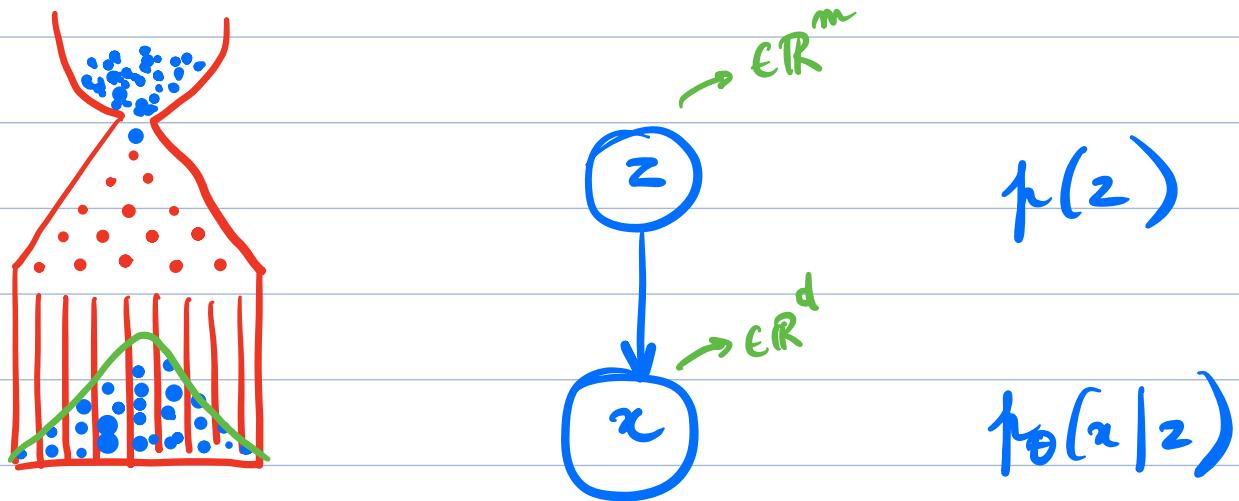


Deep Latent variable models



How to estimate θ ?

$$\begin{aligned}
 \max_{\theta} p_{\theta}(x) &= \int p_{\theta}(x, z) dz \\
 &= \int p(z) p_{\theta}(x|z) dz \\
 &= E_{p(z)} [p_{\theta}(x|z)] \\
 &\approx \frac{1}{K} \sum_k p_{\theta}(x|z_k)
 \end{aligned}$$

MC approximation

Poor approximation when m is large due to the curse of dimensionality!

I. Variational inference

$$p_{\theta}(x) = E_{q_{\phi}(z)} [p_{\theta}(x|z)]$$

↓ Bayes

$$= \mathbb{E}_{q_{\phi}(z)} \left[\frac{p(z)}{q_{\phi}(z)} p_{\theta}(z|z) \right] \quad \begin{matrix} \text{Improving} \\ \text{simplifying} \end{matrix}$$

$$\log p_{\theta}(x) = \log \mathbb{E}_{q_{\phi}(z)} \left[\frac{p(z)}{q_{\phi}(z)} p_{\theta}(z|x) \right]$$

James?

$\kappa \geq$

①

②

③

④

⑤

⑥

⑦

⑧

⑨

⑩

⑪

⑫

⑬

⑭

⑮

⑯

⑰

⑱

⑲

⑳

㉑

㉒

㉓

㉔

㉕

㉖

㉗

㉘

㉙

㉚

㉛

㉜

㉝

㉞

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

$\hat{\theta} \quad \theta^* \quad \theta$

We want $KL(q_\phi(z) || p_\theta(z|x)) \rightarrow 0$, otherwise
the gap is large and $|\theta^* - \hat{\theta}| \gg 0$.

\Rightarrow Put enough capacity in q_ϕ .

