to $\log p(x)$

Theorem:

for any q , $\log p(x) \geq \text{ELBO}(q)$
i.e. $\text{ELBO}(q)$ is a lower bound on $\log p(x)$

Proof #2:

Recall Jensen's Inequality: $f(E[x]) \geq E[f(x)]$, for concave f

$$\begin{aligned}\log p(x) &= \log \int_{\mathbf{z}} p(x, \mathbf{z}) d\mathbf{z} \quad (\text{marginal}) \\ &= \log \int_{\mathbf{z}} p(x, \mathbf{z}) \frac{q(\mathbf{z})}{q(\mathbf{z})} d\mathbf{z} \quad (\text{mult. by 1}) \\ &= \log E_{q(\mathbf{z})} \left[\frac{p(x, \mathbf{z})}{q(\mathbf{z})} \right] \quad (\text{def. of expectation}) \\ &\geq E_{q(\mathbf{z})} \left[\log \left(\frac{p(x, \mathbf{z})}{q(\mathbf{z})} \right) \right] \quad (\text{by Jensen's Ineq.}) \\ &= E_{q(\mathbf{z})} [\log p(x, \mathbf{z})] - E_{q(\mathbf{z})} [\log q(\mathbf{z})] = \text{ELBO}(q) \\ \Rightarrow \log p(x) &\geq \text{ELBO}(q)\end{aligned}$$

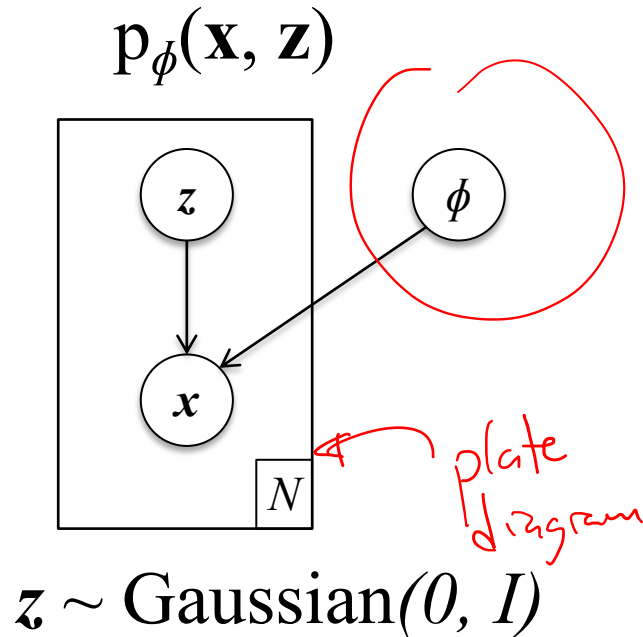
(Through the Lens of Variational Inference)

VARIATIONAL AUTOENCODERS

Why VAEs?

- **Autoencoders:**
 - learn a low dimensional representation of the input, but hard to work with as a generative model
 - one of the key limitations of autoencoders is that we have no way of **sampling** from them!
- **Variational autoencoders (VAEs)**
 - by contrast learn a continuous latent space that is **easy to sample from!**
 - can **generate** new data (e.g. images) by sampling from the learned generative model

Variational Autoencoders



The Something-like-a-VAE Model

- Consider a model $p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})$
 - where $p(\mathbf{z})$ is a $N(0, I)$
 - where $\mathbf{x} = \mathbf{z}/10 + \mathbf{z}/\|\mathbf{z}\| = g_{\phi}(\mathbf{z})$
i.e. we don't use parameters ϕ
- Trivially, we can draw samples of \mathbf{z} and directly convert them to values \mathbf{x}

The VAE Model

- The directed graphical model for VAE is the same as for the silly model above, and it's quite simple (ignoring the neural net details that give rise to \mathbf{x})
- **Key idea of VAE:** define $g_{\phi}(\mathbf{z})$ as a neural net and learn ϕ from data

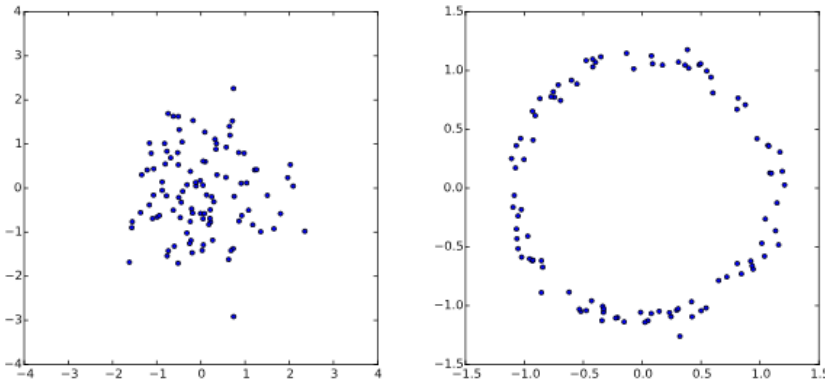


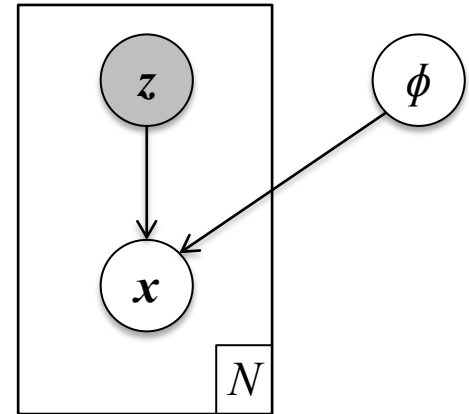
Figure from Doersch (2016)

Variational Autoencoders

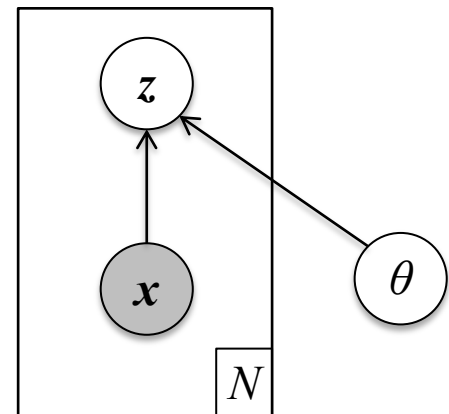
Neural Network Perspective

- We can view a variational autoencoder (VAE) as an autoencoder consisting of two neural networks
- VAEs (as encoders) define two distributions:
 - **encoder:** $q_{\theta}(\mathbf{z} | \mathbf{x})$
 - **decoder:** $p_{\phi}(\mathbf{x} | \mathbf{z})$
- Parameters θ and ϕ are neural network parameters (i.e. θ are not the variational parameters)

