



**STEVENS**  
INSTITUTE *of* TECHNOLOGY  
THE INNOVATION UNIVERSITY®

# Learning Multimodal Latent Generative Models with Energy-Based Prior

Shiyu Yuan,<sup>[1]</sup> Jiali Cui,<sup>[2]</sup> Hanao Li,<sup>[2]</sup> Tian Han<sup>[2]</sup>

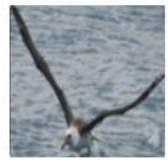
Department of Systems and Enterprises<sup>[1]</sup> Department of Computer Science<sup>[2]</sup>



# Research Problem

## 1) Joint Generation:

$$p(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)})$$

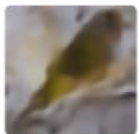


'a bird with a very long wingspan and a long pointed beak.'

inference

**Z**

generate



$\mathbf{x}^{(1)}$

an orange and small bird has brown strips on the rest of the body throughout

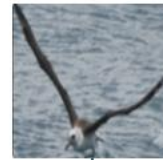
$\mathbf{x}^{(2)}$

$$\mathbf{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$$

## 2) Cross Generation:

$$p(\mathbf{x}^{(i)} | \mathbf{x}^{(j)}) \quad \text{where } (i \neq j)$$

$\mathbf{x}^{(1)}$



inference

**Z**

generate

'a bird with a very long wingspan and a long pointed beak.'

$\mathbf{x}^{(2)}$

$\mathbf{x}^{(1)} \rightarrow \mathbf{x}^{(2)}$

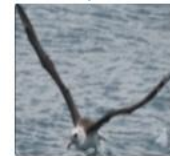
$\mathbf{x}^{(2)}$

'a bird with a very long wingspan and a long pointed beak.'

