# Variational Autoencoders & Lossy Neural Compression
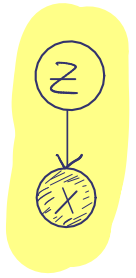
## Recap from last lecture: Variational Inference (VI)

- latent variable model: $P(Z, X) = P(Z) P(X | Z)$

- goal: approximate the posterior: $P(Z | X=x) = \dfrac{P(Z) P(X=x | Z)}{\int P(Z=z) P(X=x | Z=z) \, dz}$

→ VI turns the inference problem into an optimization problem

→ variational distribution: $Q_\phi (Z | X=x)$
  ↖ variational parameters

### Evidence lower bound (ELBO):

- $ELBO(\phi) = -\mathbb{E}_{Q_\phi(Z|X=x)} \left[ \tilde{R}_{net}^{(z)} (x) \right]$  ← how we motivated it

Problem Set 8 $\begin{cases}
= \mathbb{E}_{Q_\phi(Z|X=x)} \left[ \log P(Z, X=x) - \log Q_\phi (Z|X=x) \right] \quad \leftarrow \text{most explicit formulation ("regularized MAP")} \\
\qquad\qquad \underbrace{\phantom{-\log Q_\phi(Z|X=x)}}_{\to H[Q_\phi(Z|X=x)]} \\
= \mathbb{E}_{Q_\phi(Z|X=x)} \left[ \log P(X=x|Z) \right] - D_{KL}\left( Q_\phi(Z|X=x) \| P(Z) \right) \quad \leftarrow \text{"regularized maximum likelihood"} \\
= \log P(X=x) - D_{KL}\left( Q_\phi(Z|X=x) \| P(Z|X=x) \right) \quad \leftarrow \text{connects VI to Bayesian inference}
\end{cases}$

### How to maximize the ELBO with stochastic gradient optimization:

- reparameterization gradients: $\nabla_\phi \mathbb{E}_{Q_\phi(Z|X=x)} \left[ \ell(Z, \phi) \right] = \nabla_\phi \mathbb{E}_{\varepsilon \sim Q_0} \left[ \ell(g(\varepsilon, \phi), \phi) \right]$
  $Z = g(\varepsilon, \phi)$ where $\varepsilon \sim Q_0$ ← fixed distribution

- score function gradients: $\nabla_\phi \mathbb{E}_{Q_\phi(Z|X=x)} \left[ \ell(Z, \phi) \right] = \mathbb{E}_{Q_\phi(Z|X=x)} \left[ \left( \nabla_\phi \log Q_\phi(Z|X=x) \right) \times \ell(Z, \phi) \right]$
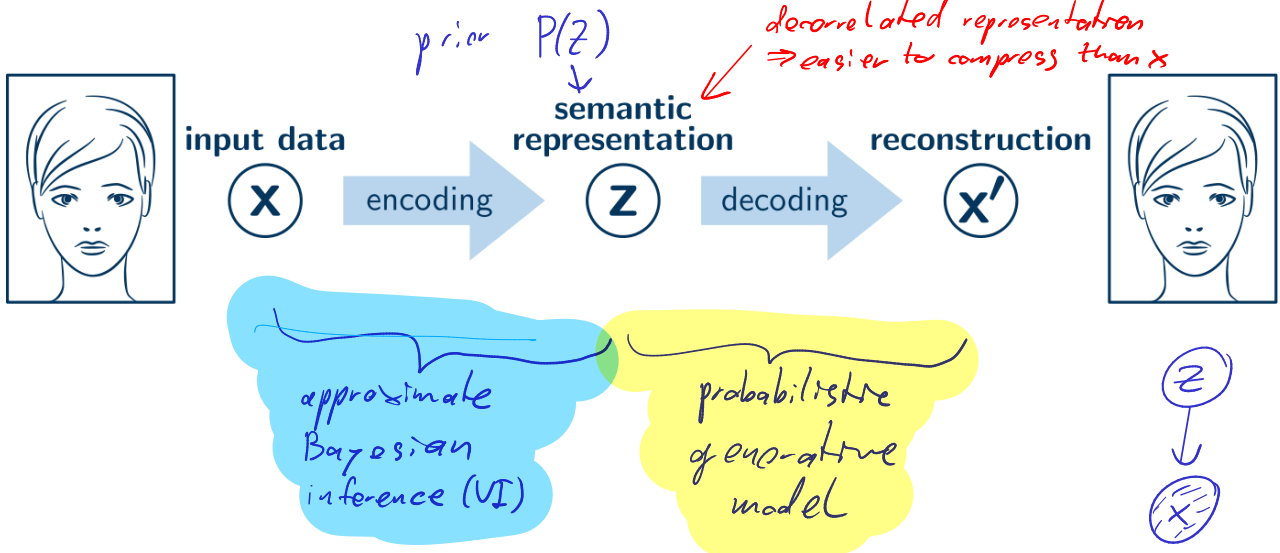  (= REINFORCE method)

