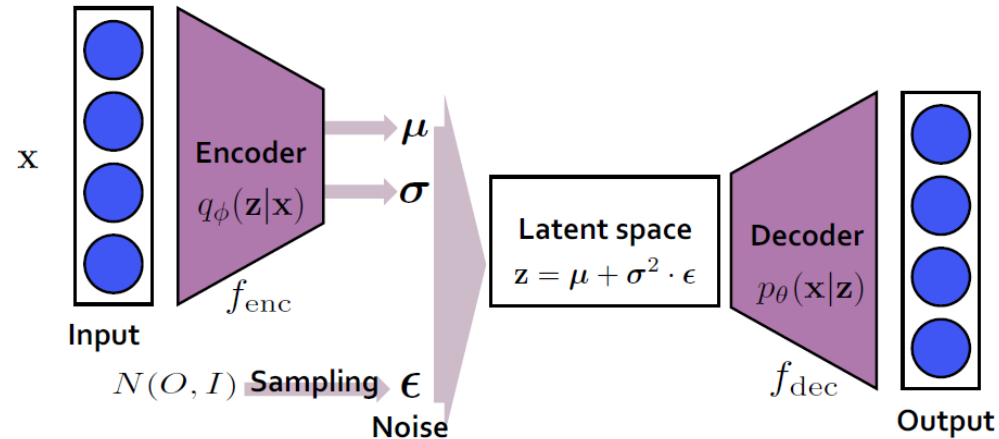


VAE: Rate-Distortion Curve & Aggregated Posterior

Presenter: Kim Seung Hwan (overnap@khu.ac.kr)

What is the Variational Autoencoder?

- Deep Maximum Likelihood Model
- Deep Latent Generative Model
- An autoencoder with gaussian noise regularization (with some fine math support!)
- A Decoder $p_{\theta}(x|z)$ and a $q_{\phi}(z|x)$ for a Data x and its latent representation z



Evidential Lower Bound?

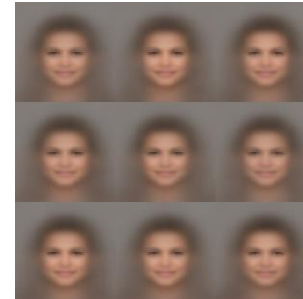
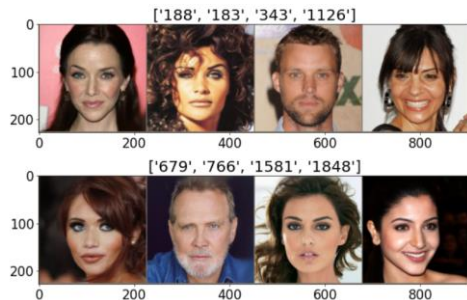
- Optimization target: ELBO (Evidence Lower Bound)
- ELBO consists of the reconstruction term and regularization term

$$-\log p_{\theta}(\mathbf{x}) \leq E_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})}[-\log p_{\theta}(\mathbf{x}|\mathbf{z})] + D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}))$$

- Since ELBO is the bound of the negative log-likelihood,
- VAE can be viewed as a deep maximum likelihood model

Pros and Cons of VAE

- Pros
 - Solid mathematical background
 - Lightweight; simple structure and implementation (compared to the Diffusion)
 - No need adversarial strategy (compared to the GAN)
 - Low-dimensional latent variable
- Cons
 - Posterior collapse (autodecoding-like behavior – always outputting the same)
 - Blurry output (bad reconstruction)
 - Poor sampling quality (samples from prior are noticeably worse than reconstruction)



Blame the ELBO

- Many VAE enhancements fix ELBO
- β -VAE

$$\mathcal{L}_\beta(\phi, \theta) = \underbrace{\mathbb{E}_{p_d(\mathbf{x})} [\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [-\log p_\theta(\mathbf{x}|\mathbf{z})]]}_{\text{Distortion}} + \beta \underbrace{\mathbb{E}_{p_d(\mathbf{x})} [D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}), p(\mathbf{z}))]}_{\text{Rate}},$$

- And...

