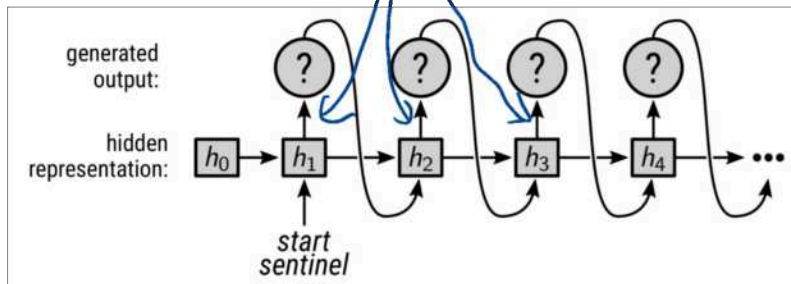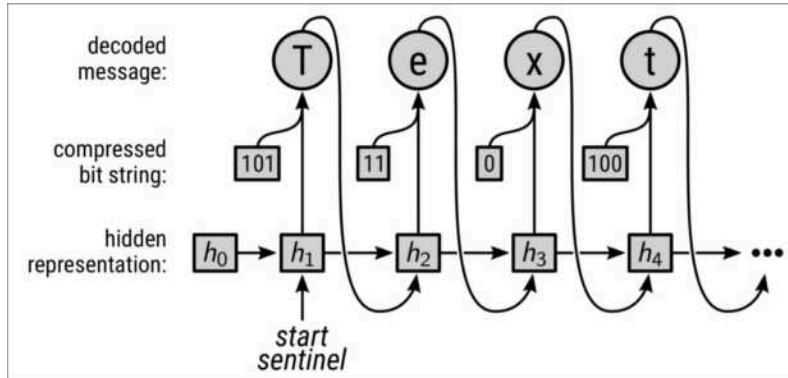# Data Compression with Deep Probabilistic Models

Reminder: Problem 3.2: compression with a learned autoregressive model

parameters as a probability dist.
$\Rightarrow$ can be used for compression



generated output: ? ? ? ?
hidden representation: $h_0$ $h_1$ $h_2$ $h_3$ $h_4$ ...
start sentinel

$\Rightarrow$ when used for compression (here: decoder side):



decoded message: T e x t
compressed bit string: 101 11 0 100
hidden representation: $h_0$ $h_1$ $h_2$ $h_3$ $h_4$ ...
start sentinel

autoregressive models: $P_\theta(\underline{X}) = P_\theta(X_1)\, P_\theta(X_2|X_1)\, P_\theta(X_3|X_1, X_2)$

model parameters (neural network weights)
$\rightarrow$ optimize $\theta$ by minimizing an empirical estimate of cross entropy $H(P_{data}, P_\theta)$

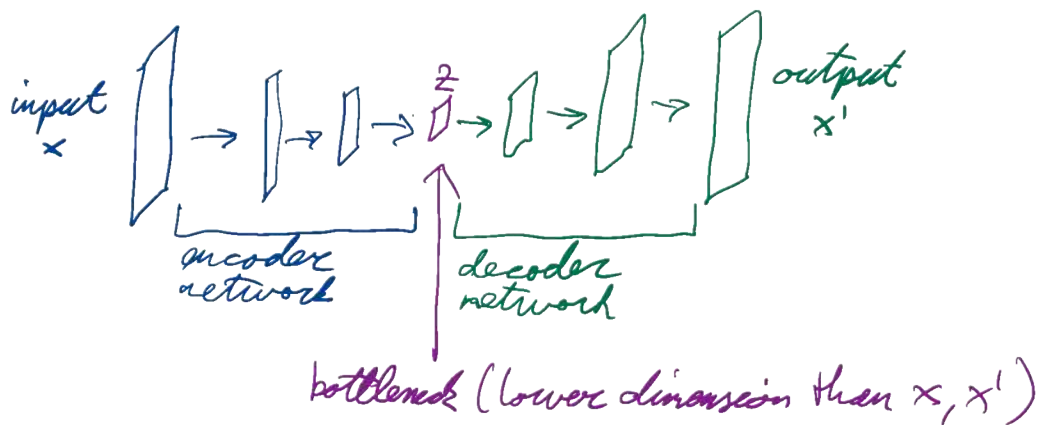$\rightarrow$ can we do the same thing with latent variable models