### Limitations so far:

(i) the generative model P(Z, X) is fixed -- and therefore limited to simple models that we can come up with manually; and

(ii) for every concrete message x that we want to compress, we have to run an expensive optimization procedure to find the optimal variational parameters φ*.

TODAY: overcoming these limitations

$\Big\}$ → Variational Expectation Maximization (Variational EM) → "learn the prob. gen. model P from training data"

$\Big\}$ → Amortized variational inference "learn how to do inference"

**Spoiler: Variational Autoencoders (VAEs)** = amortized variational expectation maximization

*prior* $P(Z)$

*decorrelated representation*
$\Rightarrow$ *easier to compress than $x$*

**input data** $X$ — encoding → **semantic representation** $Z$ — decoding → **reconstruction** $X'$

*approximate Bayesian inference (VI)*

*probabilistic generative model*

$Z \rightarrow X$

# Variational Expectation Maximization: learning a latent variable model

[Beal & Ghahramani, Bayesian statistics, 2003]

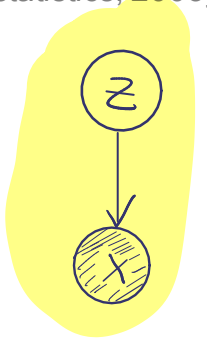Introduce free parameters into the probabilistic model P(Z, X):

$$P_\theta(Z, X) = P_\theta(Z)\, P_\theta(X \mid Z)$$

*model parameters $\theta$*

$$\mathbb{E}_{x \sim \text{train set}}\big[-\log P_\theta(X = x)\big]$$
$$= \mathbb{E}_{x \sim \text{train set}}\Big[-\log\Big(\sum_z P_\theta(Z = z, X = x)\Big)\Big]$$

↑ *prohibitively expensive*

$Z \rightarrow X$

For example, the likelihood could be parameterized by a neural network with weights θ:

$$P_\theta(X \mid Z = z) = \mathcal{N}\big(g_\theta(z),\ \sigma_0^2 I\big)$$

*"a normal distribution with mean $g_\theta(z)$ and variance $\sigma_0^2$ in each coordinate direction"*

$\mathcal{N}(x_i; \mu, \sigma_0^2)$

*$2\sigma$*

*where $g_\theta$: latent space $\rightarrow$ data space is a neural network with trainable weights $\theta$*

Thus, the ELBO now depends both on the variational parameters ϕ and on the model parameters θ:

$$\mathrm{ELBO}(\theta, \phi_x) = -\mathbb{E}\quad R_{net}$$

• *maximize over both of these jointly (at training time)*

$$= \mathbb{E}\quad \big[\log P - \log Q\big]$$

• *at compression time: keep $\theta$ fixed & maximize only over $\phi$*

$$= \log P_\theta(X = x) - D_{KL}$$

*— optimal bitrate with model $P_\theta$*

*overhead due to VI*

minimize the expected $\overset{not}{\text{bitrate}}$ of modified bits-back

$\Rightarrow$ maximize: $\text{ELBO}(\phi, \vartheta) = - \mathbb{E}_{Q_\phi(z|x \to)}\left[\tilde{R}_{not}^{(z, \vartheta)}(x)\right]$

$\underset{\substack{\text{variational} \\ \text{params}}}{\uparrow} \quad \underset{\substack{\text{model} \\ \text{params}}}{\uparrow}$

Alternative:

- store $\phi_x \quad \forall x \in$ train set on disk

- training loop:

    for training_step in $\{1, 2, 3, \ldots, n\}$:

        sample $\boxed{x} \sim$ train set

        look up $\phi_x$ on disk

        calculate $g_\vartheta = \nabla_\vartheta \text{ELBO}(\vartheta, \phi_x, x)$

        $\qquad \qquad g_\phi = \nabla_{\phi_x} \text{ELBO}(\vartheta, \phi_x, x) \Leftarrow$

        update $\vartheta \leftarrow \vartheta + g_\vartheta$

        $\qquad \qquad \phi \leftarrow \phi + g_\phi \Leftarrow$

        store $\phi_x$ back to disk

**Data compression with learnt latent variable models (try 1: without amortization):**

1) When designing the compression method:
   - collect large (unlabeled) data set of training samples (e.g., a large collection of images)
   - come up with a model architecture for the generative model $P_\theta$ (that still has free parameters)
   - train the model by maximizing the ELBO jointly over both θ and φ.