$$p_{\phi}(\mathbf{x} | \mathbf{z})$$



$$q_{\theta}(\mathbf{z} | \mathbf{x})$$

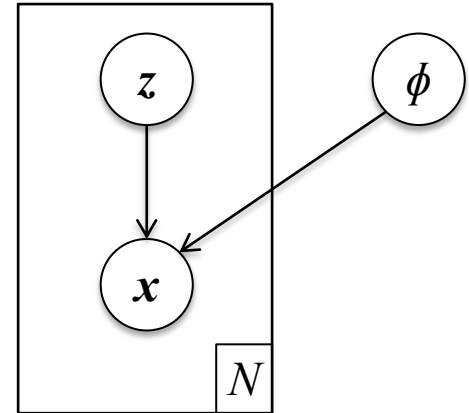


Variational Autoencoders

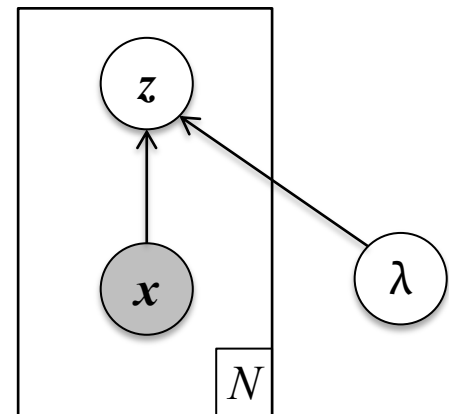
Graphical Model Perspective

- We can also view the VAE from the perspective of variational inference
- In this case we have two distributions:
 - **model:** $p_{\phi}(\mathbf{x}, \mathbf{z}) = p_{\phi}(\mathbf{x} | \mathbf{z}) p(\mathbf{z})$
 - **variational approximation:**
 $q_{\lambda=f(\mathbf{x}; \theta)}(\mathbf{z} | \mathbf{x})$
- We have the same model parameters ϕ
- The variational parameters λ are a function of NN parameters θ

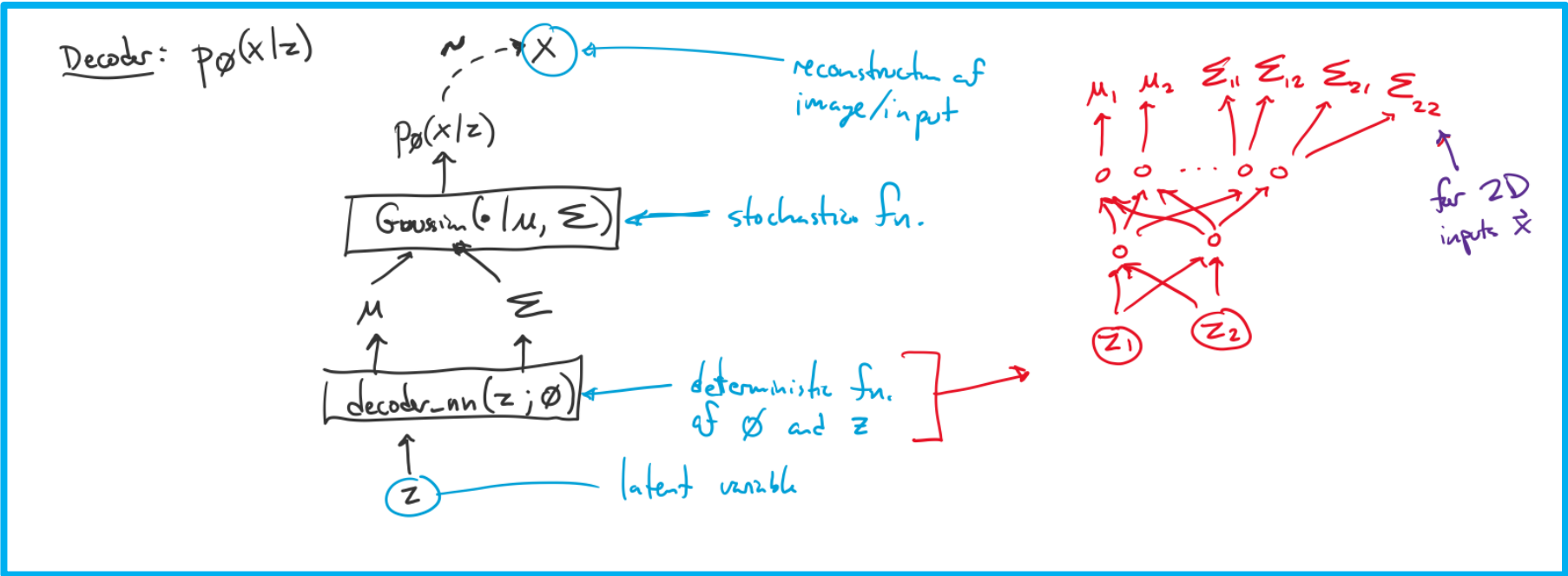
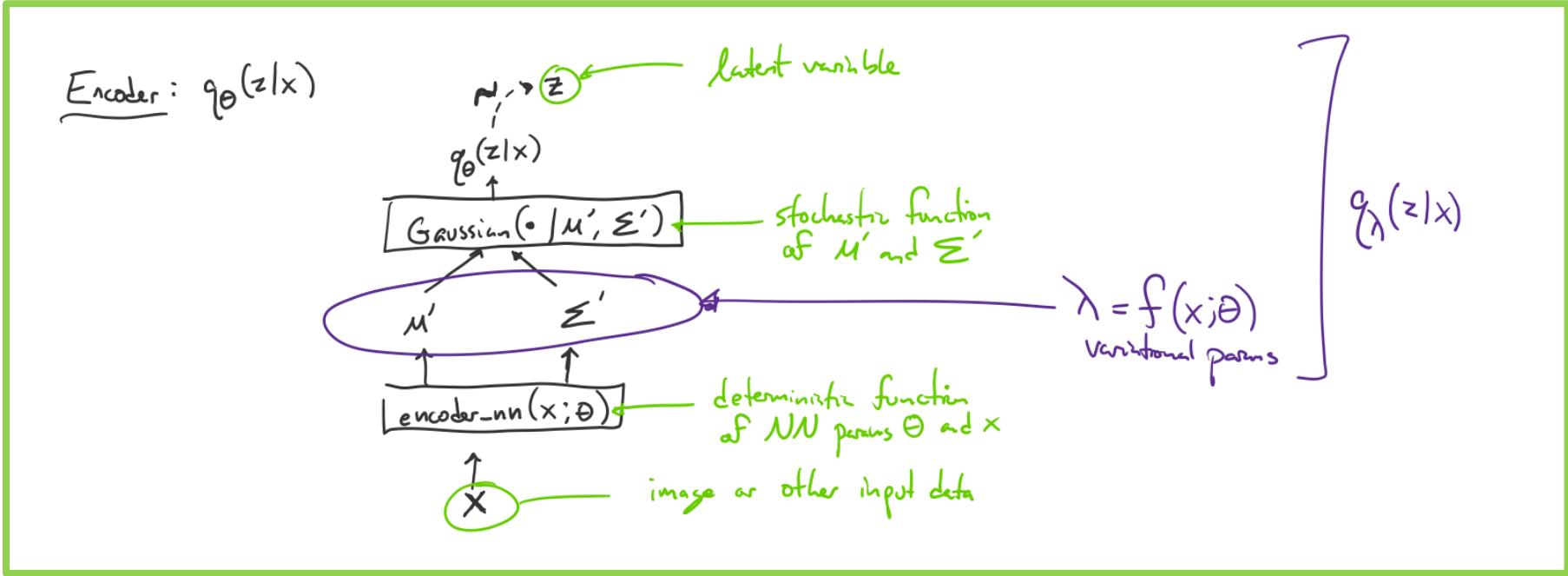
$$p_{\phi}(\mathbf{x}, \mathbf{z})$$
$$\mathbf{z} \sim \text{Gaussian}(\mathbf{0}, \mathbf{I})$$