inference

**Z**

generate



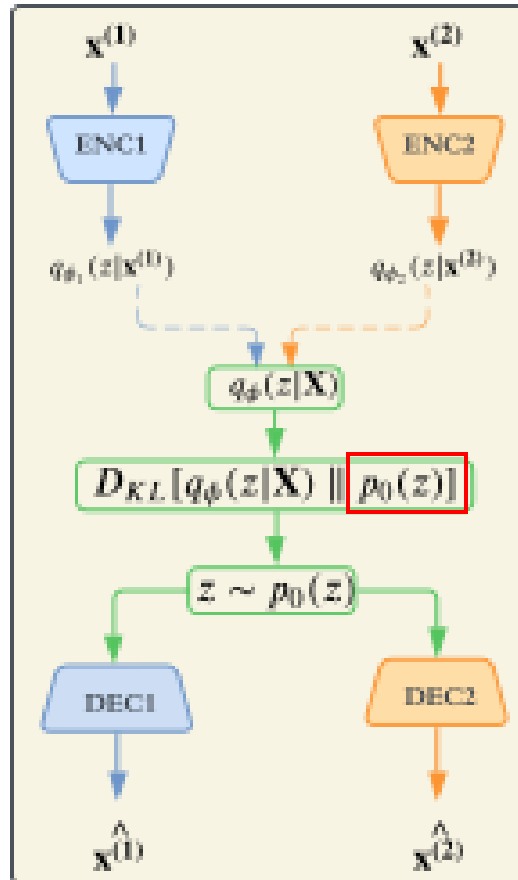
$\mathbf{x}^{(1)}$

$\mathbf{x}^{(2)} \rightarrow \mathbf{x}^{(1)}$

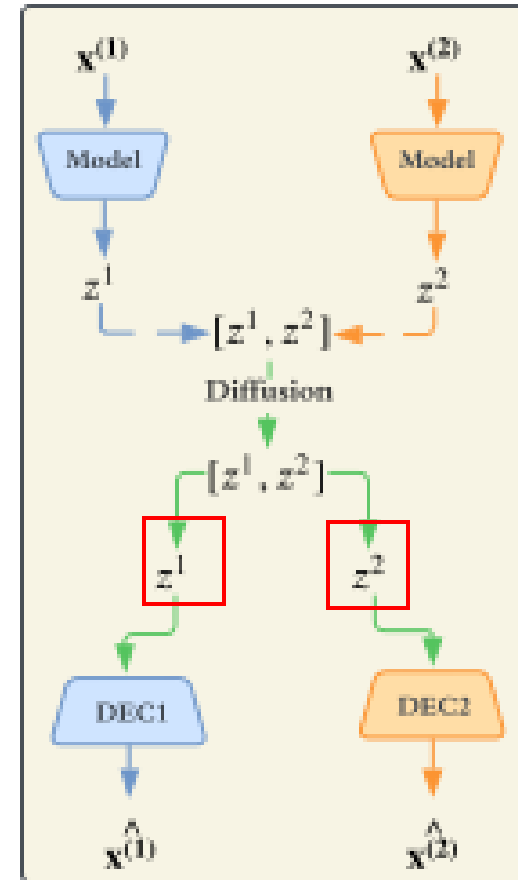


# Related Work

Variational Autoencoder-based <sup>[1,2,3,4]</sup>



Diffusion-based <sup>[5,6,7]</sup>



[1] Shi, Y et al. Variational mixture-of-experts autoencoders for multi-modal deep generative models. **NeurIPS2019**  
[2] Wu, M et al. Multimodal generative models for scalable weakly-supervised learning. **NeurIPS2018**  
[3] Sutter, T. M., et al. Generalized multimodal ELBO. **ICLR2021**  
[4] Palumbo, E., et al. MMVAE+: Enhancing the generative quality of multimodal VAEs without compromises. **ICLR2023**

[5] Bao, F., et al. One transformer fits all distributions in multi-modal diffusion at scale. **ICML2023**  
[6] Hu, M., Zheng, et al. UniD3: unified discrete diffusion for simultaneous vision-language generation. **ICLR2023**  
[7] Ramesh, A., et al. DALL-E 2



# Multimodal Latent Generative Model

---

$$p_{\theta}(\mathbf{X}, \mathbf{z}) = p_{\beta_{(1)}}(\mathbf{x}^{(1)} | \mathbf{z}) p_{\beta_{(2)}}(\mathbf{x}^{(2)} | \mathbf{z}) \cdots p_{\beta_{(m)}}(\mathbf{x}^{(m)} | \mathbf{z}) p(\mathbf{z})$$

$\mathbf{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$   $m$  : number of modalities

$\mathbf{z}$  shared representation among modalities

$p_{\beta_{(m)}}(\mathbf{x}^{(m)} | \mathbf{z})$  generation model : modality-specific information

$p(\mathbf{z})$  prior: common information



# Generation Model

generation model : modality-specific information



$$p_{\theta}(\mathbf{X}, \mathbf{z}) = p_{\beta_{(1)}}(\mathbf{x}^{(1)} | \mathbf{z}) p_{\beta_{(2)}}(\mathbf{x}^{(2)} | \mathbf{z}) \cdots p_{\beta_{(m)}}(\mathbf{x}^{(m)} | \mathbf{z}) p(\mathbf{z})$$

assume:  $\mathbf{x}^{(1)} \perp\!\!\!\perp \mathbf{x}^{(2)} \perp\!\!\!\perp \cdots \perp\!\!\!\perp \mathbf{x}^{(m)} | \mathbf{z}$

$$p_{\beta_{(m)}}(\mathbf{x}^{(m)} | \mathbf{z}) \sim \mathcal{N}(G_{\beta_{(m)}}(\mathbf{z}), I_{D^{(m)}})$$

$$\mathbf{x}^{(m)} = G_{\beta_{(m)}}(\mathbf{z}) + \epsilon \quad \text{where} \quad \epsilon \sim \mathcal{N}(0, I_{D^{(m)}})$$



modality-specific generator



# Towards more Expressive Prior: Energy-based prior

$$p_{\theta}(\mathbf{X}, \mathbf{z}) = p_{\beta_{(1)}}(\mathbf{x}^{(1)} | \mathbf{z}) p_{\beta_{(2)}}(\mathbf{x}^{(2)} | \mathbf{z}) \cdots p_{\beta_{(m)}}(\mathbf{x}^{(m)} | \mathbf{z}) p_{\alpha}(\mathbf{z})$$