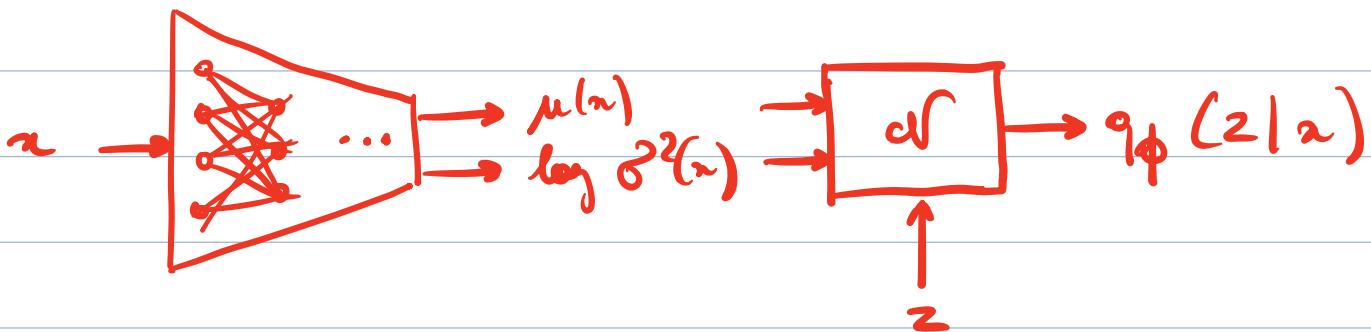
I. VAEs

Amortize inference for any x .

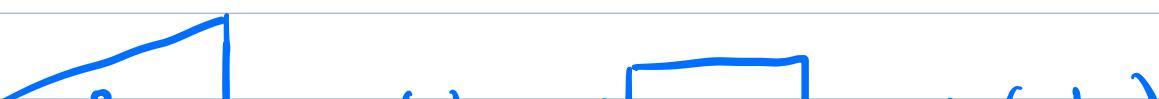
$$q_\phi(z) \longrightarrow q_\phi(z|x)$$

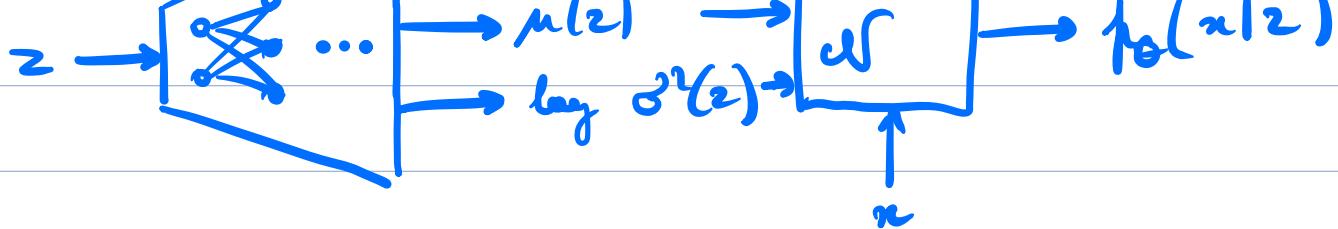
Now a DNN
that outputs the
parameters of the
variational dist.

Encoder



Decoder





Training

$$\underset{\theta, \phi}{\text{max}} \mathbb{E}_{p(x)} [\text{ELBO}(x; \theta, \phi)]$$

$$= \mathbb{E}_{p(x)} \left[\mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - \text{KL}(q_{\phi}(z|x) \parallel p(z)) \right]$$



Issue: posterior collapse
when $\text{KL} = 0$.