$$\text{in detail:} \quad \theta^*, \{\phi^*_x\} := \arg\max_{\theta, \{\phi_x\}} \mathbb{E}_{x \sim \text{trainset}}\left[ELBO(\theta, \phi_x, x)\right]$$

   $\phi_x$ are the variational params for training point $x$

   - throw away $\{\phi^*_x\}$ and share θ* between sender and receiver

   for training step $\in \{1, 2, ..., N\}$
   find $\phi^*_x := \arg\max ELBO(\theta, \phi_x)$
   set $\hat{g} := \nabla_\theta ELBO(...)$
   update $\theta \leftarrow \theta + \varsigma \hat{g}$

2) When compressing some given data x (i.e., on the sender side):
   - perform variational inference, i.e., maximize ELBO(θ*, φ, x) over φ but keep θ* fixed at the agreed-upon values.
   - use probabilistic generative model $P_{\theta^*}(Z, X)$ and the resulting variational distribution $Q_{\phi^*}(Z \mid X{=}x)$ to compress x.

3) When decompressing data (i.e., on the receiver side):
   - needs the exact same probabilistic generative model $P_{\theta^*}(Z, X)$ that the sender used for compresion.
   - if the data was compressed with bits-back coding, then the receiver also needs to perform variational inference once it has reconstructed x (i.e., maximize ELBO(θ, φ) over φ but keep θ*).

compress $z$
using prior $P_{\theta^*}(z)$



input data **X** — encoding → semantic representation **Z** — decoding → reconstruction **X'**

$Q_{\phi^*}(Z \mid X{=}x)$
where $\phi^* := \arg\max_\phi ELBO(\phi, \theta^*, x)$
expensive operation

$P_{\theta^*}(X \mid Z)$
where $\theta^*$ is known at (de-)compression time

**Amortized Variational Inference: learn how to do inference**

Variational inference maps data $x$ to a variational distribution $Q_{\phi^*}(Z \mid X{=}x)$:

Idea: learn this mapping from x to $Q_{\phi^*}(Z \mid X{=}x)$:

   - rename the parameteres of $Q_\lambda(Z)$ from φ to λ

   for example: $Q_\lambda(Z) = \mathcal{N}(\mu, \text{diag}(\sigma^2))$
   $\lambda = (\mu, \sigma^2)$

   - rather than optimizing over λ, learn a function f that maps x to λ (and that is parameterized by some neural network weights φ:

$\lambda = f_\phi(x)$   where $f_\phi$ is a neural network with weights φ

$Q_\phi(z \mid x{=}x) := Q_\lambda(Z)$
where $\lambda = f_\phi(x)$

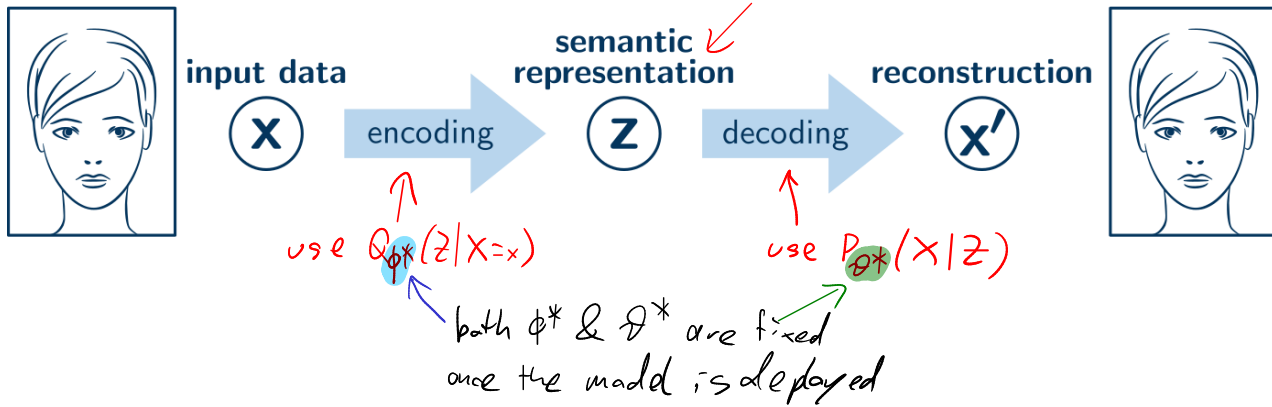⇒ variational parameters φ are now <u>shared</u> between all data points x

$$ELBO(\phi, \theta, x) = \mathbb{E}_{Q_\phi(Z \mid X{=}x)}\left[\log P(Z, X{=}x) - \log Q_\phi(Z \mid X{=}x)\right]$$

