

# Incomplete Multimodal Learning

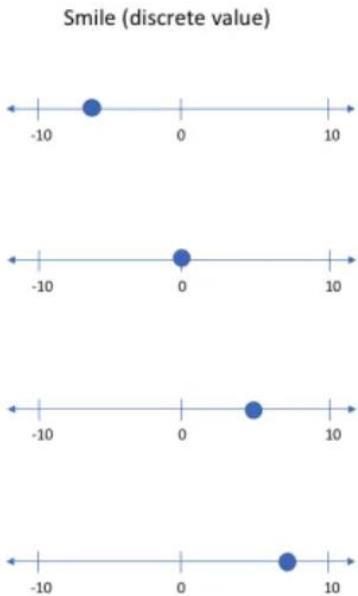
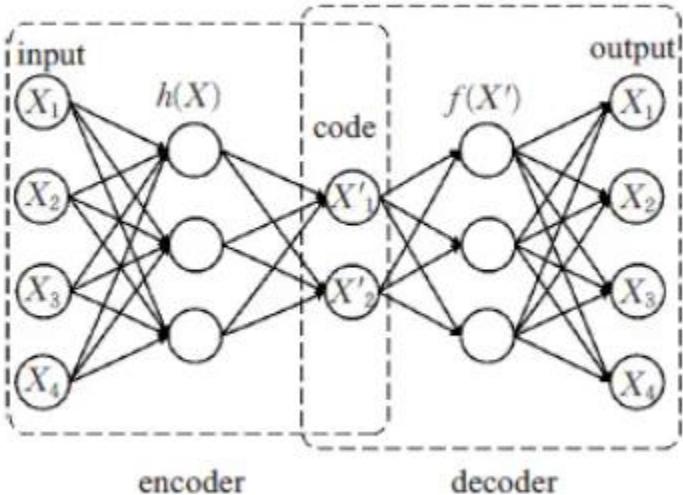
LI Xiaolou