$$\begin{aligned} \mathcal{L}(\theta, \phi) &= \mathbb{E}_{q_\phi(\mathbf{z}, \mathbf{x})} \left[\log \frac{p_\theta(\mathbf{x}, \mathbf{z})}{p_\theta(\mathbf{x})p(\mathbf{z})} + \log \frac{q_\phi(\mathbf{z})q(\mathbf{x})}{q_\phi(\mathbf{z}, \mathbf{x})} + \log \frac{p_\theta(\mathbf{x})}{q(\mathbf{x})} + \log \frac{p(\mathbf{z})}{q_\phi(\mathbf{z})} \right], \\ &= \mathbb{E}_{q_\phi(\mathbf{z}, \mathbf{x})} \left[\underbrace{\log \frac{p_\theta(\mathbf{x} | \mathbf{z})}{p_\theta(\mathbf{x})}}_{\textcircled{1}} - \underbrace{\log \frac{q_\phi(\mathbf{z} | \mathbf{x})}{q_\phi(\mathbf{z})}}_{\textcircled{2}} \right] - \underbrace{\text{KL}(q(\mathbf{x}) || p_\theta(\mathbf{x}))}_{\textcircled{3}} - \underbrace{\text{KL}(q_\phi(\mathbf{z}) || p(\mathbf{z}))}_{\textcircled{4}}. \end{aligned}$$

Paper	Objective
Kingma and Welling [2013], Rezende et al. [2014]	$\textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4}$
Higgins et al. [2016]	$\textcircled{1} + \textcircled{3} + \beta (\textcircled{2} + \textcircled{4})$
Kumar et al. [2017]	$\textcircled{1} + \textcircled{2} + \textcircled{3} + \lambda \textcircled{4}$
Zhao et al. [2017]	$\textcircled{1} + \textcircled{3} + \lambda \textcircled{4}$
Alemi et al. [2018], Burgess et al. [2018]	$\textcircled{1} + \textcircled{3} + \gamma (\textcircled{2} + \textcircled{4}) - C $
Gao et al. [2018]	$\textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} - \lambda \textcircled{2}^a$
Achille and Soatto [2018]	$\textcircled{1} + \textcircled{3} + \beta \textcircled{2} + \gamma \textcircled{A}^*$
Kim and Mnih [2018], Chen et al. [2018]	$\textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{6} + \beta \textcircled{A}^*$
HFVAE (this paper)	$\textcircled{1} + \textcircled{3} + \textcircled{ii} + \alpha \textcircled{2} + \beta \textcircled{A} + \gamma \textcircled{1}$

Table 1: Comparison of objectives in autoencoding deep generative models. The asterisk \textcircled{A}^* indicates that the prior factorizes, i.e. $p(\mathbf{z}) = \prod_d p(z_d)$. The notation $\textcircled{2}^a$ refers to restriction of the mutual information $\textcircled{2}$ to a subset of "Anchor" variables \mathbf{z}_a .

Limitations of β -VAE

- Let us look at the VAE as a lossy compression
- For a RD-curve possible with the model structure,
- The VAE and the β -VAE only get one point on it
- This means one must train a new model according to
- Architecture, dataset, purpose, class, condition... etc.
- Sensitive hyperparameters tire us all the time

$$\mathcal{L}_\beta(\phi, \theta) = \underbrace{\mathbb{E}_{p_d(\mathbf{x})}[\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[-\log p_\theta(\mathbf{x}|\mathbf{z})]]}_{\text{Distortion}} + \beta \underbrace{\mathbb{E}_{p_d(\mathbf{x})}[D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}), p(\mathbf{z}))]}_{\text{Rate}},$$