# Deep Latent Variable Models & Scalable Approximate Bayesian Inference

<u>Spoiler:</u> variational autoencoders (VAEs)

→ a form of <u>representation learning</u>

→ often introduced with the following explanation:

"learn to map data to itself while squeezing it through a bottleneck"



bottleneck (lower dimension than $x, x'$)

uses cases of VAEs for compression

↳ lossless compression

   1) map $x$ to $z$ & encode $z$

   2) map $z$ to $x'$ & encode residual $x - x'$

↳ lossy compression: leave out residual

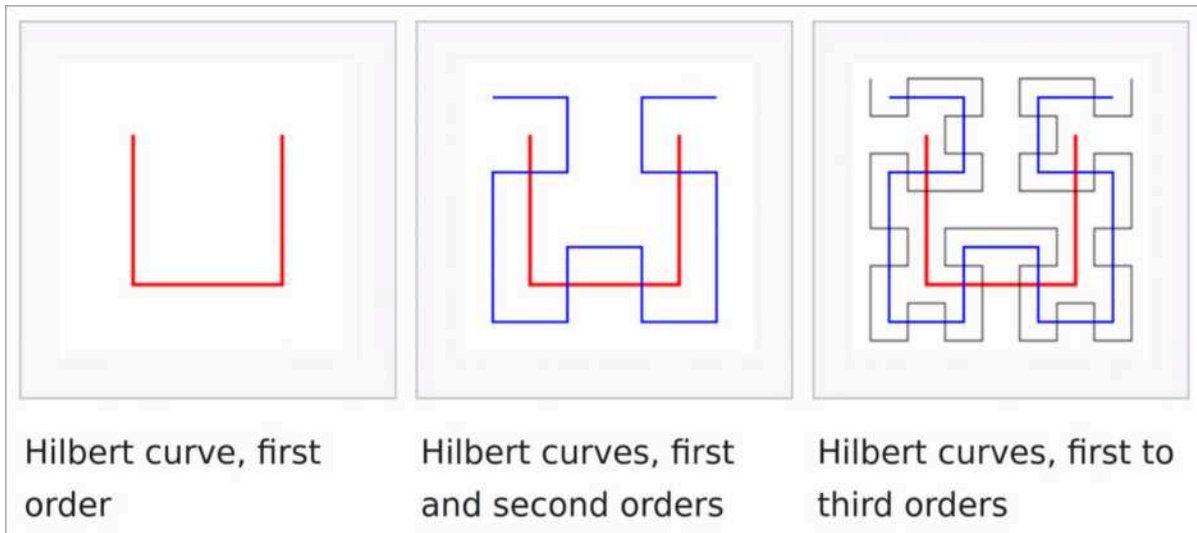⇒ 3 training objectives

   (i) decoder network should reconstruct the data well

     ($⇒$ residual $x' - x$ small / low entropy)

   (ii) encoder network decorrelates data

     → need probabilistic model   (we want $P(z) = \prod_{i=1}^{k} P(z_i)$)

Note: just squeezing data through a lower-dimensional bottleneck does not in itself imply compression
→ think about information theoretical measures rather than dim.
(iii) keep M(Z) low to enable effective compression



Hilbert curve, first order

Hilbert curves, first and second orders

Hilbert curves, first to third orders

"Hilbert curve" (drawings from Wikipedia)

**A Compiler for 3D Machine Knitting**

James McCann[1]    Lea Albaugh[1]    Vidya Narayanan[1]    April Grow[1,2]
Wojciech Matusik[3]    Jen Mankoff[1,4]    Jessica Hodgins[1]

[1]Disney Research    [2]UC Santa Cruz    [3]Massachusetts Institute of Technology    [4]Carnegie Mellon University

# Deep Latent Variable Models

· look at decoder network only

→ interpret as a latent variable model:

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



learned model parameters
(e.g. neural network weights)

common example:

↳ prior is fully factorized, i.e., $p_\theta(Z) = \prod_{i=1}^{k} p_\theta(z_i)$

↳ likelihood: $p_\theta(x|z) = \mathcal{N}(x; \, f_\theta(z), \, \sigma^2 I)$

lower case:
density fct

normal
dist. (= gaussian)

neural
network

fixed or
learned



Goal: minimize $H[P_{data}(X), \, P_\theta(X)] = \mathbb{E}_{P_{data}(X)}[-\log P_\theta(X)]$

"evidence"

Problem $P_\theta(X=x) = \int p_\theta(x, z) \, dz$

prohibitively
expensive

high dimensional

We want to maximize evidence $P_\theta(X=x)$ when evaluated on data $x$ from the training set.