Combining amortized inference with variational expectation maximization results in:

**Variational Autoencoers (VAEs):**   prior $P_{\theta^*}(z)$   [Kingma and Welling, 2015]



input data **X** → encoding → semantic representation **Z** → decoding → reconstruction **X'**

use $Q_{\phi^*}(z|X=x)$      use $P_{\theta^*}(X|Z)$

both $\phi^*$ & $\theta^*$ are fixed
once the model is deployed

training objective: $ELBO(\theta, \phi) = \mathbb{E}_{x \sim P_{train}(x)} \left[ \mathbb{E}_{Q_\phi(z|x=x)} \left[ \log P_\theta(z, x=x) - \log Q_\phi(z|x=x) \right] \right]$
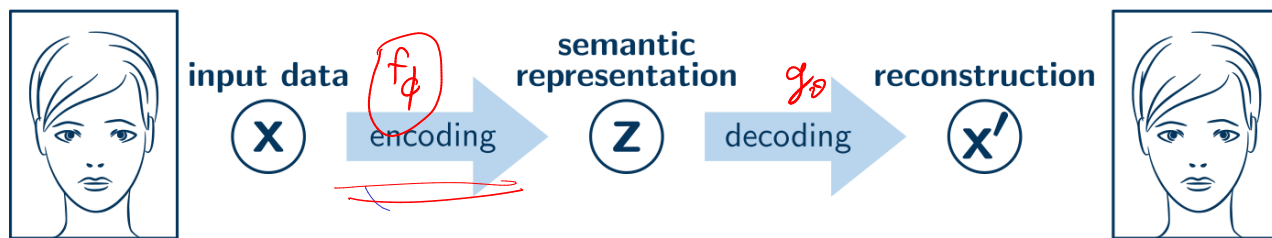
Problem Set: implement a variational autoencoder for simple images (MNIST) $= \mathbb{E}_{Q_\phi(z|x=x)} \left[ \log P_\theta(x=x|z) \right]$
$- D_{KL}(Q_\phi(z|x=x) \| P_\theta(z))$

**Note:** Variational expectation maximization (EM) is not limited to VAEs. Even without amortized inference, variational EM is a very useful algorithm that is very simple and allows you to treat some model parameters (Z) probabilistically while using point estimates for others (θ).


# Lossy Compression with VAEs -- a Brief Pragmatic Introduction

**(more details: Problem 10.1 on Problem Set 10 and lectures 11 and 12)**

- VAEs are popular models for lossy data compression
- here's a brief overview; we'll dive deeper into lossy compression starting next week



input data **X** $f_\phi$ encoding → semantic representation **Z** $g_\theta$ decoding → reconstruction **X'**

Simplest (yet surprisingly powerful) method [Ballé et al., 2016]:     part of model params $\theta$

- generative model:  Prior must have a learnable variance: e.g., $P_\theta(z) = \mathcal{N}(0, \sigma^2 I)$   (for each coordinate)

- variational distribution:  $Q_\phi(z|x=x) = Uniform\left(\left[ f_\phi(x) - \frac{1}{2}, f_\phi(x) + \frac{1}{2} \right]\right)$

- encoding:  • given input $x$, calculate $f_\phi(x)$, then round each component to nearest integer $\to z$
  • encode $z$ using $P(z)$

- decoding:  • decode losslessly $z$ using Prior $P(z)$   ← $z$ needs to be discrete rand. var.
  • reconstruction $\hat{x} := \arg\max_x P_\theta(X|Z=z)$

Problem: we can only compress z losslessly if it comes from a discrete distribution.

Approximation:
• assume at compression time that $z \in \mathbb{Z}^k \Rightarrow$ get $z$ by rounding each coordinate of $f_\phi(x)$
• during training: approximate rounding of $f_\phi(x)$ by adding uniform random noise $[-\frac{1}{2}, \frac{1}{2}]^k$