2023/4/28

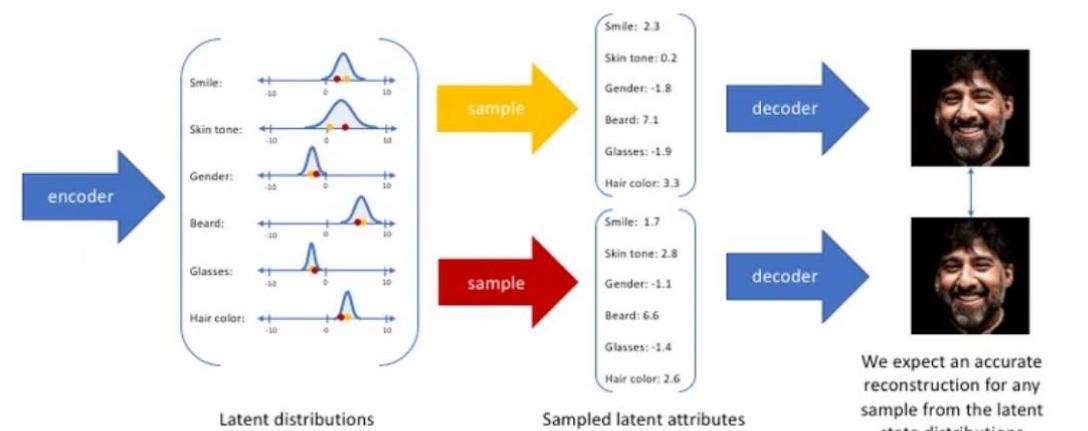
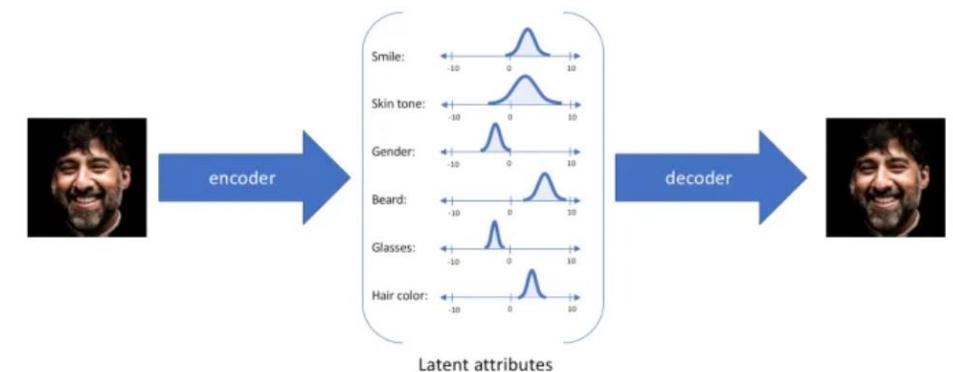
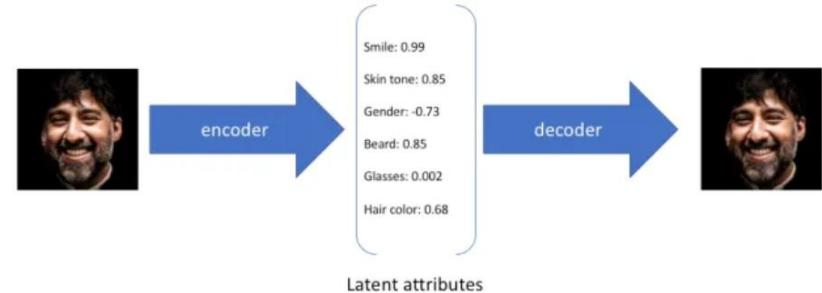
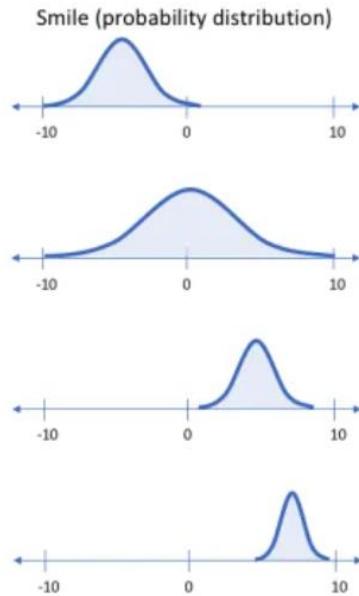
# Paper

1. (VAE) Doersch, C., 2016. Tutorial on variational autoencoders. arXiv preprint arXiv:1606.05908.
2. (MVAE) Wu, M. and Goodman, N., 2018. Multimodal generative models for scalable weakly-supervised learning. Advances in neural information processing systems, 31.
3. (IWAE) Burda, Y., Grosse, R. and Salakhutdinov, R., 2015. Importance weighted autoencoders. arXiv preprint arXiv:1509.00519.
4. (MMVAE) Shi, Y., Paige, B. and Torr, P., 2019. Variational mixture-of-experts autoencoders for multi-modal deep generative models. Advances in Neural Information Processing Systems, 32.

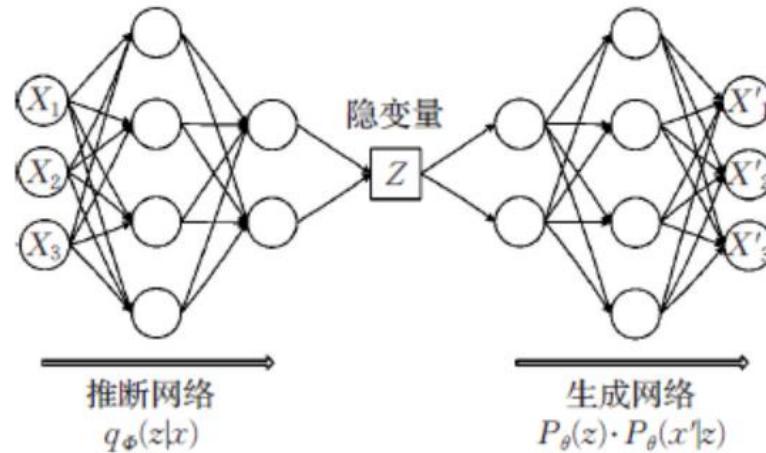
# VAE



vs.