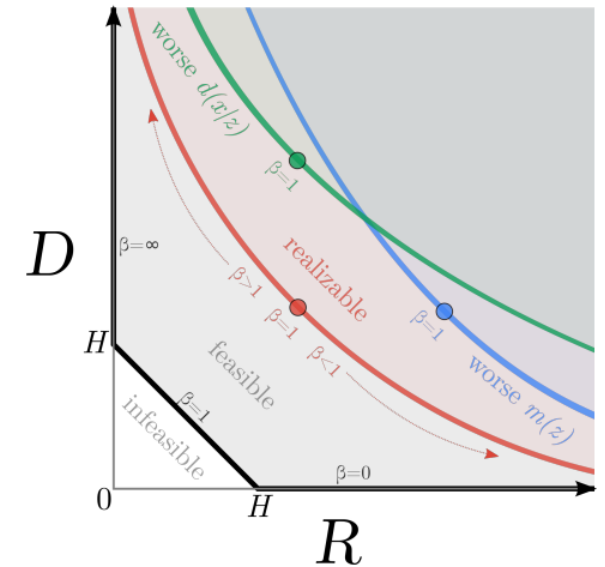
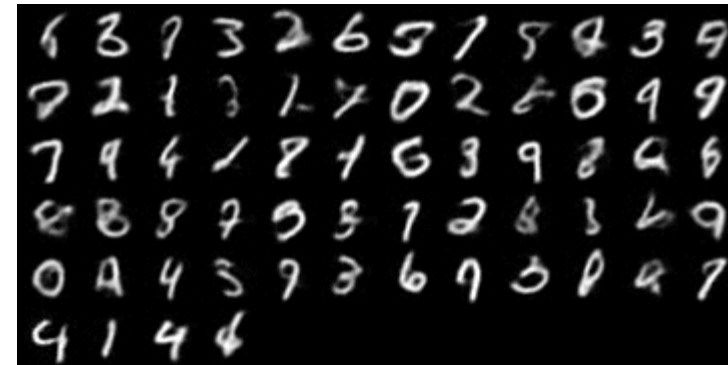
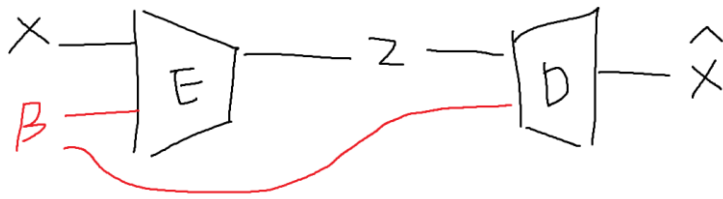


Figure 1. Schematic representation of the phase diagram in the RD -plane. The *distortion* (D) axis measures the reconstruction error of the samples in the training set. The *rate* (R) axis measures the relative KL divergence between the encoder and our own marginal approximation. The thick black lines denote the feasible boundary in the infinite model capacity limit.

Multi- β Model

- I believed that models with different β s would have strong associations
- So I implemented one VAE instance that can handle multiple β s
- This worked great but lacked novelty due to the existence of the prior work!
- I tried some analysis and proposed a novel heuristic, but it did not work well
- This is the study I did this summer



Rate-Distortion Curve

- This is the achieved RD-Curve, and the right side is from the references
- Why are these in log-like form?

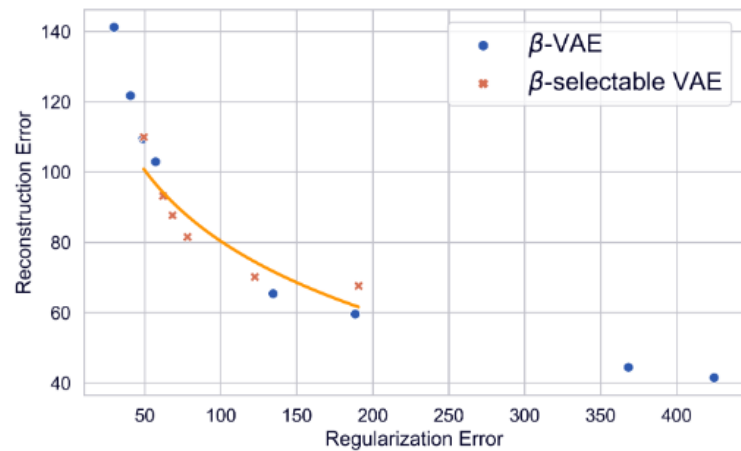


Figure 3: Reconstruction and regularization errors of ELBO for various β s. For ELBO, i.e. the learning objective, a single β -Selectable VAE is capable of approximate many β -VAEs. This graph can be interpreted as a rate-distortion curve.