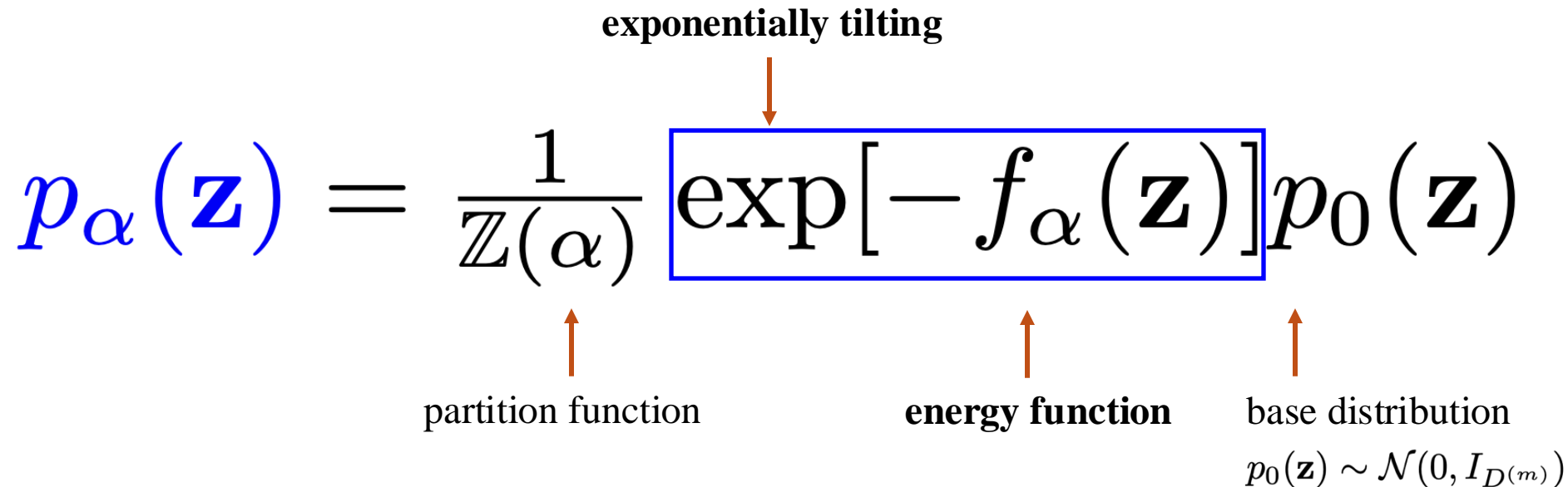
- **prior in existing work:** less-informative unimodal distribution such as Gaussian, Laplacian
- **EBM prior:** expressive prior to capture complexity of multimodal shared information
- **Exponential Tilting:** exponentially tilt modification of base distribution via energy function

$$p_{\alpha}(\mathbf{z}) = \frac{1}{\mathcal{Z}(\alpha)} \exp[-f_{\alpha}(\mathbf{z})] p_0(\mathbf{z})$$

↑                                  ↑                                  ↑

partition function                                  energy function                                  base distribution  
 $p_0(\mathbf{z}) \sim \mathcal{N}(0, I_{D^{(m)}})$

exponentially tilting





# Learning: Maximum Likelihood Estimation

---

$$\max_{\theta=\{\beta,\alpha\}} L(\theta) = \log \int_{\mathbf{z}} p_{\beta_{(1)}}(\mathbf{x}^{(1)}|\mathbf{z}) \cdots p_{\beta_{(m)}}(\mathbf{x}^{(m)}|\mathbf{z}) p_{\alpha}(\mathbf{z}) d\mathbf{z}$$

with sufficient data



$$\min_{\theta} \text{KL}(p_{\text{data}}(\mathbf{X}) \parallel p_{\theta}(\mathbf{X}))$$

$\beta$  : generator parameter

$\alpha$  : EBM prior parameter

---

$$\frac{\partial}{\partial \theta} L(\theta) = \mathbb{E}_{p_{\theta}(\mathbf{z}|\mathbf{X})} \left[ \frac{\partial}{\partial \theta} \log p_{\theta}(\mathbf{X}, \mathbf{z}) \right]$$