# VAE



Original Data:  $X = \{x_i\}_{i=1}^N$

Why  $q(z|x)$ ?

Generated Data:  $X' = \{x'_i\}_{i=1}^N$

$$P(x_1, x_2, x_3, \dots, x_n) = P(x_1)P(x_2|x_1)P(x_3|x_2, x_1) \dots P(x_n|x_{n-1}, x_{n-2}, x_{n-3}, \dots, x_1)$$

↓ Mean Field Theory

Inference Net:  $q_\phi(z|x)$

$$Q(x_1, x_2, x_3, \dots, x_n) = Q(x_1)Q(x_2)Q(x_3) \dots Q(x_n)$$

Generation Net:  $P_\theta(z)P_\theta(x'|z)$

# VAE

$$q(z|x) \rightarrow p(z|x)?$$

minKL( $q(z|x)||p(z|x)$ )

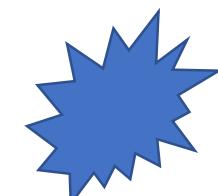
$L^v$ : evidence lower bound

$$\begin{aligned} L &= \log(p(x)) \\ &= \sum_z q(z|x) \log(p(x)) \\ &= \sum_z q(z|x) \log\left(\frac{p(z,x)}{p(z|x)}\right) \\ &= \sum_z q(z|x) \log\left(\frac{p(z,x)}{q(z|x)} \frac{q(z|x)}{p(z|x)}\right) \\ &= \sum_z q(z|x) \log\left(\frac{p(z,x)}{q(z|x)}\right) + \sum_z q(z|x) \log\left(\frac{q(z|x)}{p(z|x)}\right) \\ &= L^v + D_{KL}(q(z|x)||p(z|x)) \end{aligned}$$

$$\begin{aligned} L^v &= \sum_z q(z|x) \log\left(\frac{p(z,x)}{q(z|x)}\right) \\ &= \sum_z q(z|x) \log\left(\frac{p(x|z)p(z)}{q(z|x)}\right) \\ &= \sum_z q(z|x) \log\left(\frac{p(z)}{q(z|x)}\right) + \sum_z q(z|x) \log(p(x|z)) \end{aligned}$$

$$L^v = -D_{KL}(q(z|x)||p(z)) + \mathbb{E}_{q(z|x)}(\log(p(x|z)))$$

$$L_1 = \frac{1}{2} \sum_{j=1}^J [1 + \log((\sigma_j)^2) - (u_j)^2 - (\sigma_j)^2]$$



$$z \sim N(\mu, \sigma^2)$$

$$L_2 = \mathbb{E}_{q(z|x)}(\log(p(x|z))) \approx \frac{1}{L} \sum_{l=1}^L \log p(x|z^{(l)})$$

# MVAE

## Motivation:

The incompleteness of the mode is very common, and the work done by the predecessors either does not learn the joint distribution or requires additional calculations. And this paper uses a multimodal VAE

