



Trading Information between Latents in Hierarchical Variational Autoencoders

Tim Z. Xiao, Robert Bamler

Department of Computer Science,

Cluster of Excellence "Machine Learning for Science", Tübingen AI Center,

University of Tübingen

Presentation on 11th International Conference on Learning Representations (ICLR), 2023



SCAN ME

Controlling Information in β -VAEs



$$\underbrace{\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})]}_{\text{-- distortion } D} - \underbrace{\beta D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}))}_{\text{rate } R}$$

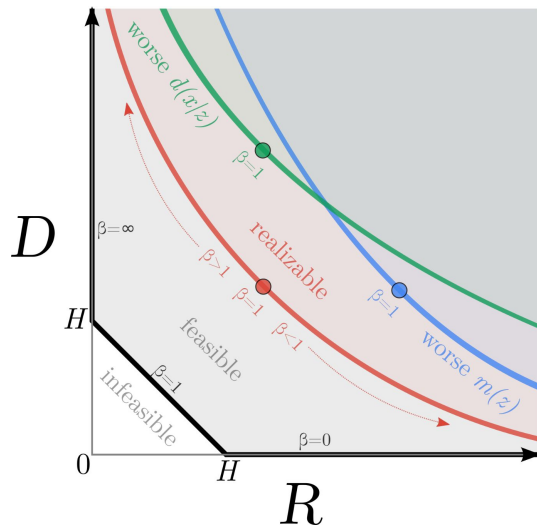
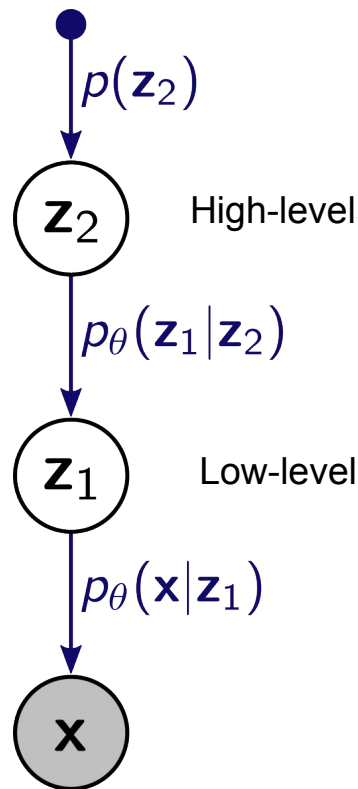


Figure taken from Alemi et al.,
Fixing a Broken ELBO, ICML 2018.



Defining Layer-Wise Bit Rates

For one architecture, total bit rate separates into:

$$R = R(z_L) + R(z_{L-1}|z_L) + R(z_{L-2}|z_{L-1}, z_L) + \dots + R(z_1|z_{\geq 2})$$

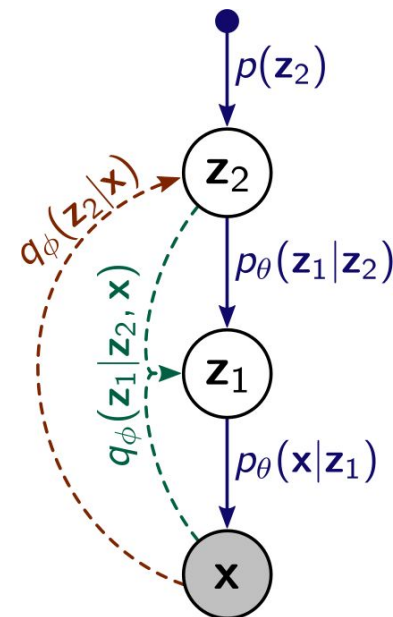
where:

$$R(z_\ell|z_{\geq \ell+1}) = \mathbb{E}_{q(z_{\geq \ell+1}|\mathbf{x})} [D_{\text{KL}}[q_\phi(z_\ell | z_{\geq \ell+1}, \mathbf{x}) \| p_\theta(z_\ell | z_{\geq \ell+1})]]$$