<u>Recall</u>: bits-back coding

$$R_{net}(x) = -\log P_\vartheta(X=x)$$
$$= -\log P_\vartheta(Z=z) - \log P_\vartheta(X=x|Z=z) + \log \underbrace{P_\vartheta(Z=z|X=x)}_{posterior}$$

problem: $P_\vartheta(Z|X) = \dfrac{P_\vartheta(X,Z)}{P_\vartheta(X) \leftarrow intractable}$

· Idea: replace posterior with some other dist. $Q_{\lambda_x}(Z)$

(e.g.: $Q_{\lambda_x}(Z) = \prod\limits_{i=1}^{k} \mathcal{N}(z_i ; \underset{\underset{make\ up\ \lambda_x}{\uparrow\quad\uparrow}}{\mu_i, \sigma_i^2})$ )

$$\Rightarrow \tilde{R}_{net}^{(z)}(x) = -\log P_\vartheta(X=x, Z=z) + \log Q_{\lambda_x}(Z=z)$$

$$\mathbb{E}_{z \sim Q_{\lambda_x}(z)}\left[\tilde{R}_{net}^{(z)}(x)\right] \geq R_{net}(x) = \underbrace{-\log P_\vartheta(X=x)}_{\substack{we\ want\ to\ minimize\\ this}}$$

$$\underset{equality\ if\ Q_{\lambda_x}(Z)=P_\vartheta(Z|X=x)}{\Big\uparrow}$$

<u>Notation & Naming Conventions</u>

· $\log P_\vartheta(X=x)$ is called <u>evidence</u> (we want this to be high)

· $-\mathbb{E}_{z \sim Q_{\lambda_x}(z)}\left[\tilde{R}_{net}^{(z)}(x)\right] = \mathbb{E}_{z \sim Q_\lambda(z)}\left[\log P_\vartheta(X=x, Z=z) - \log Q_{\lambda_x}(Z=z)\right]$

  is called the <u>evidence lower bound</u> (ELBO)

$$\Rightarrow ELBO(\vartheta, \lambda_x) \leq \underbrace{\log P_\vartheta(X=x)}_{evidence}$$

· parameters $\lambda_x$ of the distribution $Q_{\lambda_x}(Z)$ are called

  "<u>variational parameters</u>"

· $Q_{\lambda_x}(Z)$ is called "<u>variational distribution</u>"

· Variational Inference (VI): approximate evidence

  $\log P_\vartheta(X=x)$ by $ELBO(\vartheta, \lambda_x^*)$ where

$$\lambda_x^* := \underset{\lambda_x}{arg\ max}\ ELBO(\vartheta, \lambda_x)$$

→ observation: this typically leads to a $Q_{d_x^*}(z)$
which is "close" to true posterior $P_\theta(z|X=x)$.
(Reviews: Blei et al. 2016, Zhang et al. 2018)

→ we now can approximate $\log P_\theta(X=x)$, but we still
have to maximize it over $\theta$.
→ idea: maximize our approximation $ELBO(\theta, d_x^*)$
over $\theta$.

Pseudocode:

for t in training_steps:
    sample a minibatch B of training points
    initialize $d_x$ randomly $\forall x \in B$
    for t' in inner_training_steps:
        perform gradient step for $d_x \; \forall x \in B$
    perform gradient step for $\theta$ on $ELBO(\theta, d_x^*)$

nested loop → extremely expensive

VI, finds $d_x^*$

Remember: · model params $\theta$ are <u>global</u> (i.e., the same
        for all data points $x$)
    · variational params $d_x$ parameterize an approximation
        of $P(z|X=x)$ ⇒ they are <u>local</u> (i.e., different
        for all data points $x$)
    · we want to maximize $\mathbb{E}_{x \sim P_{data}}\{\log P_\theta(X=x)\}$
      ⇒ we have to sample a new minibatch in
      each iteration of outer loop
        → invalidates $d_x^*$ from previous iteration of outer loop.