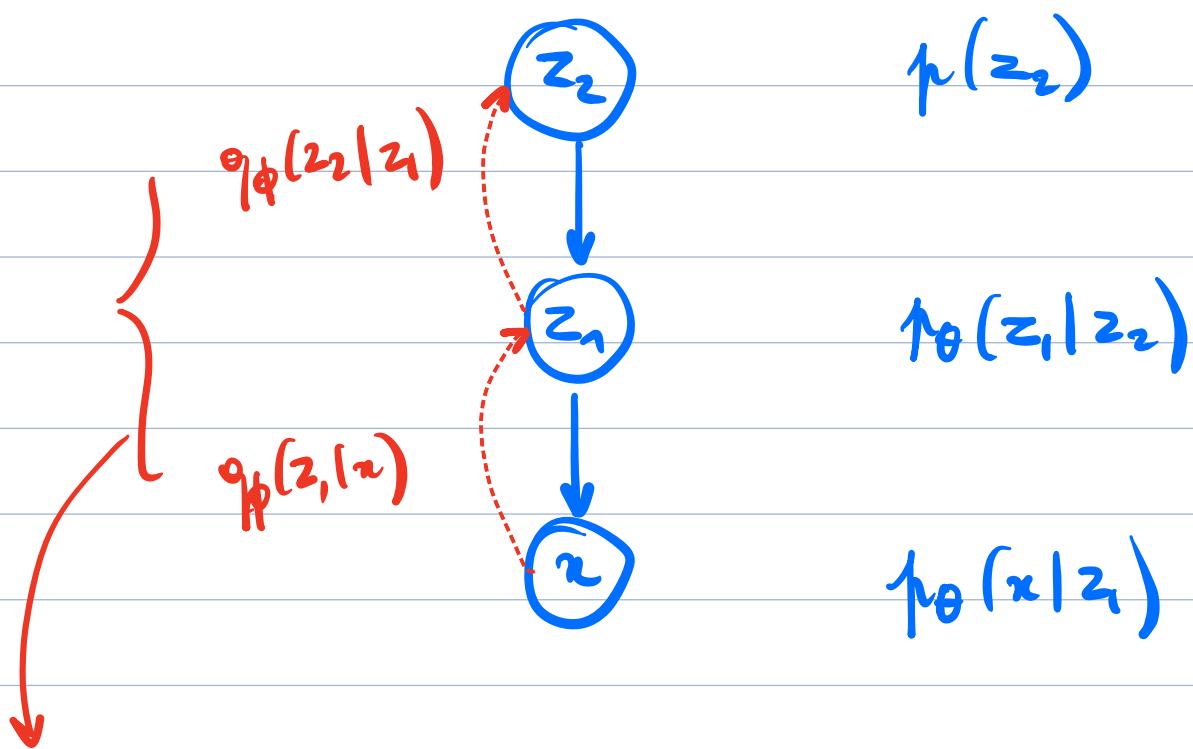
- x does not carry
any info about z
- z is poor noise



Tension between $\text{KL}(q_{\phi}(z|x) \parallel p(z))$
and $\text{KL}(q_{\phi}(z|x) \parallel p_{\theta}(z|x))$

→ Slides + Code.

II. Hierarchical VAEs



$$q_\phi(z_1, z_2 | x) = q_\phi(z_1 | x) q_\phi(z_2 | z_1)$$

Training

$$\max_{\phi, \theta} \mathbb{E}_{p(x)} [\text{ELBO}(x)]$$

$$\rightarrow p(x|z_1) p(z_1|z_2) p(z_2)$$

$$= \mathbb{E}_{p(x)} \mathbb{E}_{q_\phi(z_1, z_2 | x)} \left[\log \frac{p(x, z_1, z_2)}{q(z_1, z_2 | x)} \right]$$

do it!

$$= \mathbb{E}_{p(x)} \mathbb{E}_{q_\phi} (z_1, z_2 | x) \left[\log p_\theta(x | z_1) - \text{KL}(q_\phi(z_1 | x) || p(z_1 | z_2)) - \text{KL}(q_\phi(z_2 | z_1) || p(z_2)) \right]$$

Some
or done
but with
this now
team

$\rightarrow 0$ when $q_\phi(z_2 | z_1)$ has
too much capacity

$$\Rightarrow q_\phi(z_2 | z_1) = p(z) = \mathcal{N}$$

\rightarrow much lower in

me info about
 x in z_2

\Rightarrow the second layer is
not used!

\Rightarrow Some \rightarrow regular VAEs!