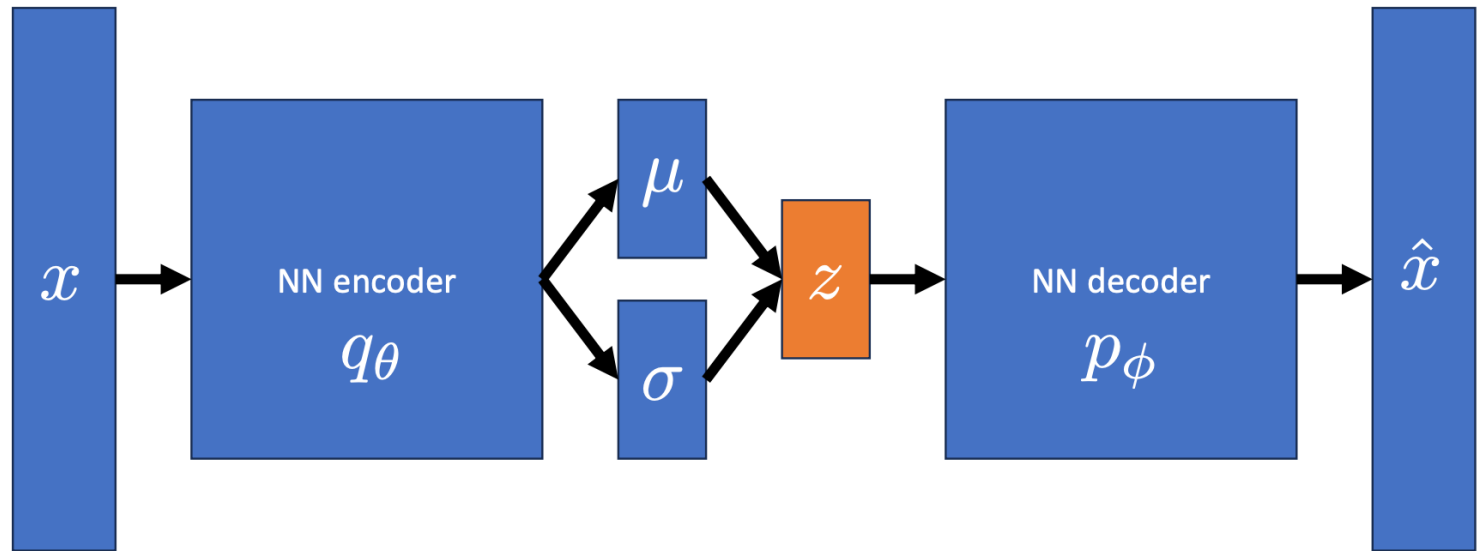
$$q_{\lambda}(\mathbf{z} | \mathbf{x})$$
$$\lambda = f(\mathbf{x}; \theta)$$



VAEs: Neural Network View



Variational Autoencoder: Network Perspective

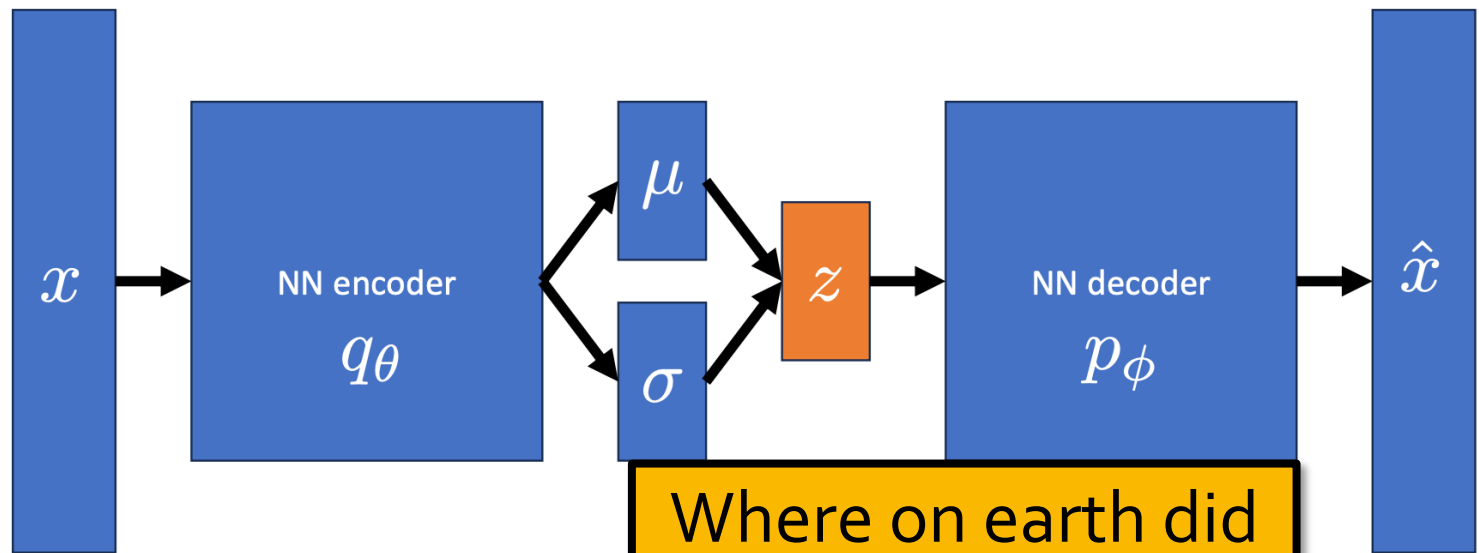


- Objective: minimize the negative log-likelihood of the dataset plus a *regularization term* that encourages a dense latent space

$$J(\theta, \phi) = \sum_{i=1}^N \ell_i(\theta, \phi)$$

$$\ell_i(\theta, \phi) = -\mathbb{E}_{q_{\theta}(z|x^{(i)})}[\log p_{\phi}(x^{(i)}|z)] + KL(q_{\theta}(z|x^{(i)}) \parallel p(z))$$

Variational Autoencoder: Network Perspective



- Objective: minimize the negative log-likelihood of the dataset plus a *regularization term* that constrains the latent space

Where on earth did this objective function come from?