$\rightarrow$ "Variational Expectation Maximization"

(Dempster et al 1977, Beal & Ghahramani 2003)

Final additional trick: learn how to do variational inference

i.e., learn a function $g_\phi: x \mapsto \lambda_x$

set $\lambda_x = g_\phi(x)$ in the ELBO

notation: $Q_\phi(Z|x) = Q_{\lambda_x}(Z)$ with $\lambda_x = g_\phi(x)$

$$ELBO(\vartheta, \phi) = \mathbb{E}_{z \sim Q_\phi(z|x)}\left[\log P_\vartheta(X=x, Z) - \log Q_\phi(Z|x)\right]$$
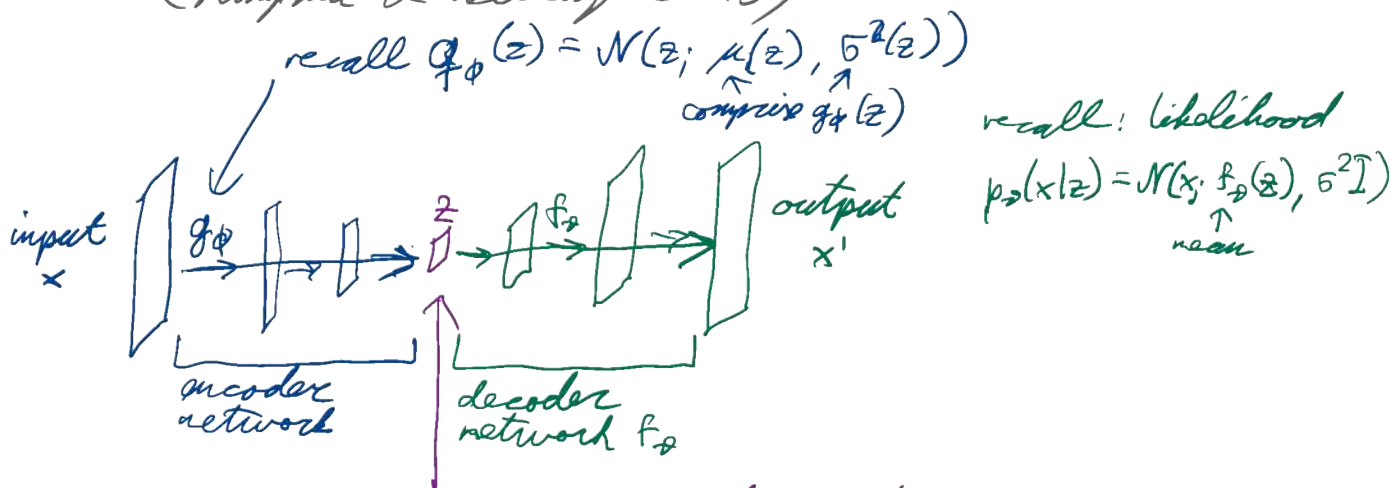
now **both** are global params

$\leq \log P_\vartheta(X=x)$

evidence

$\Rightarrow$ maximize $\mathbb{E}_{x \sim P_{data}}\left[ELBO(\vartheta, \phi)\right]$ over both $\vartheta, \phi$

often also just called "ELBO"

$\rightarrow$ "Amortized Variational Expectation Maximization"

= "Variational Autoencoders" (VAEs)

(Kingma & Welling 2013)

recall $Q_\phi(z) = \mathcal{N}(z; \mu(z), \sigma^2(z))$

comprise $g_\phi(z)$

recall: likelihood
$p_\vartheta(x|z) = \mathcal{N}(x; f_\vartheta(z), \sigma^2 I)$

mean



input x

$g_\phi$   $z$   $f_\vartheta$   output x'

encoder network

decoder network $f_\vartheta$

· minimise entropy of this part
· input wise here because we sample $z \sim Q_\phi(z|x)$