Top-down VAEs

! Swap the directional dependencies between
the latents

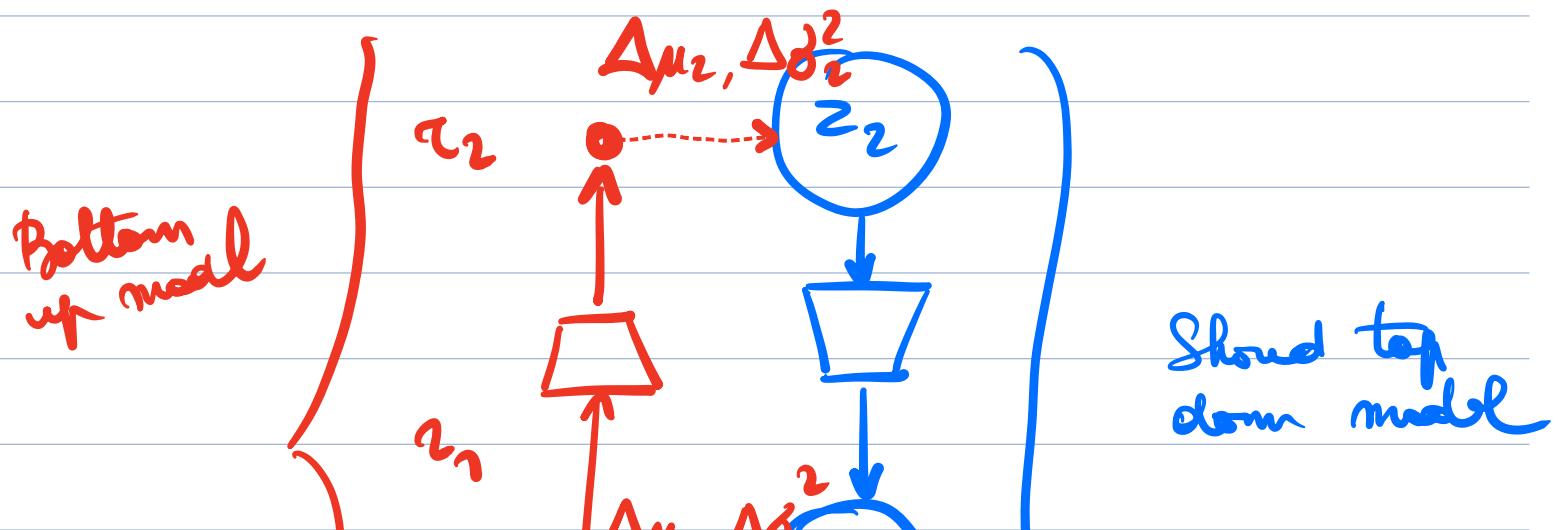
$$q_{\phi}(z_1, z_2 | x) = q_{\phi}(z_2 | x) q_{\phi}(z_1 | z_2, x)$$

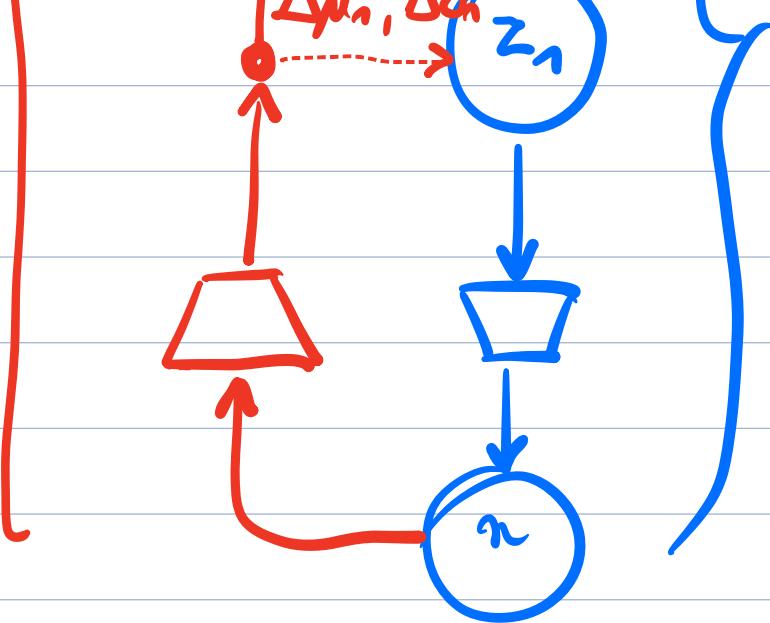


! Some directed
dependencies or
in $p_{\theta}(z_1 | z_2)$

INDUCTIVE
BIAS

\Rightarrow Share a common
top-down path.