.....  
↓

$$q_{\phi}(\mathbf{z}|\mathbf{X}) = \frac{1}{m} \sum_{m=1}^M q_{\phi_{(m)}}(\mathbf{z}|\mathbf{x}^{(m)})^{[1]} \quad \text{Mixture of Expert (MOE)}$$

$\phi$  : inference model parameter 

# Final Objective for Joint Learning

---

$$\min_{\theta} \text{KL}(p_{\text{data}}(\mathbf{X}) \parallel p_{\theta}(\mathbf{X}))$$

$q_{\phi}(\mathbf{z}|\mathbf{X})$  ↓

$$\min_{\beta, \phi, \alpha} \text{KL}(p_{\text{data}}(\mathbf{X}) q_{\phi}(\mathbf{z}|\mathbf{X}) \parallel p_{\beta}(\mathbf{X}|\mathbf{z}) p_{\alpha}(\mathbf{z}))$$

↓

$$\max_{\beta, \phi, \alpha} \frac{1}{m} \sum_{m=1}^M \mathbb{E}_{q_{\phi(m)}(\mathbf{z}|\mathbf{x}^{(m)})} \left[ \log \frac{p_{\beta, \alpha}(\mathbf{X}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{X})} \right]$$

$\beta$  : generator parameter

$\phi$  : inference model parameter

$\alpha$  : EBM prior parameter





# Learning Generator and Inference Model

$$\frac{\partial}{\partial \beta, \phi} L(\beta, \phi, \alpha) = \frac{\partial}{\partial \beta, \phi} \frac{1}{m} \sum_{m=1}^M \left[ \mathbb{E}_{q_{\phi(m)}(\mathbf{z}|\mathbf{x}^{(m)})} [\log p_{\beta(m)}(\mathbf{x}^{(m)}|\mathbf{z})] \right] \leftarrow \text{Modality-Specific Reconstruction}$$

Cross-Modality Generation  $\rightarrow$   $+ \sum_{n=1, n \neq m}^M \mathbb{E}_{q_{\phi(m)}(\mathbf{z}|\mathbf{x}^{(m)})} [\log p_{\beta(n)}(\mathbf{x}^{(n)}|\mathbf{z})]$

Regularization with EBM Prior  $\rightarrow$   $+ \mathbb{E}_{q_{\phi(m)}(\mathbf{z}|\mathbf{x}^{(m)})} \left[ \log \frac{p_{\alpha}(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{X})} \right]$



# Learning EBM

---

$$\frac{\partial}{\partial \alpha} L(\beta, \phi, \alpha) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{X})} \left[ \frac{\partial}{\partial \alpha} f_{\alpha}(\mathbf{z}) \right] - \mathbb{E}_{p_{\alpha}(\mathbf{z})} \left[ \frac{\partial}{\partial \alpha} f_{\alpha}(\mathbf{z}) \right]$$

MOE: variation inference

Langevin dynamics

---

## Sampling from EBM: Langevin dynamics

$$\mathbf{z}_{\tau+1} = \mathbf{z}_{\tau} + \frac{s^2}{2} \frac{\partial}{\partial \mathbf{z}} [\log p_{\alpha}(\mathbf{z}_{\tau})] + s \cdot \epsilon_{\tau} \quad \text{where} \quad \epsilon_{\tau} \sim \mathcal{N}(0, I_d)$$

latent variable at  $\tau$    step size   sampled distribution   Gaussian noise

$\tau$  : time step                          EBM prior



# Model Generalization

## EBM Prior: Base Version

$$p_{\theta}(\mathbf{X}, \mathbf{z}) = p_{\beta}(\mathbf{X}|\mathbf{z})p_{\alpha}(\mathbf{z})$$

## MOE with modality prior

$$p_{\theta}(\mathbf{X}, \mathbf{z}, \mathbf{W}) = p_{\beta}(\mathbf{X}|\mathbf{z}, \mathbf{W})p_0(\mathbf{z})p_0(\mathbf{W}) \quad \mathbf{W} = \{\mathbf{w}^{(1)}, \dots, \mathbf{w}^{(m)}\}$$