## Method:

1. MVAE
2. PoE (Product-of-expert inference network)
3. Sub-sampled training

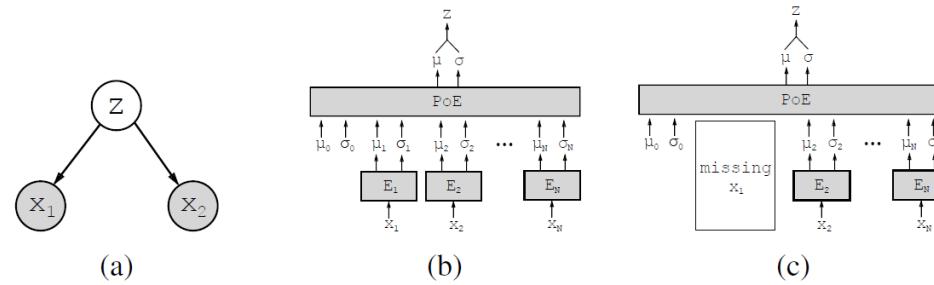


Figure 1: (a) Graphical model of the MVAE. Gray circles represent observed variables. (b) MVAE architecture with  $N$  modalities.  $E_i$  represents the  $i$ -th inference network;  $\mu_i$  and  $\sigma_i$  represent the  $i$ -th variational parameters;  $\mu_0$  and  $\sigma_0$  represent the prior parameters. The product-of-experts (PoE) combines all variational parameters in a principled and efficient manner. (c) If a modality is missing during training, we drop the respective inference network. Thus, the parameters of  $E_1, \dots, E_N$  are shared across different combinations of missing inputs.

$$\text{ELBO}(x) \triangleq \mathbb{E}_{q_\phi(z|x)}[\lambda \log p_\theta(x|z)] - \beta \text{KL}[q_\phi(z|x), p(z)]$$

conditional independence  
 $p(x_1, x_2, \dots, x_n, z) = p(z)p(x_1|z)p(x_2|z) \dots p(x_n|z)$

$$\text{ELBO}(X) \triangleq \mathbb{E}_{q_\phi(z|X)}\left[\sum_{x_i \in X} \lambda_i \log p_\theta(x_i|z)\right] - \beta \text{KL}[q_\phi(z|X), p(z)].$$

# MVAE

## Approximating the joint posterior

$$\begin{aligned} p(z|x_1, \dots, x_N) &= \frac{p(x_1, \dots, x_N|z)p(z)}{p(x_1, \dots, x_N)} = \frac{p(z)}{p(x_1, \dots, x_N)} \prod_{i=1}^N p(x_i|z) \\ &= \frac{p(z)}{p(x_1, \dots, x_N)} \prod_{i=1}^N \frac{p(z|x_i)p(x_i)}{p(z)} = \frac{\prod_{i=1}^N p(z|x_i)}{\prod_{i=1}^{N-1} p(z)} \cdot \frac{\prod_{i=1}^N p(x_i)}{p(x_1, \dots, x_N)} \quad (3) \\ &\propto \frac{\prod_{i=1}^N p(z|x_i)}{\prod_{i=1}^{N-1} p(z)} \end{aligned}$$

Alternatively, if we approximate  $p(z|x_i)$  with  $q(z|x_i) \equiv \tilde{q}(z|x_i)p(z)$ , where  $\tilde{q}(z|x_i)$  is the underlying inference network, we can avoid the quotient term:

$$p(z|x_1, \dots, x_N) \propto \frac{\prod_{i=1}^N p(z|x_i)}{\prod_{i=1}^{N-1} p(z)} \approx \frac{\prod_{i=1}^N [\tilde{q}(z|x_i)p(z)]}{\prod_{i=1}^{N-1} p(z)} = p(z) \prod_{i=1}^N \tilde{q}(z|x_i). \quad (4)$$



## Why Sub-sampled Training?

1. Can't handle modality lost
2. Unable to get inter-modal info

- the ELBO using the product of all N Gaussians
- all ELBO terms using a single modality
- k ELBO terms using k randomly chosen subsets

$$\text{ELBO}(x_1, \dots, x_N) + \sum_{i=1}^N \text{ELBO}(x_i) + \sum_{j=1}^k \text{ELBO}(X_j)$$