$$J(\theta, \phi) = \sum_{i=1}^N \ell_i(\theta, \phi)$$

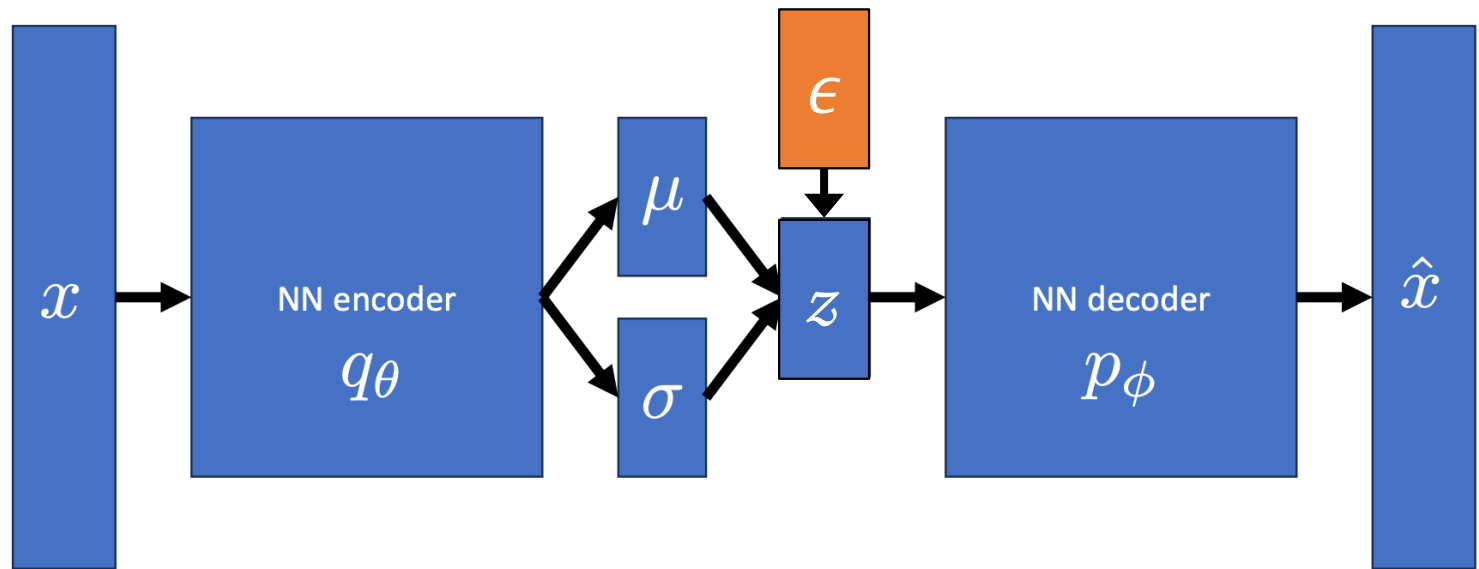
$$\ell_i(\theta, \phi) = -\mathbb{E}_{q_\theta(z|x^{(i)})} [\log p_\phi(x^{(i)}|z)] + KL(q_\theta(z|x^{(i)}) \parallel p(z))$$

VAE Objective Function

$$\begin{aligned} - \text{ELBO}(q_\theta) &= \left(E_{q_\theta(\mathbf{z}|\mathbf{x})} [\log p_\alpha(\mathbf{x}, \mathbf{z})] - E_{q_\theta(\mathbf{z}|\mathbf{x})} [\log q_\theta(\mathbf{z} | \mathbf{x})] \right) \\ &= \left(E_q [\log p(\mathbf{x}|\mathbf{z})] + \underbrace{E_q [\log p(\mathbf{z})] - E_q [\log q(\mathbf{z}|\mathbf{x})]} \right) \\ &= \left(E_q [\log p(\mathbf{x}|\mathbf{z})] + -E_q \left[\frac{\log q(\mathbf{z}|\mathbf{x})}{\log p(\mathbf{z})} \right] \right) \\ &= \end{aligned}$$

$$\ell_i(\boldsymbol{\theta}, \boldsymbol{\phi}) = -\mathbb{E}_{q_\theta(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_\phi(\mathbf{x}^{(i)}|\mathbf{z})] + KL(q_\theta(\mathbf{z}|\mathbf{x}^{(i)}) \parallel p(\mathbf{z}))$$

Reparameterization Trick



- Objective: minimize the negative log-likelihood of the dataset plus a *regularization term* that encourages a dense latent space

$$J(\theta, \phi) = \sum_{i=1}^N \ell_i(\theta, \phi)$$

$$\ell_i(\theta, \phi) \approx - \left(\sum_{s=1}^S \log p_\phi(\mathbf{x}^{(i)} | \mathbf{z}_s(\theta)) \right) + KL(q_\theta(\mathbf{z} | \mathbf{x}^{(i)}) \parallel p(\mathbf{z}))$$

for $\mathbf{z}_s(\theta) = \mu_\theta(\mathbf{x}^{(i)}) + \sigma_\theta(\mathbf{x}^{(i)}) \odot \epsilon_s$ where $\epsilon_s \sim N(\mathbf{0}, I)$

VAE RESULTS

VAEs for Image Generation

Kingma & Welling (2014)

- introduced VAEs
- applied to image generation

Model

- $p_{\phi}(\mathbf{z}) \sim N(\mathbf{z}; \mathbf{0}, \mathbf{I})$
- $p_{\phi}(\mathbf{x} | \mathbf{z})$ is a multivariate Gaussian with mean and variance computed by an MLP, fully connected neural network with a single hidden layer with parameters ϕ
- $q_{\theta}(\mathbf{z} | \mathbf{x})$ is a multivariate Gaussian with diagonal covariance structure and with mean and variance computed by an MLP with parameters θ

Auto-Encoding Variational Bayes

Diederik P. Kingma
Machine Learning Group
Universiteit van Amsterdam
dpkingma@gmail.com

Max Welling
Machine Learning Group
Universiteit van Amsterdam
welling.max@gmail.com

Abstract

How can we perform efficient inference and learning in directed probabilistic models, in the presence of continuous latent variables with intractable posterior distributions, and large datasets? We introduce a stochastic variational inference and learning algorithm that scales to large datasets and, under some mild differentiability conditions, even works in the intractable case. Our contributions is two-fold. First, we show that a reparameterization of the variational lower bound yields a lower bound estimator that can be straightforwardly optimized using standard stochastic gradient methods. Second, we show that for i.i.d. datasets with continuous latent variables per datapoint, posterior inference can be made especially efficient by fitting an approximate inference model (also called a recognition model) to the intractable posterior using the proposed lower bound estimator. Theoretical advantages are reflected in experimental results.