Forces a close connection between q_{θ} and p_{θ}
and helps encode information about x
in z_2, z_1 .

→ Slides.

IV . Deep diffusion probabilistic models

Diffusion

④ W encoder
 $\uparrow q(z_t | z_{t-1})$



$\alpha \in \mathbb{R}$

$$\pi_\phi(z_t | z_{t+1}) \downarrow$$

DNN
decoder

$$\mu_z := w$$

 $T = 0(1000)$

Big first-order
Markov chain

Generation = Reverse diffusion

→ Show diffusion.

$$f_\theta(z_{0:T}) = \left[\prod_{t=0}^{T-1} \pi_\phi(z_t | z_{t+1}) \right] \mu(z_T)$$

$$q_\phi(z_{1:T} | z_0) = \prod_{t=1}^T q_\phi(z_t | z_{t-1})$$

$$q_\phi(z_t | z_{t-1}) = \mathcal{W}(z_t | \sqrt{1-\beta_t} z_{t-1}; \beta_t \mathcal{I})$$

$$\Leftrightarrow z_t := \sqrt{1-\beta_t} z_{t-1} + \beta_t \varepsilon, \quad \varepsilon \sim \mathcal{W}(0, 1)$$

Training

$$\max_{\phi, \theta} \mathbb{E}_{p(z_0)} [\text{ELBO}(\pi)]$$

$$= \mathbb{E}_{\mu(z_0)} \mathbb{E}_{q_\phi(z_{1:T} | z_0)} \left[\log \frac{f_\theta(z_{0:T})}{q_\phi(z_{1:T} | z_0)} \right]$$

$$\begin{aligned}
 &= \mathbb{E}_{p(z_0)} \mathbb{E}_{q_{\phi}(z_{1:T} | z_0)} \left[\log \prod_{t=1}^T p(z_t | z_1) \right. \\
 &\quad - \sum_{t=1}^T \text{KL}(q_{\phi}(z_t | z_{t-1}) \| p(z_t | z_{t+1})) \\
 &\quad \left. - \text{KL}(q_{\phi}(z_T | z_{T-1}) \| p(z_T)) \right]
 \end{aligned}$$

Han et al., 2020 :

[skip]
if time is short.

① Since q is linear Gaussian, we have

$$q(z_t | z_0) = \mathcal{N}(z_t | \sqrt{\bar{\alpha}_t} z_0, (1 - \bar{\alpha}_t) I)$$

$$\text{where } \bar{\alpha}_t = 1 - \beta_t$$

$$\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$$

$\Rightarrow z_t$ can be sampled without all intermediate steps!

$$\begin{aligned}
 \textcircled{2} \quad q(z_{t-1} | z_t, z_0) &= \frac{q(z_t | z_{t-1}, z_0) q(z_{t-1} | z_0)}{q(z_t | z_0)} \\
 &\Downarrow \text{ar} \quad \Downarrow \text{ar} \\
 &\quad (\text{conjugate prior})
 \end{aligned}$$

$$\begin{aligned}
 &= \mathcal{N}(z_{t-1} | \tilde{\mu}_t(z_t, z_0), \tilde{\beta}_t I) \\
 &\quad \xrightarrow{\text{ar}} \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t
 \end{aligned}$$

$$\frac{\sqrt{\bar{\alpha}_{t-1}} \beta_t}{1 - \bar{\alpha}_t} z_0 + \frac{\sqrt{\alpha_t} (1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} z_t$$