# MVAE

Model	BinaryMNIST	MNIST	FashionMNIST	MultiMNIST	CelebA
Estimated log $p(x_1)$					
VAE	-86.313	-91.126	-232.758	-152.835	-6237.120
BiVCCA	-87.354	-92.089	-233.634	-202.490	-7263.536
JMVAE	-86.305	-90.697	-232.630	-152.787	-6237.967
MVAE-Q	-91.665	-96.028	-236.081	-166.580	-6290.085
MVAE	<b>-86.026</b>	<b>-90.619</b>	<b>-232.535</b>	<b>-152.761</b>	-6236.923
MVAE19	-	-	-	-	<b>-6236.109</b>
Estimated log $p(x_1, x_2)$					
JMVAE	-86.371	<b>-90.769</b>	<b>-232.948</b>	<b>-153.101</b>	-6242.187
MVAE-Q	-92.259	-96.641	-236.827	-173.615	-6294.861
MVAE	<b>-86.255</b>	-90.859	-233.007	-153.469	-6242.034
MVAE19	-	-	-	-	<b>-6239.944</b>
Estimated log $p(x_1 x_2)$					
CVAE	<b>-83.448</b>	<b>-87.773</b>	<b>-229.667</b>	-	-6228.771
JMVAE	-83.985	-88.696	-230.396	<b>-145.977</b>	<b>-6231.468</b>
MVAE-Q	-90.024	-94.347	-234.514	-163.302	-6311.487
MVAE	-83.970	-88.569	-230.695	-147.027	-6234.955
MVAE19	-	-	-	-	-6233.340

Table 2: Estimates (using  $q(z|x_1)$ ) for marginal probabilities on the average test example. MVAE and JMVAE are roughly equivalent in data log-likelihood but as Table 1 shows, MVAE uses far fewer parameters. The CVAE is often better at capturing  $p(x_1|x_2)$  but does not learn a joint distribution.

# IWAW

## IMPORTANCE WEIGHTED AUTOENCODERS

LOSS of VAE:  $E_{z \sim q(z|x)} [\log \frac{p(x,z)}{q(z|x)}]$

LOSS of IWAE:  $E_{z_1, z_2, \dots, z_k \sim q(z|x)} [\log \frac{1}{k} \sum_{i=1}^k \frac{p(x,z_i)}{q(z_i|x)}]$



Why weighted?

Gradient:

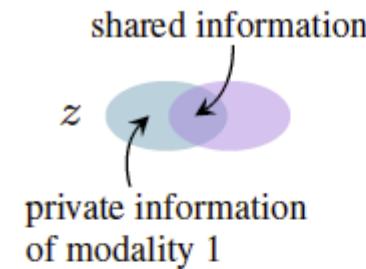
**VAE:**  $\frac{1}{k} \sum_{i=1}^k \nabla_{\theta} \log w_i$        $w: \frac{p(x,z|\theta)}{q(z|x,\theta)}$

**IWAE:**  $\sum_{i=1}^k \tilde{w}_i \nabla_{\theta} \log w_i$        $\frac{w_i}{\sum_{j=1}^k w_j} = \tilde{w}_i$

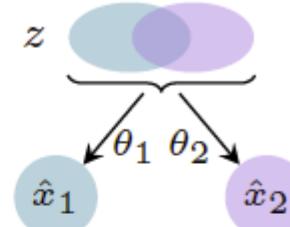
# MMVAE

## 4 Standard:

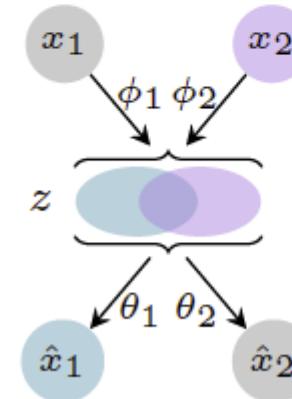
- Latent Factorisation
- Coherent Joint Generation
- Coherent Cross Generation
- Synergy