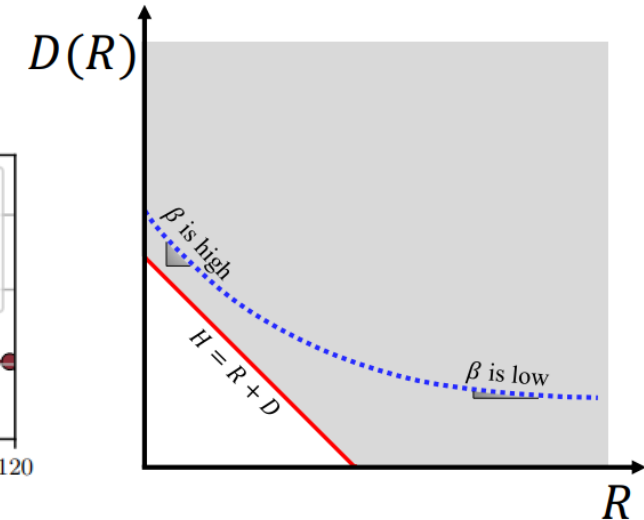
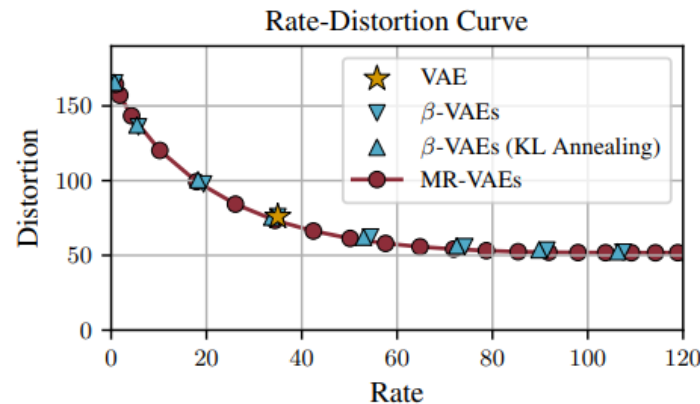


Figure 2: Schematic view of distortion-rate function. A red line corresponds to the theoretical lower bound of the rate and distortion. By varying β of β -VAE, we could achieve the points on a blue dashed curve, the sub-optimal distortion-rate function, which is best achievable with VAEs.

Decoupling σ_x and β

- The β -VAE has been popular and there are numerous analyses on it
- What I doubt is how they dealt with the variance of the decoder σ_x
- In terms of formula, $\sigma_x(z)$ has a similar meaning to β
- The difference is that the σ_x itself remains the likelihood

- But the typical beta-VAE is pretty much the same

- Original VAE loss:

$$E_{z \sim q_\phi(z|x)} [-\log p_\theta(x|z)] + D_{KL}(q_\phi(z|x) || p_\theta(z))$$

- Beta-VAE loss :

$$E_{z \sim q_\phi(z|x)} [MSE] + \beta D_{KL}(q_\phi(z|x) || p_\theta(z))$$

- Implemented correct ELBO :

$$E_{z \sim q_\phi(z|x)} \left[\frac{MSE}{2\sigma^2(z)} + \frac{1}{2} \log 2\pi\sigma^2(z) \right] + D_{KL}(q_\phi(z|x) || p_\theta(z))$$

Decoupling σ_x and β

- Many implementations keep σ_x constant and exclude it from optimization
- My claim is: studies so far have mixed the effects of σ_x and β
- Separating these explicitly will make the effect of each clearer
- This is the *first research direction*

- But the typical beta-VAE is pretty much the same

- Original VAE loss:

$$E_{z \sim q_\phi(z|x)} [-\log p_\theta(x|z)] + D_{KL}(q_\phi(z|x) || p_\theta(z))$$

- Beta-VAE loss :

$$E_{z \sim q_\phi(z|x)} [MSE] + \beta D_{KL}(q_\phi(z|x) || p_\theta(z))$$

- Implemented correct ELBO :

$$E_{z \sim q_\phi(z|x)} \left[\frac{MSE}{2\sigma^2(z)} + \frac{1}{2} \log 2\pi\sigma^2(z) \right] + D_{KL}(q_\phi(z|x) || p_\theta(z))$$

Decoupling σ_x and β

- Perhaps the analytically optimal $\sigma_x = \sqrt{MSE}$ works
- In this case, the reconstruction term becomes $\log MSE$
- The toy example had an unusual result: it has linear RD-Curve...?