VAEs for Image Generation

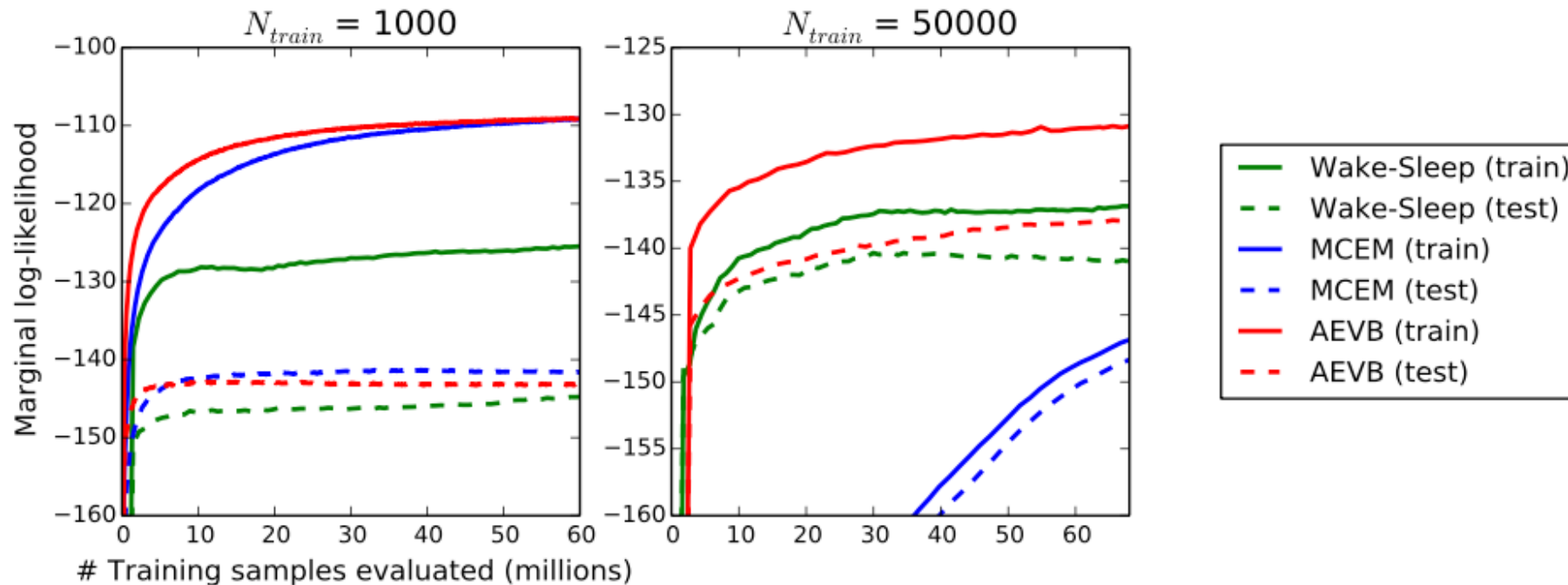
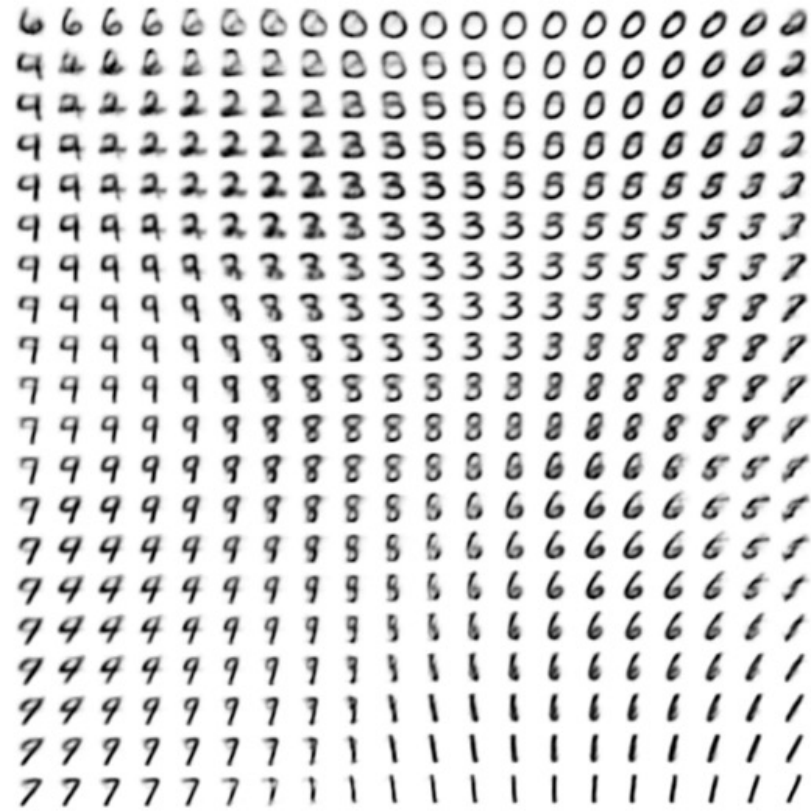


Figure 3: Comparison of AEVB to the wake-sleep algorithm and Monte Carlo EM, in terms of the estimated marginal likelihood, for a different number of training points. Monte Carlo EM is not an on-line algorithm, and (unlike AEVB and the wake-sleep method) can't be applied efficiently for the full MNIST dataset.

VAEs for Image Generation



(a) Learned Frey Face manifold



(b) Learned MNIST manifold

Figure 4: Visualisations of learned data manifold for generative models with two-dimensional latent space, learned with AEVB. Since the prior of the latent space is Gaussian, linearly spaced coordinates on the unit square were transformed through the inverse CDF of the Gaussian to produce values of the latent variables \mathbf{z} . For each of these values \mathbf{z} , we plotted the corresponding generative $p_{\theta}(\mathbf{x}|\mathbf{z})$ with the learned parameters θ .

VAEs for Image Generation

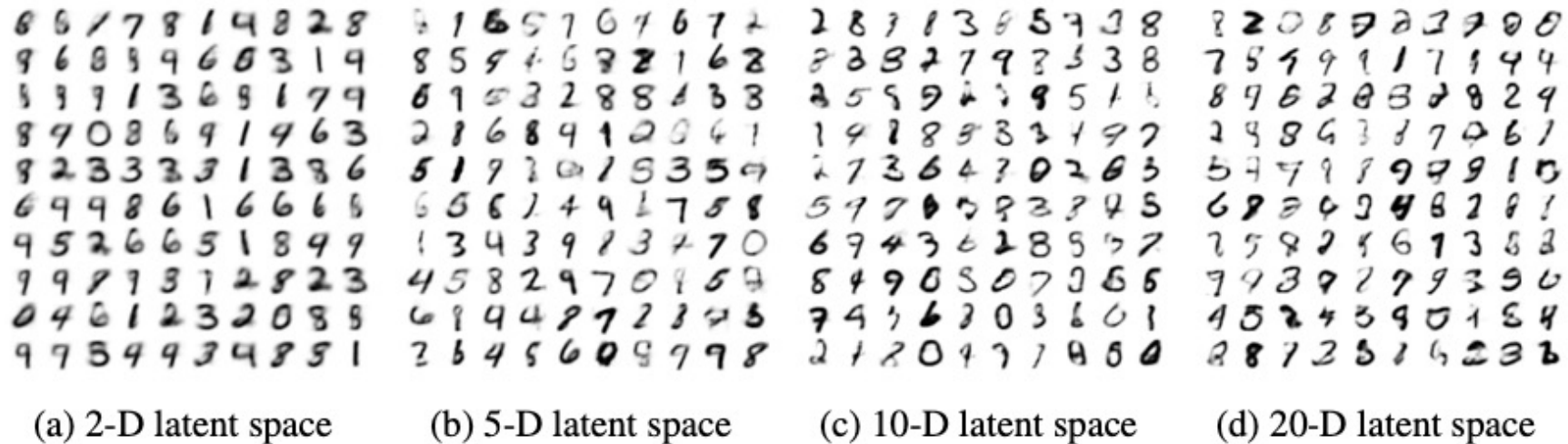


Figure 5: Random samples from learned generative models of MNIST for different dimensionalities of latent space.

VAEs for Text Generation

Bowman et al. (2015)

- example of an application of VAEs to discrete data
- built on the sequence-to-sequence framework:
 - input is read in by an LSTM
 - output is generated by an LSTM-LM

Model

- $p_{\phi}(\mathbf{z}) \sim N(\mathbf{z}; \mathbf{0}, \mathbf{I})$
- $p_{\phi}(\mathbf{x} | \mathbf{z})$ is an LSTM Language Model with parameters ϕ
- $q_{\theta}(\mathbf{z} | \mathbf{x})$ is a multivariate Gaussian with mean and variance computed by an LSTM with parameters θ

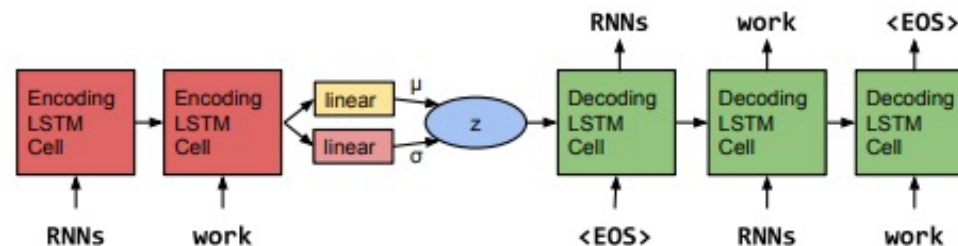


Figure 1: The core structure of our variational autoencoder language model. Words are represented using a learned dictionary of embedding vectors.