Then, the ELBO can be expressed as

$$\mathbb{E}_{\mu(z_0)} \mathbb{E}_{q(z_{1:T}|z_0)} \left[\log \prod_{t=1}^T p_\theta(z_t|z_1) \rightarrow \mathcal{L}_t \right. \\ \left. - \sum_{t=1}^T \text{KL}(q_\phi(z_{t-1}|z_t, z_0) \| p_\theta(z_{t-1}|z_t)) \right. \\ \left. - \text{KL}(q_\phi(z_T|z_0) \| p_\theta(z_T)) \right]$$

Training

$\text{KL}(\omega \| \tilde{\omega})$,
closed form!

$$\mathbb{E}_{t \sim [1..T]} \mathbb{E}_{\mu(z_0)} \mathbb{E}_{q(z_t|z_0)} [\mathcal{L}_t]$$

- Update the layer one at a time!
- much more memory efficient
- Scale to $T = O(1000)$

→ Slides

$$\mathcal{L}_{t-1} = \mathbb{E}_q \left[\frac{1}{2\sigma_t^2} \| \tilde{\mu}_t(x_t, z_0) - \mu_\theta(x_t, t) \|^2 \right]$$

II . VAE prior

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

Reconstruction Ω

$$\begin{aligned} \mathcal{L} &= \mathbb{E}_{p(x)} \mathbb{E}_{q_{\phi}(z|x)} \left[\log p(z) - \log q_{\phi}(z|x) \right] \\ &= \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{q_{\phi}(z|x_n)} \left[\log p(z) - \log q_{\phi}(z|x_n) \right] \\ &= \mathbb{E}_{\overset{N}{\underset{\sim}{p}}(z)} [\log p(z)] - \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{q_{\phi}(z|x_n)} [-\log q_{\phi}(z|x_n)] \\ &= CE(q_{\phi}(z) || p(z)) + H_{\overset{p(z)}{q_{\phi}(z)}} \end{aligned}$$

non-entropy entropy

$$q_{\phi}(z) = \frac{1}{N} \sum_{n=1}^N q_{\phi}(z|x_n)$$

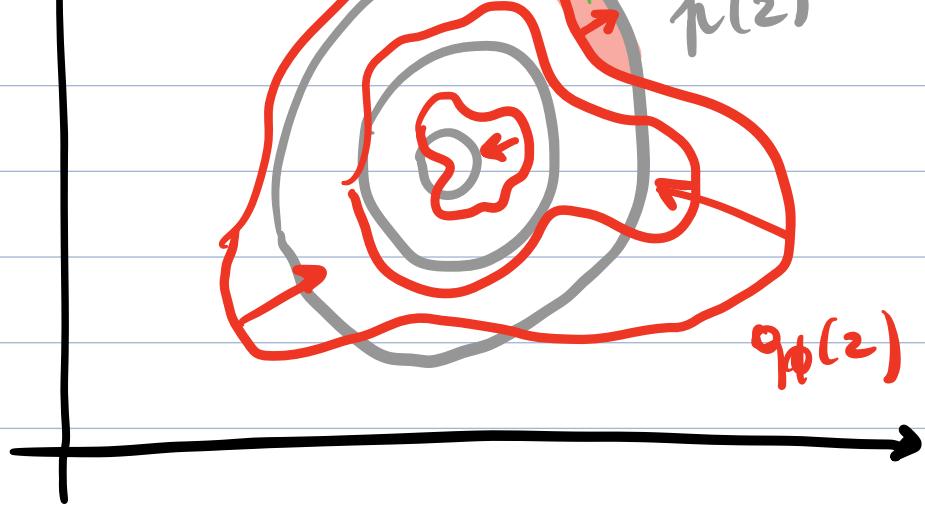
aggregated posterior

Make the aggregated posterior match with the prior.

Make the posterior variance $\rightarrow \infty$, but counter-balanced with RE.



Difficult when



"the prior is fixed since the decoder also forces the encoder to be peaky."

=> Learn the prior with another generative model!

$$p(z) \rightarrow p_x(z)$$

MoG
NF_j
DDPM,
...

→ Slides.