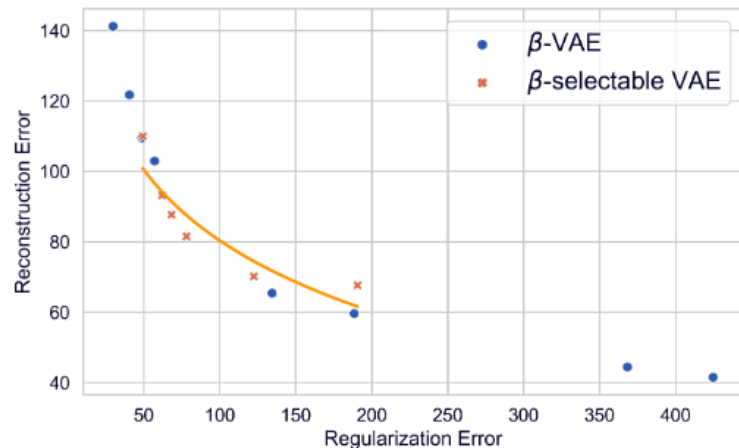


Figure 3: Reconstruction and regularization errors of ELBO for various β s. For ELBO, i.e. the learning objective, a single β -Selectable VAE is capable of approximate many β -VAEs. This graph can be interpreted as a rate-distortion curve.

- But the typical beta-VAE is pretty much the same

- Original VAE loss:

$$E_{z \sim q_\phi(z|x)} [-\log p_\theta(x|z)] + D_{KL}(q_\phi(z|x) || p_\theta(z))$$

- Beta-VAE loss :

$$E_{z \sim q_\phi(z|x)} [MSE] + \beta D_{KL}(q_\phi(z|x) || p_\theta(z))$$

- Implemented correct ELBO :

$$E_{z \sim q_\phi(z|x)} \left[\frac{MSE}{2\sigma^2(z)} + \frac{1}{2} \log 2\pi\sigma^2(z) \right] + D_{KL}(q_\phi(z|x) || p_\theta(z))$$

Aggregated Posterior

- This means $q_\phi(z) = \frac{1}{N} \sum_x q_\phi(z|x)$
- Sampling would be good if $q_\phi(z) \approx p(z)$
- In general, as β increases,
- i.e. the regularization term decreases,
- The above property is satisfied
- But the shape of the posterior is different
- For the same Rate...?

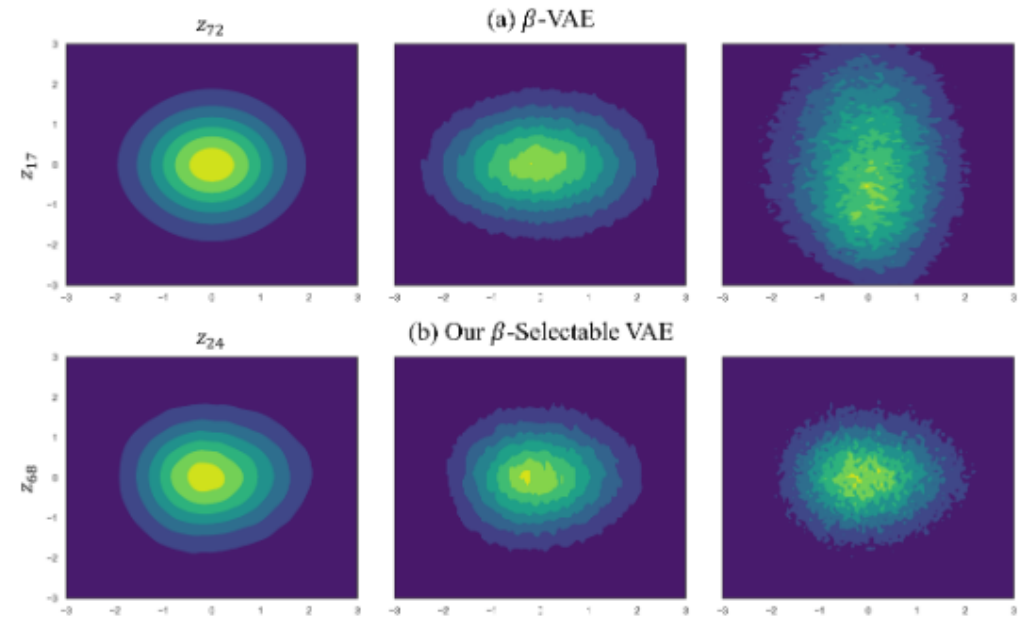


Figure 2: Aggregated posteriors for arbitrarily chosen latent channels. Each row shares a model and channels, and distributions are placed in order of decreasing β . The posteriors of our model (b) are *aligned* relative to their counterparts (a).