VAEs for Text Generation

INPUT	we looked out at the setting sun .	i went to the kitchen .	how are you doing ?
MEAN	<i>they were laughing at the same time .</i>	<i>i went to the kitchen .</i>	<i>what are you doing ?</i>
SAMP. 1	<i>ill see you in the early morning .</i>	<i>i went to my apartment .</i>	<i>“ are you sure ?</i>
SAMP. 2	<i>i looked up at the blue sky .</i>	<i>i looked around the room .</i>	<i>what are you doing ?</i>
SAMP. 3	<i>it was down on the dance floor .</i>	<i>i turned back to the table .</i>	<i>what are you doing ?</i>

Table 7: Three sentences which were used as inputs to the VAE, presented with greedy decodes from the mean of the posterior distribution, and from three samples from that distribution.

“ i want to talk to you . ”
“i want to be with you . ”
“i do n’t want to be with you . ”
i do n’t want to be with you .
she did n’t want to be with him .

he was silent for a long moment .
he was silent for a moment .
it was quiet for a moment .
it was dark and cold .
there was a pause .
it was my turn .

Table 8: Paths between pairs of random points in VAE space: Note that intermediate sentences are grammatical, and that topic and syntactic structure are usually locally consistent.

VQ-VAE

- Vector Quantized VAE (VQ-VAE) learns a continuous codebook, but the encoder outputs discrete codes
- Decoder takes a code and generates a sample conditioned on it

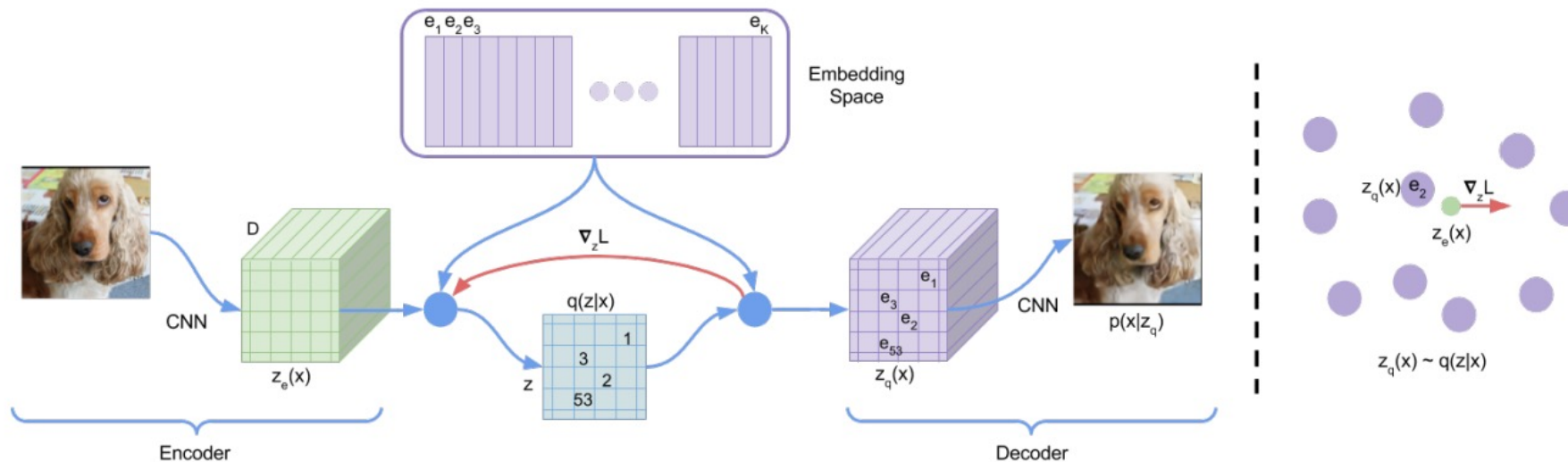
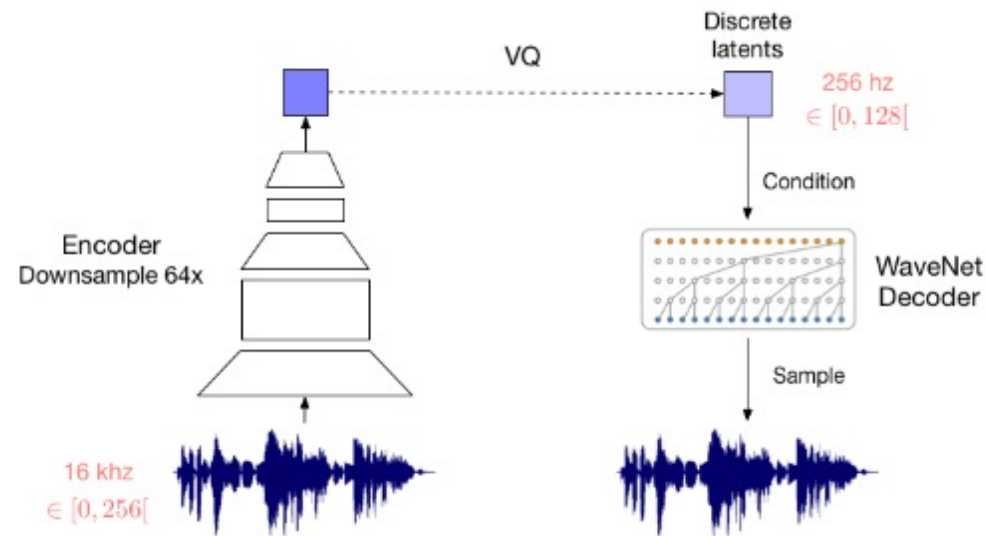


Figure 1: Left: A figure describing the VQ-VAE. Right: Visualisation of the embedding space. The output of the encoder $z(x)$ is mapped to the nearest point e_2 . The gradient $\nabla_z L$ (in red) will push the encoder to change its output, which could alter the configuration in the next forward pass.

VQ-VAE

- Vector Quantized VAE (VQ-VAE) learns a continuous codebook, but the encoder outputs discrete codes
- Decoder takes a code and generates a sample conditioned on it

Example: Generating Audio



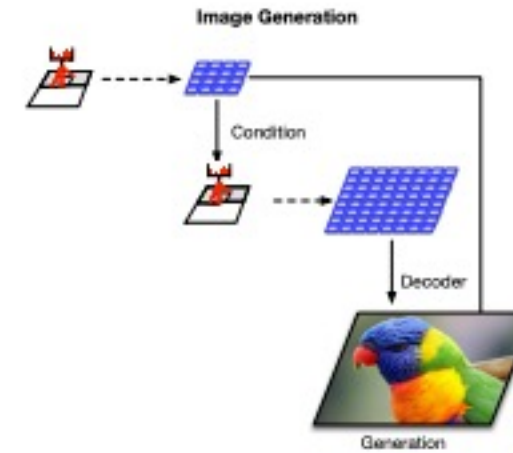
<https://avdnoord.github.io/homepage/vqvae>

VQ-VAE

- VQ-VAE-2 extended the original idea by learning two levels (bottom and top) and a strong prior over the latent space
- Samples from this new model can be convincing even at high-fidelity



(a) Overview of the architecture of our hierarchical VQ-VAE. The encoders and decoders consist of deep neural networks. The input to the model is a 256×256 image that is compressed to quantized latent maps of size 64×64 and 32×32 for the *bottom* and *top* levels, respectively. The decoder reconstructs the image from the two latent maps.