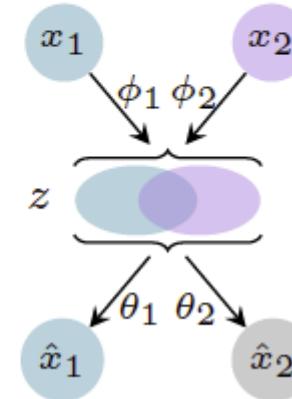
(a) Latent Factorisation



(b) Joint Generation



(c) Cross Generation



(d) Synergy

# MMVAE

ELBO of VAE:  $\mathcal{L}_{\text{ELBO}}(\mathbf{x}_{1:M}) = \mathbb{E}_{\mathbf{z} \sim q_{\Phi}(\mathbf{z} | \mathbf{x}_{1:M})} \left[ \log \frac{p_{\Theta}(\mathbf{z}, \mathbf{x}_{1:M})}{q_{\Phi}(\mathbf{z} | \mathbf{x}_{1:M})} \right]$

ELBO of IWAE:  $\mathcal{L}_{\text{IWAE}}(\mathbf{x}_{1:M}) = \mathbb{E}_{\mathbf{z}^{1:K} \sim q_{\Phi}(\mathbf{z} | \mathbf{x}_{1:M})} \left[ \log \sum_{k=1}^K \frac{1}{K} \frac{p_{\Theta}(\mathbf{z}^k, \mathbf{x}_{1:M})}{q_{\Phi}(\mathbf{z}^k | \mathbf{x}_{1:M})} \right]$  Why MoE?

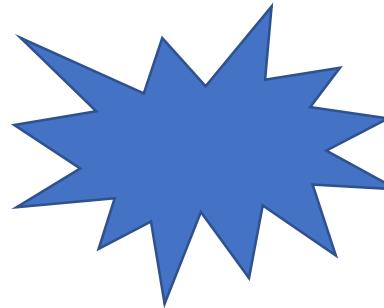
MoE(mixture of experts)

$$q_{\Phi}(\mathbf{z} | \mathbf{x}_{1:M}) = \sum_m \alpha_m \cdot q_{\phi_m}(\mathbf{z} | \mathbf{x}_m)$$

ELBO of MMVAE:  $\mathcal{L}_{\text{IWAE}}^{\text{MoE}}(\mathbf{x}_{1:M}) = \frac{1}{M} \sum_{m=1}^M \mathbb{E}_{\mathbf{z}_m^{1:K} \sim q_{\phi_m}(\mathbf{z} | \mathbf{x}_m)} \left[ \log \frac{1}{K} \sum_{k=1}^K \frac{p_{\Theta}(\mathbf{z}_m^k, \mathbf{x}_{1:M})}{q_{\Phi}(\mathbf{z}_m^k | \mathbf{x}_{1:M})} \right],$

# MMVAE

$$\begin{aligned}\mathcal{L}_{\text{ELBO}}(\mathbf{x}_{1:M}) &= \mathbb{E}_{q_{\Phi}(\mathbf{z} | \mathbf{x}_{1:M})} \left[ \log \frac{p_{\Theta}(\mathbf{z}, \mathbf{x}_{1:M})}{q_{\Phi}(\mathbf{z} | \mathbf{x}_{1:M})} \right] = \frac{1}{M} \sum_{m=1}^M \mathbb{E}_{\mathbf{z}_m \sim q_{\phi_m}(\mathbf{z} | \mathbf{x}_m)} \left[ \log \frac{p_{\Theta}(\mathbf{z}_m, \mathbf{x}_{1:M})}{q_{\Phi}(\mathbf{z}_m | \mathbf{x}_{1:M})} \right] \\ &\leq \frac{1}{M} \sum_{m=1}^M \mathbb{E}_{\mathbf{z}_m^{1:K} \sim q_{\phi_m}(\mathbf{z} | \mathbf{x}_m)} \left[ \log \frac{1}{K} \sum_{k=1}^K \frac{p_{\Theta}(\mathbf{z}_m^k, \mathbf{x}_{1:M})}{q_{\Phi}(\mathbf{z}_m^k | \mathbf{x}_{1:M})} \right] = \mathcal{L}_{\text{IWAE}}^{\text{MoE}}(\mathbf{x}_{1:M}).\end{aligned}$$



Multimodal IWAE!!!!