Aggregated Posterior

- If this is theoretically true, RD-Curve will not be able to express all of VAEs
- A new axis for some “sampling ease” should be added
- The candidate emerges from the decomposition mentioned earlier
- This is the *second research direction*

Table 1: FID scores of VAEs for both

	β (value or range)	CelebA
β -VAE	1.0	65.33
β -VAE	0.1	57.19
β -VAE	0.01	63.42
Ours	1.0	58.41
Ours	0.1	55.43
Ours	0.01	55.28

$$\mathcal{L}_\beta(\phi, \theta) = \underbrace{\mathbb{E}_{p_d(\mathbf{x})}[\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[-\log p_\theta(\mathbf{x}|\mathbf{z})]]}_{\text{Distortion}} + \beta \underbrace{\mathbb{E}_{p_d(\mathbf{x})}[D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}), p(\mathbf{z}))]}_{\text{Rate}},$$

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{q_\phi(\mathbf{z}, \mathbf{x})} \left[\log \frac{p_\theta(\mathbf{x}, \mathbf{z})}{p_\theta(\mathbf{x})p(\mathbf{z})} + \log \frac{q_\phi(\mathbf{z})q(\mathbf{x})}{q_\phi(\mathbf{z}, \mathbf{x})} + \log \frac{p_\theta(\mathbf{x})}{q(\mathbf{x})} + \log \frac{p(\mathbf{z})}{q_\phi(\mathbf{z})} \right],$$

$$= \mathbb{E}_{q_\phi(\mathbf{z}, \mathbf{x})} \left[\underbrace{\log \frac{p_\theta(\mathbf{x} | \mathbf{z})}{p_\theta(\mathbf{x})}}_{\textcircled{1}} - \underbrace{\log \frac{q_\phi(\mathbf{z} | \mathbf{x})}{q_\phi(\mathbf{z})}}_{\textcircled{2}} - \underbrace{\text{KL}(q(\mathbf{x}) || p_\theta(\mathbf{x}))}_{\textcircled{3}} - \underbrace{\text{KL}(q_\phi(\mathbf{z}) || p(\mathbf{z}))}_{\textcircled{4}} \right].$$

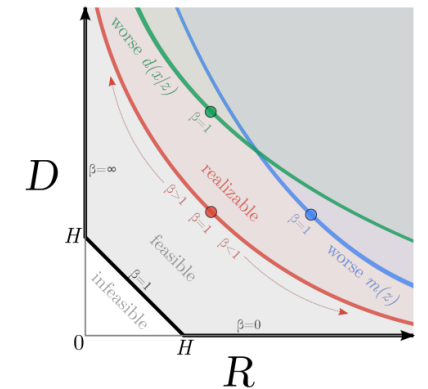


Figure 1. Schematic representation of the phase diagram in the RD-plane. The distortion (D) axis measures the reconstruction error of the samples in the training set. The rate (R) axis measures the relative KL divergence between the encoder and our own marginal approximation. The thick black lines denote the feasible boundary in the infinite model capacity limit.

Aggregated Posterior

- I might find a practical use for the observation: the posteriors are aligned
- The recursive heuristic I tried would be one of them (albeit with bad results)
- It is an interesting property, but I do not know how to utilize it
- This may be the third research direction

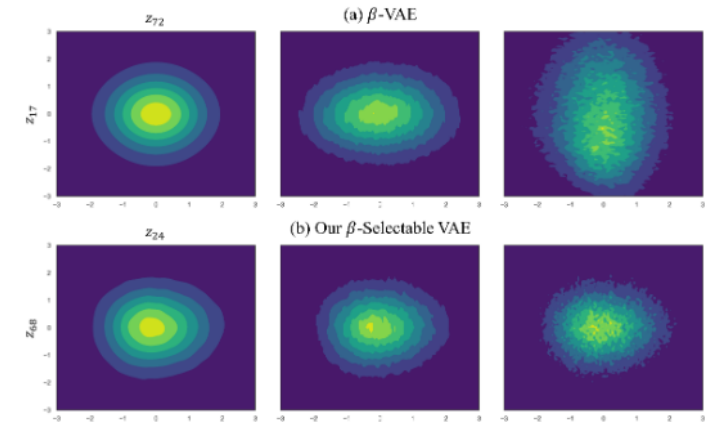


Figure 2: Aggregated posteriors for arbitrarily chosen latent channels. Each row shares a model and channels, and distributions are placed in order of decreasing β . The posteriors of our model (b) are *aligned* relative to their counterparts (a).

Thank you for listening

Presenter: Kim Seung Hwan (overnap@khu.ac.kr)