# Interpretations of the ELBO (i.e. the objective function)

$$ELBO(\theta, \phi) = \mathbb{E}_{z \sim Q_\phi(z|x)} \left[ \log p_\theta(z) + \log p_\theta(x|z) - \log q_\phi(z|x) \right]$$

we maximize this ↑

$$= + \mathbb{E}_{z \sim Q_\phi(z|x)} \left[ \log p_\theta(x|z) \right] - D_{KL}\left( Q_\phi(z|x) \| P_\theta(z) \right)$$

maximizing only this part would be maximum likelihood estimation (MLE)
→ it would make $Q_\phi(z|x)$ collapse to a $\delta$-function peaked at the MLE $= \arg\max_z \log p_\theta(x|z)$

think of this as a regularizer
→ tries to make $Q_\phi(z|x)$ similar to $P_\theta(z)$
→ at compression: want to encode $z$ using $P_\theta(z)$; this term ensures that $z$'s obtained from encoder have high $P_\theta(z)$

$$= \log P_\theta(X=x) - D_{KL}\left( Q_\phi(z|x) \| P_\theta(z|X=x) \right)$$

evidence
→ maximizing this minimizes the inf. content of $x$ under our model $P_\theta$, i.e., the theoretical lower bound of the bitrate

minimizing this makes the variational dist. $Q_\phi$ similar to the true posterior
⇒ $Q_\phi$ can be called the "approximate posterior"

→ **Goal:** maximize ELBO over $\theta$ & $\phi$

· issue: $ELBO(\theta, \phi) = \mathbb{E}_{z \sim Q_\phi(z|x)} \left\{ \cdots \right\}$

distribution from which we have to sample depends on $\phi$, by which we want to differentiate
→ see Problem set
(reparametrization grad: Kingma & Welling 2013
REINFORCE-gradients: Ranganath et al. 2014)

# Why all this fuss?

ongoing research on VI & related methods may be applicable to compression - or it may not be
⇒ look into that literature & try out if it improves compression methods

**Examples:**
- lots of research on tighter bounds of the evidence (tighter than the standard ELBO):
  → e.g. importance weighted VI, recently applied to compression by Theis & Ho 2021
- iterative amortized inference
  ↪ Marino et al 2018
  ↪ Campos et al 2019
- other approximate Bayesian inference methods (alternatives to VI) exist (in particular: Markov Chain Monte Carlo = MCMC)
  → nontrivial how to use these for compression (pioneering work: Havasi et al., 2018)

# References

[Beal and Ghahramani, 2003] Beal, M. J. and Ghahramani, Z. (2003). The variational bayesian em algorithm for incomplete data: with application to scoring graphical model structures. *Bayesian statistics*, 7(453-464):210.

[Blei et al., 2017] Blei, D. M., Kucukelbir, A., and McAuliffe, J. D. (2017). Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518).

[Campos et al., 2019] Campos, J., Meierhans, S., Djelouah, A., and Schroers, C. (2019). Content adaptive optimization for neural image compression. *arXiv preprint arXiv:1906.01223*.

[Dempster et al., 1977] Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1).

[Havasi et al., 2018] Havasi, M., Peharz, R., and Hernández-Lobato, J. M. (2018). Minimal random code learning: Getting bits back from compressed model parameters. *arXiv preprint arXiv:1810.00440*.

[Kingma and Welling, 2014] Kingma, D. P. and Welling, M. (2014). Auto-encoding variational Bayes. In *International Conference on Learning Representations*.

[Marino et al., 2018] Marino, J., Yue, Y., and Mandt, S. (2018). Iterative amortized inference. In *International Conference on Machine Learning*.

[Ranganath et al., 2014] Ranganath, R., Gerrish, S., and Blei, D. (2014). Black box variational inference. In *Artificial intelligence and statistics*, pages 814–822.

[Theis and Ho, 2021] Theis, L. and Ho, J. (2021). Importance weighted compression. In *Neural Compression: From Information Theory to Applications–Workshop@ ICLR 2021*.

[Zhang et al., 2018] Zhang, C., Bütepage, J., Kjellström, H., and Mandt, S. (2018). Advances in variational inference. *IEEE transactions on pattern analysis and machine intelligence*, 41(8).