↑ Laplacian prior    ↑ modality-specific prior

## EBM Prior: Generalized Version

$$p_{\theta}(\mathbf{X}, \mathbf{z}, \mathbf{W}) = p_{\beta}(\mathbf{X}|\mathbf{z}, \mathbf{W})p_{\alpha}(\mathbf{z})p_0(\mathbf{W})$$

↑ EBM prior    ↑ modality-specific prior



# PolyMNIST<sup>[3]</sup> : Coherence

- **Joint Coherence:** generated samples modalities alignment and mutually consistent
- **Cross Coherence:** capacity of one modality infer other modalities

## EBM Prior: Base Version

Model	Joint Coherence $\uparrow$	Cross Coherence $\uparrow$
	PolyMNIST	
Ours	0.746	0.853
MMVAE <sub>[1]</sub>	0.232	0.844

## EBM Prior: Generalized Version

Model	Joint Coherence $\uparrow$	Cross Coherence $\uparrow$
	PolyMNIST	
Ours	0.878	0.897
MMVAE <sub>+[2]</sub>	0.344	0.869
MoPoE <sub>[3]</sub>	0.141	0.720
MVTCAE <sub>[4]</sub>	0.003	0.591
mmJSD <sub>[5]</sub>	0.060	0.778

[1] Shi, Y et al. Variational mixture-of-experts autoencoders for multi-modal deep generative models. **NeurIPS2019**

[2] Palumbo, E., et al. MMVAE+: Enhancing the generative quality of multimodal VAEs without compromises. **ICLR2023**

[3] Sutter, T. M., et al. Generalized multimodal ELBO. **ICLR2021**

[4] Hwang, H., et al. Multi-view representation learning via total correlation objective. **NeurIPS2021**

[5] Sutter, T., et al. Multimodal generative learning utilizing jensen-shannon-divergence. **NeurIPS2020**



# PolyMNIST: Joint Generation Visual Result

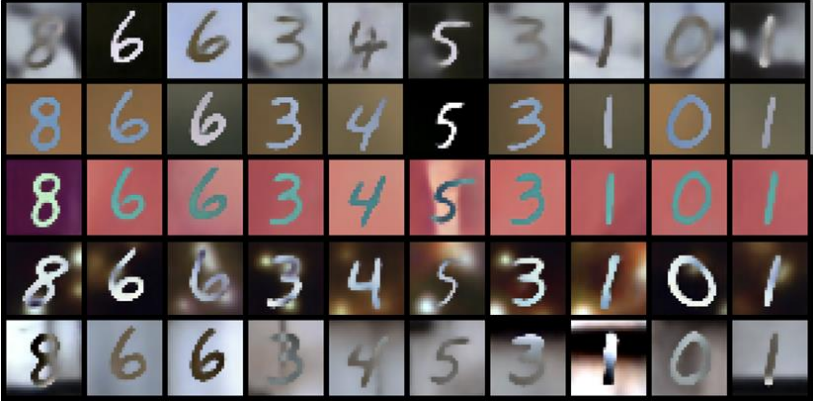
EBM Prior: Base Version



MMVAE



EBM Prior: Generalized Version



MMVAE+



# CUB.<sup>[1]</sup> : Markov Transition

$$z \sim p_0(\mathbf{z}) \xrightarrow{z_{\tau+1} = z_{\tau} + \frac{s^2}{2} \frac{\partial}{\partial \mathbf{z}} [\log p_{\alpha}(\mathbf{z}_{\tau})] + s \cdot \epsilon_{\tau}} z \sim p_{\alpha}(\mathbf{z})$$



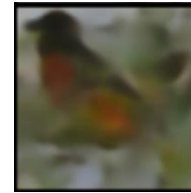
this bird is a black and are and and a very red ..



this bird has wings that are black and are very red beak



this small has has a are breast and white belly



this small has a orange breast and has a black belly

$z \sim p_0(\mathbf{z})$



$z \sim p_{\alpha}(\mathbf{z})$



Questions please reach out to:  
**syuan14@stevens.edu**