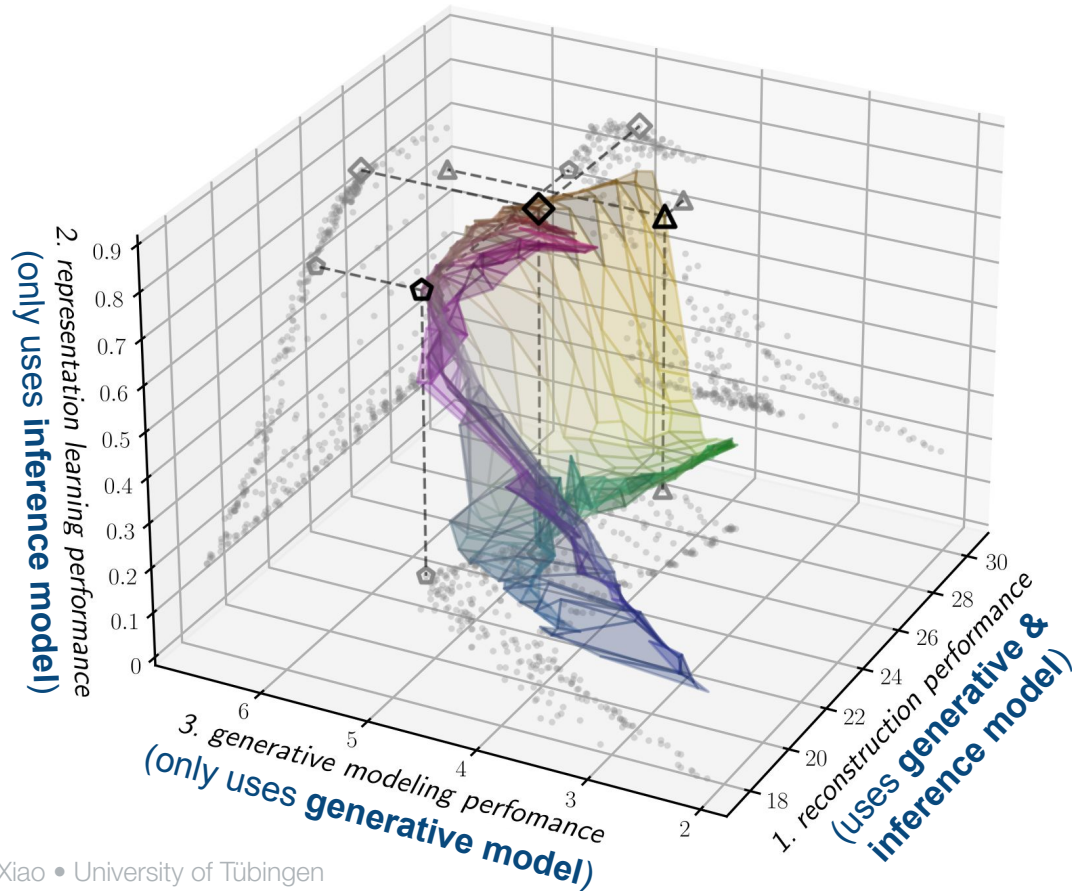
⇒ Proposed training objective:

$$\mathbb{E}_{\mathbf{x} \sim \mathbb{X}_{\text{train}}} [D + \beta_L R(z_L) + \beta_{L-1} R(z_{L-1}|z_L) + \dots + \beta_1 R(z_1|z_{\geq 2})]$$

L independent
Lagrange multipliers



There is no “One VAE Fits All”



diverse application domains



need fine-grained control
(no one-size-fits-all hierarchical VAE)

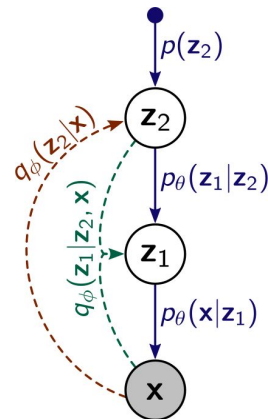
Application Type 1: Reconstruction



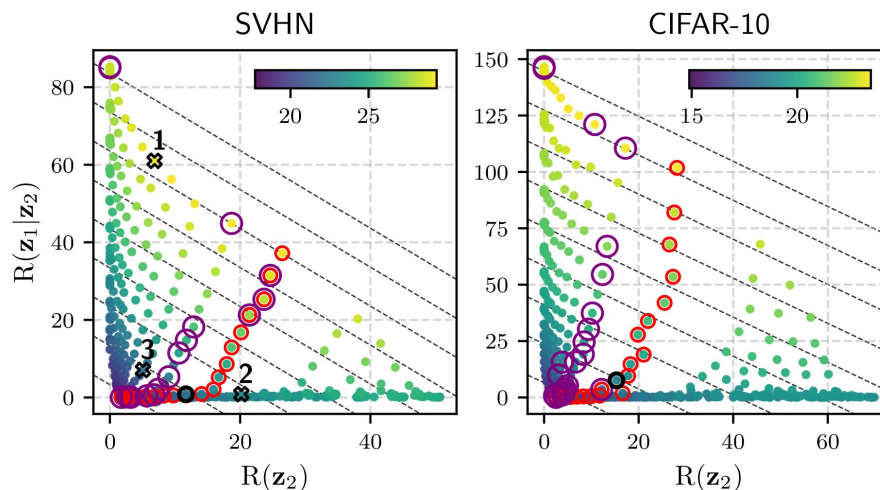
Theory:

$$\mathbb{E}_{\mathbf{x} \in p_{\text{data}}} [D] \geq H[p_{\text{data}}] - E_{\mathbf{x} \in p_{\text{data}}} [R(\mathbf{z}_L) + R(\mathbf{z}_{L-1}|\mathbf{z}_L) + \dots + R(\mathbf{z}_1|\mathbf{z}_{\geq 2})]$$

distortion



Experiment:



Application Type 2: Rep. Learning

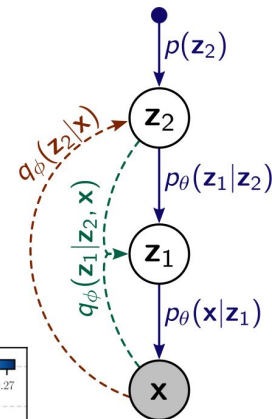
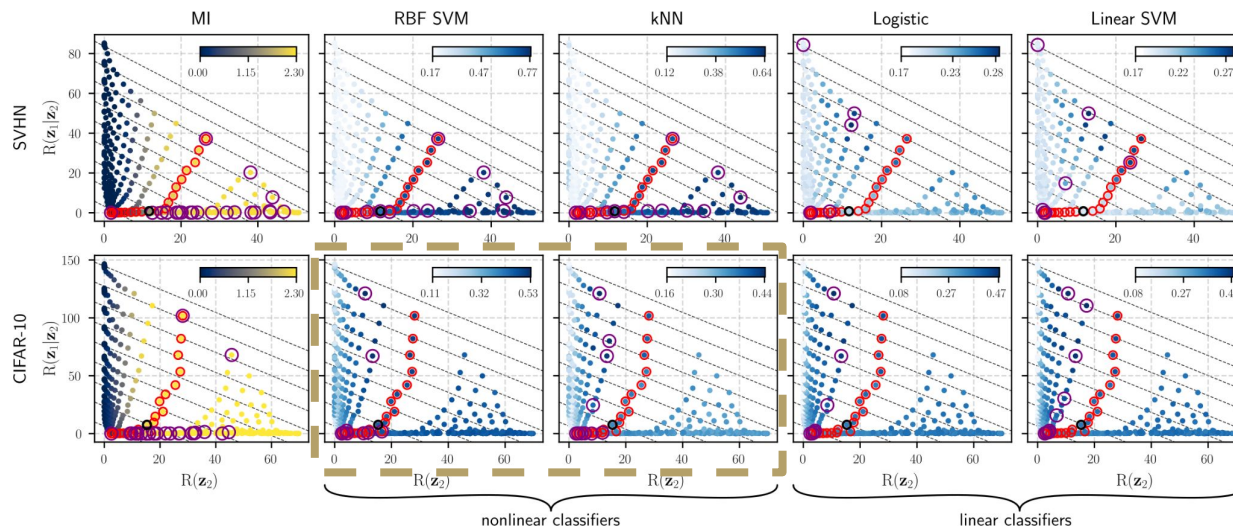


Theory: consider classifier operating on \mathbf{z}_2

$$\Rightarrow \text{accuracy} \leq f^{-1}(I_q(\text{label}; \mathbf{z}_2)) \leq f^{-1}(\mathbb{E}_{p_{\text{data}}}[R(\mathbf{z}_2)])$$

$$f(\alpha) = H[p_{\text{data}}(y)] + \alpha \log \alpha + (1 - \alpha) \log \frac{1 - \alpha}{M - 1} \quad [\text{analogous to Meyen, 2016}]$$

Experiment:



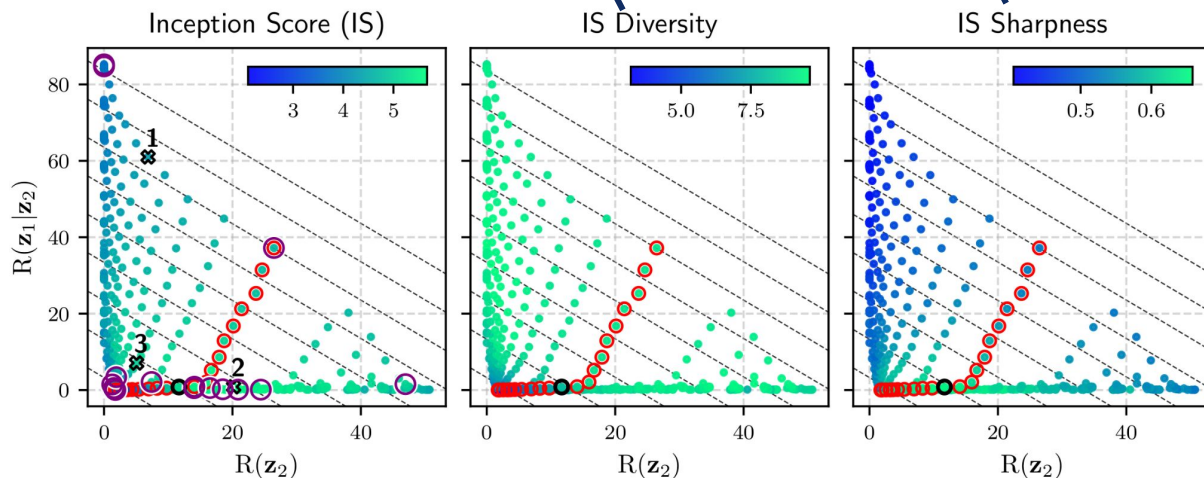
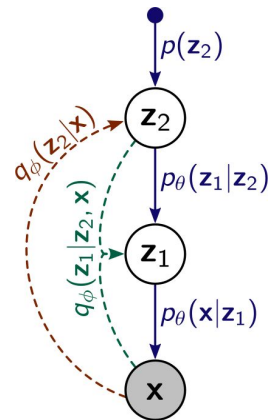
Application Type 3: Generation



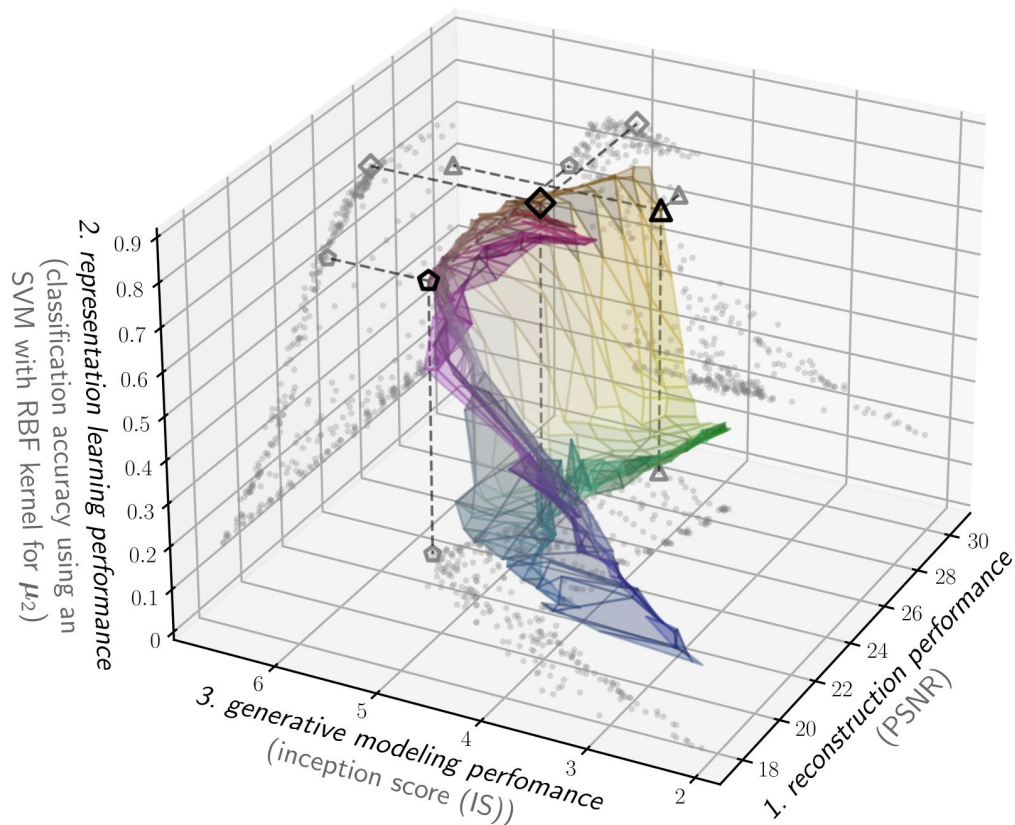
Theory: expect best generative performance when all $\beta = 1$

Experiment:
$$\text{IS} = \exp \left\{ \mathbb{E}_{p_{\theta}(\mathbf{x})} \left[D_{\text{KL}}[p_{\text{cls.}}(y|\mathbf{x}) \| p_{\text{cls.}}(y)] \right] \right\}$$

$$= e^{H[p_{\text{cls.}}(y)]} \times e^{-\mathbb{E}_{p_{\theta}(\mathbf{x})} [H[p_{\text{cls.}}(y|\mathbf{x})]]}$$
 [Salimans et al., 2016]



Summary



diverse application domains



need fine-grained control



control layer-wise rates