(b) Multi-stage image generation. The top-level PixelCNN prior is conditioned on the class label, the bottom level PixelCNN is conditioned on the class label as well as the first level code. Thanks to the feed-forward decoder, the mapping between latents to pixels is fast. (The example image with a parrot is generated with this model).

- VQ-VAE-2 extended the original idea by learning two levels (bottom and top) and a strong prior over the latent space
- Samples from this new model can be convincing even at high-fidelity



Figure 4: Class conditional random samples. Classes from the top row are: 108 sea anemone, 109 brain coral, 114 slug, 11 goldfinch, 130 flamingo, 141 redshank, 154 Pekinese, 157 papillon, 97 drake, and 28 spotted salamander.

- VQ-VAE-2 extended the original idea by learning two levels (bottom and top) and a strong prior over the latent space
- Samples from this new model can be convincing even at high-fidelity

Figure from Razavi et al. (2019)



- VQ-VAE-2 extended the original idea by learning two levels (bottom and top) and a strong prior over the latent space
- Samples from this new model can be convincing even at high-fidelity

Figure from Razavi et al. (2019)



- VQ-VAE-2 extended the original idea by learning two levels (bottom and top) and a strong prior over the latent space
- Samples from this new model can be convincing even at high-fidelity

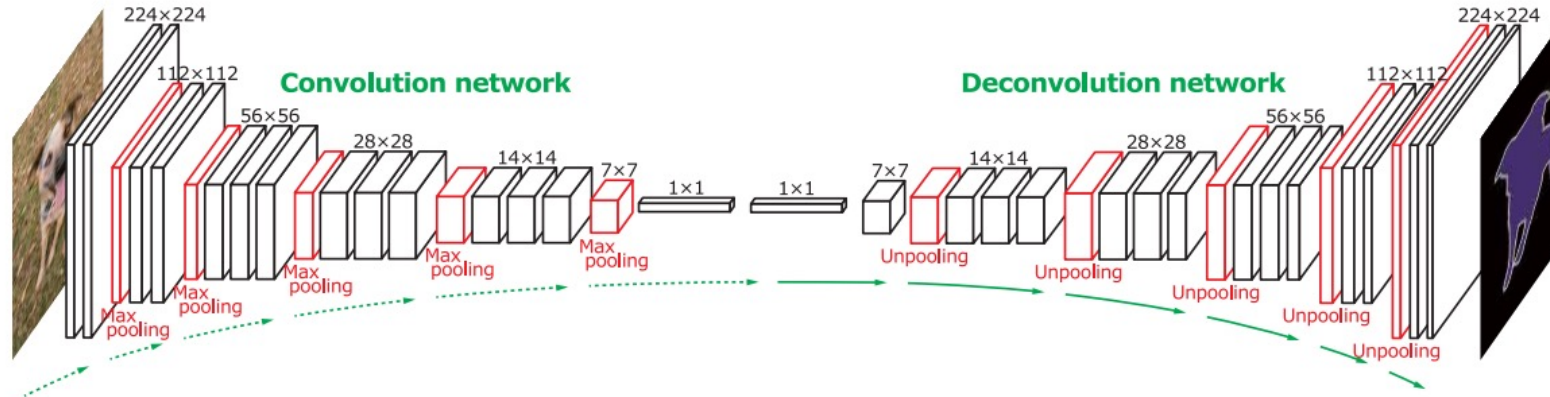
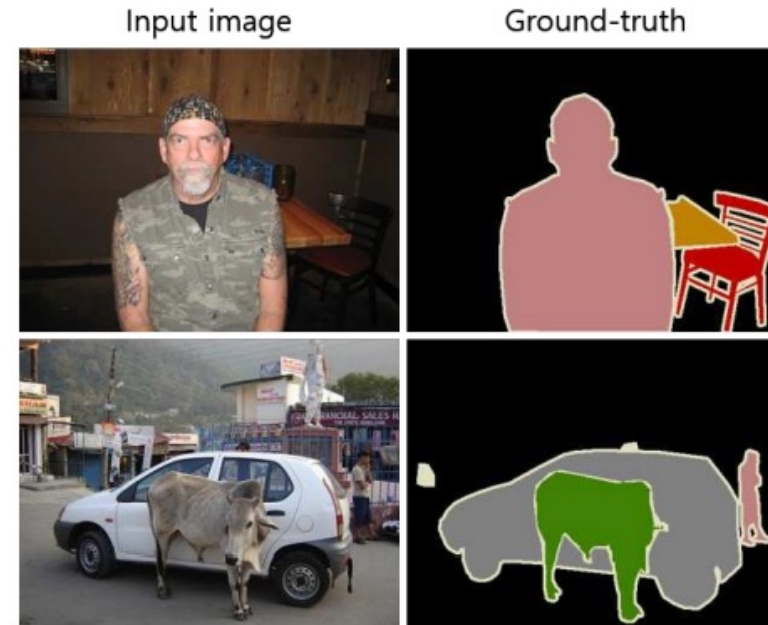
Figure from Razavi et al. (2019)



U-NET

Semantic Segmentation

- Given an image, predict a label for every pixel in the image
- Not merely a classification problem, because there are strong correlations between pixel-specific labels



Instance Segmentation

- Predict per-pixel labels as in semantic segmentation, but differentiate between different instances of the same label
- *Example:* if there are two people in the image, one person should be labeled **person-1** and one should be labeled **person-2**

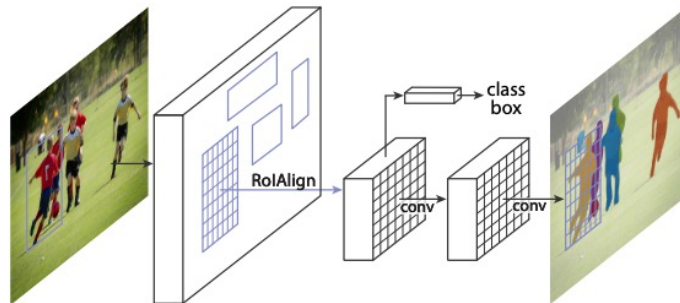
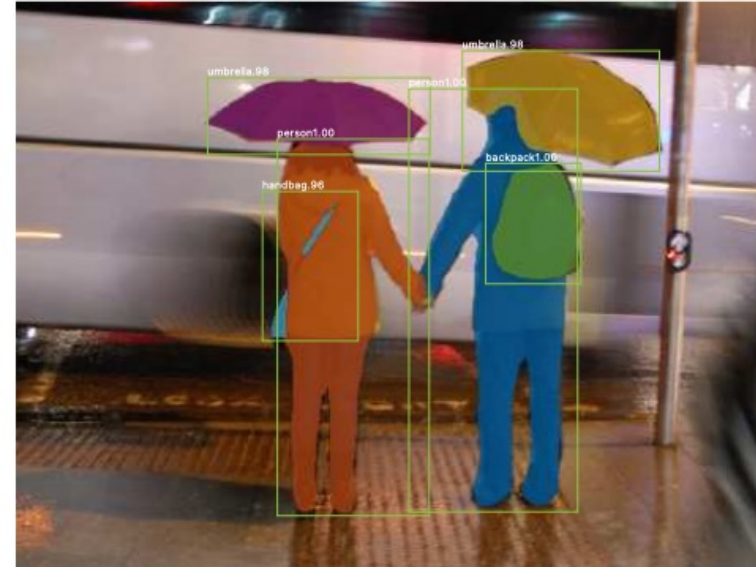


Figure 1. The **Mask R-CNN** framework for instance segmentation.