# MMVAE

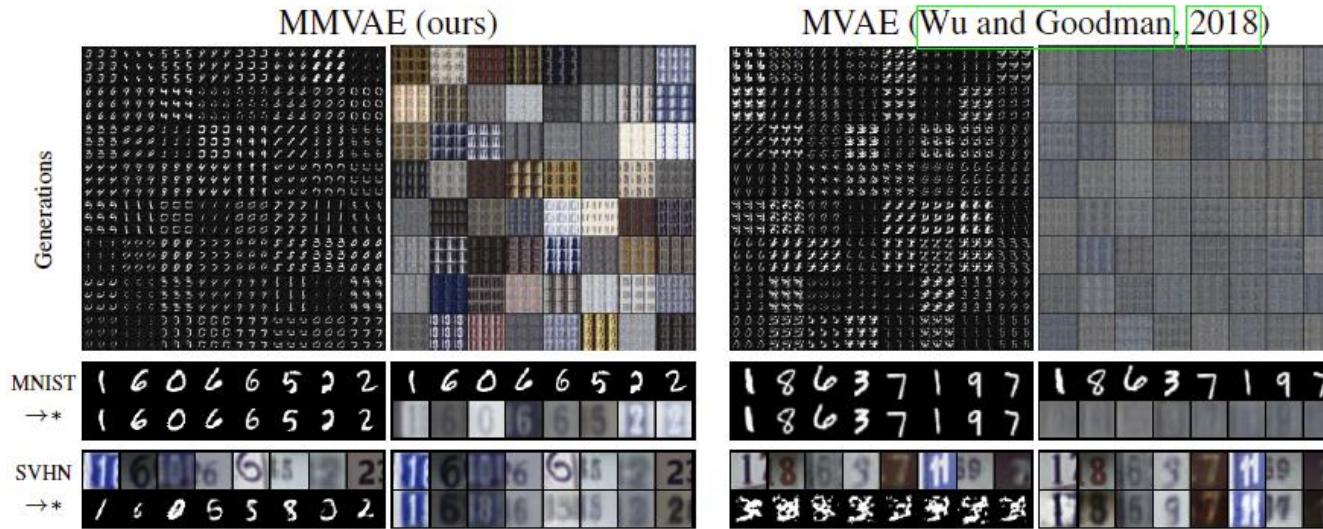


Table 1: Digit classification accuracy (%) of latent variables in different models.

	MMVAE	MVAE (single)	MVAE (both)	single-VAE
MNIST	91.3	95.7	94.9	85.3
SVHN	68.0	9.1	90.1	20.7

Table 2: Probability of digit matching (%) for joint and cross generation

	Joint	Cross (M→S)	Cross (S→M)
MMVAE	42.1	86.4	69.1
MVAE	12.7	—	9.5

Table 3: Evaluating the different log likelihoods for different arrangements of MNIST and SVHN.

		$\log p(\mathbf{x}_m, \mathbf{x}_n)$	$\log p(\mathbf{x}_m   \mathbf{x}_m, \mathbf{x}_n)$	$\log p(\mathbf{x}_m   \mathbf{x}_m)$	$\log p(\mathbf{x}_m   \mathbf{x}_n)$
$m = \text{MNIST}$ , $n = \text{SVHN}$	MMVAE	6261.40	868.76	868.37	628.31
$m = \text{SVHN}$ , $n = \text{MNIST}$	MVAE	2961.80	-176.68	-107.46	-778.20
$m = \text{SVHN}$ , $n = \text{MNIST}$	MMVAE	6261.40	3441.01	3441.01	2337.56
$m = \text{SVHN}$ , $n = \text{MNIST}$	MVAE	2961.80	3395.12	3536.86	-12747.50