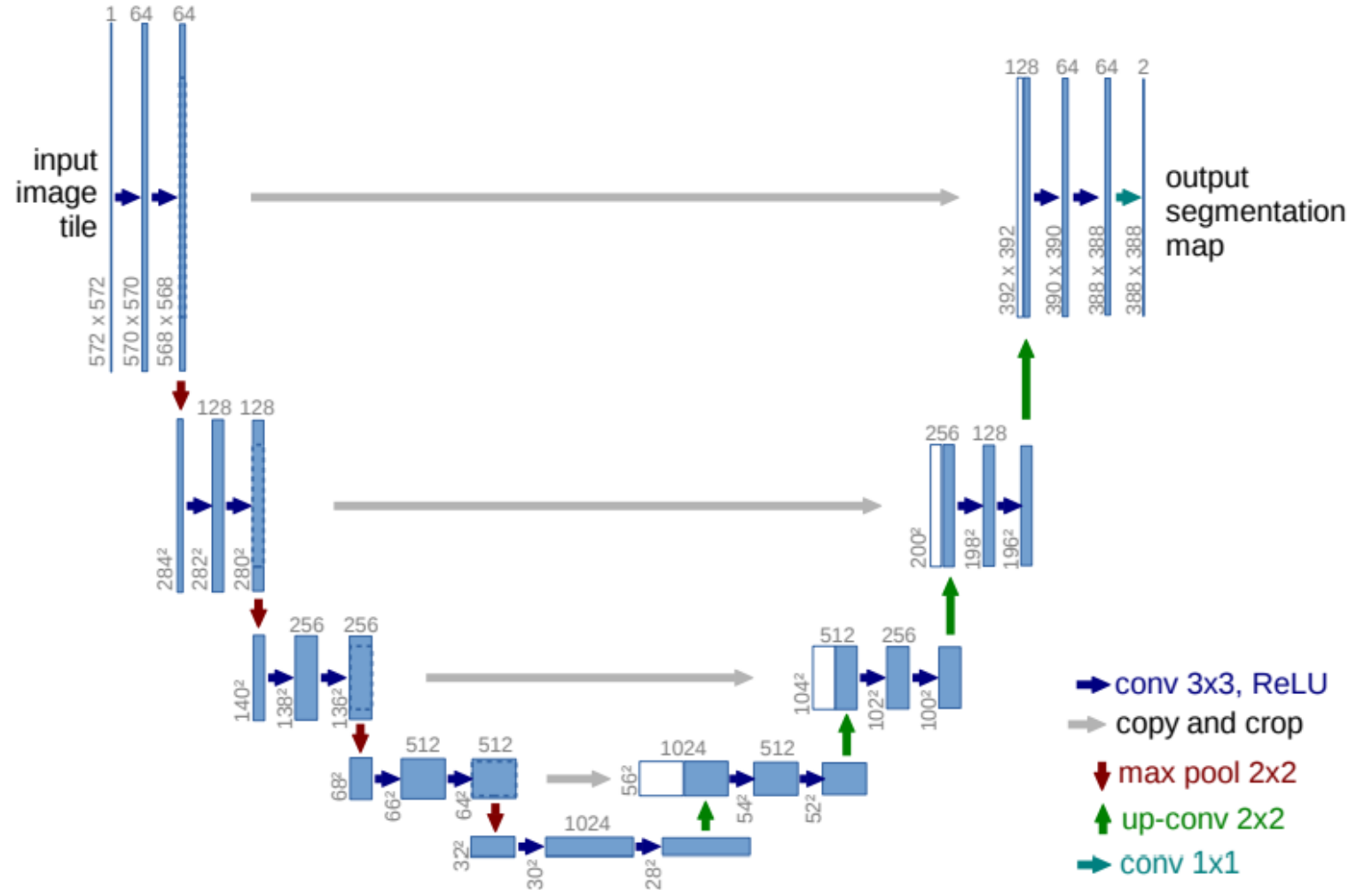
U-Net

Contracting path

- block consists of:
 - 3x3 convolution
 - 3x3 convolution
 - ReLU
 - max-pooling with stride of 2 (downsample)
- repeat the block N times, doubling number of channels

Expanding path

- block consists of:
 - 2x2 convolution (upsampling)
 - concatenation with contracting path features
 - 3x3 convolution
 - 3x3 convolution
 - ReLU
- repeat the block N times, halving the number of channels



U-Net

- Originally designed for applications to biomedical segmentation
- Key observation is that the output layer has the **same** dimensions as the input image (possibly with different number of channels)

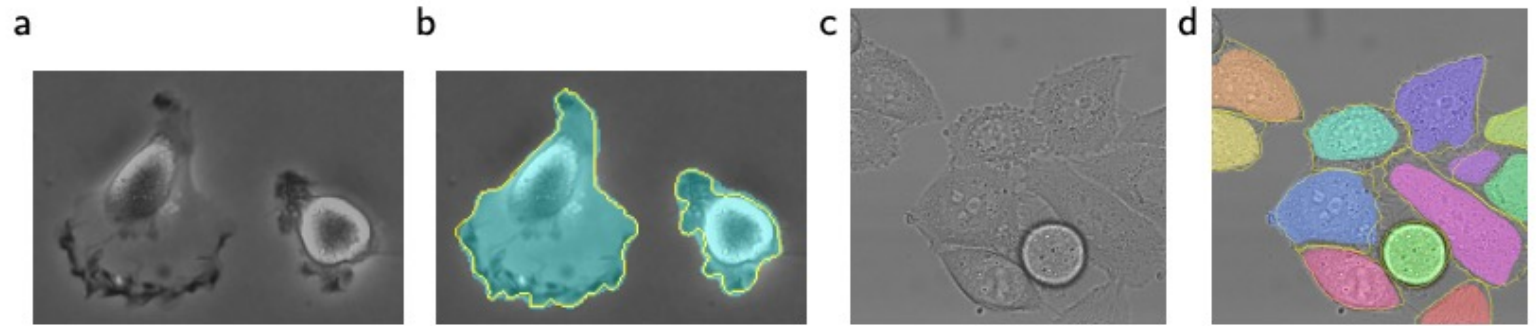


Fig. 4. Result on the ISBI cell tracking challenge. (a) part of an input image of the “PhC-U373” data set. (b) Segmentation result (cyan mask) with manual ground truth (yellow border) (c) input image of the “DIC-HeLa” data set. (d) Segmentation result (random colored masks) with manual ground truth (yellow border).