

# Multivariate-Information Adversarial Ensemble for Scalable Joint Distribution Matching

Ziliang Chen\*, Zhanfu Yang\*, Xiaoxi Wang\*, Xiaodan Liang, Xiaopeng Yan,  
Guanbin Li, Liang Lin

Sun Yat-sen University, Purdue University



# Motivation

---

Implicit Generative Models (IGM), e.g., GAN [1] and CycleGAN [2], boil down to an  $m$ -domain joint distribution matching (JDM):

$$\begin{aligned} p(\mathbf{x}_1, \dots, \mathbf{x}_m) &:= p(\{\mathbf{x}_j\}_{j \in [m] \& j \neq i} | \mathbf{x}_i) p(\mathbf{x}_i) \\ &= p_{\Theta}(\{\mathbf{x}_j\}_{j \in [m] \& j \neq i} | \mathbf{x}_i) p(\mathbf{x}_i) \end{aligned}$$

[1] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, Generative Adversarial Nets. Proceedings Neural Information Processing Systems Conference, 2014

[2] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks. ArXiv, 2017

# Motivation

---

Implicit Generative Models (IGM), e.g., GAN [1] and CycleGAN [2], boil down to an  $m$ -domain joint distribution matching (JDM):

$$\begin{aligned} p(\mathbf{x}_1, \dots, \mathbf{x}_m) &:= p(\{\mathbf{x}_j\}_{j \in [m] \& j \neq i} | \mathbf{x}_i) p(\mathbf{x}_i) \\ &= p_{\Theta}(\{\mathbf{x}_j\}_{j \in [m] \& j \neq i} | \mathbf{x}_i) p(\mathbf{x}_i) \end{aligned}$$

However, if  $m > 2$ ,

[1] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, Generative Adversarial Nets. Proceedings Neural Information Processing Systems Conference, 2014

[2] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks. ArXiv, 2017

# Motivation

---

Implicit Generative Models (IGM), e.g., GAN [1] and CycleGAN [2], boil down to an  $m$ -domain joint distribution matching (JDM):

$$\begin{aligned} p(\mathbf{x}_1, \dots, \mathbf{x}_m) &:= p(\{\mathbf{x}_j\}_{j \in [m] \& j \neq i} | \mathbf{x}_i) p(\mathbf{x}_i) \\ &= p_{\Theta}(\{\mathbf{x}_j\}_{j \in [m] \& j \neq i} | \mathbf{x}_i) p(\mathbf{x}_i) \end{aligned}$$

However, if  $m > 2$ ,

Solution 1: CycleGAN, JointGAN[3]

They suffer combinatorial explosion in their parameters

Solution 2: StarGAN [4] and its variants

The domain-shared model lacks theoretical support, fragile in model collapse

[1] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, Generative Adversarial Nets. Proceedings Neural Information Processing Systems Conference, 2014

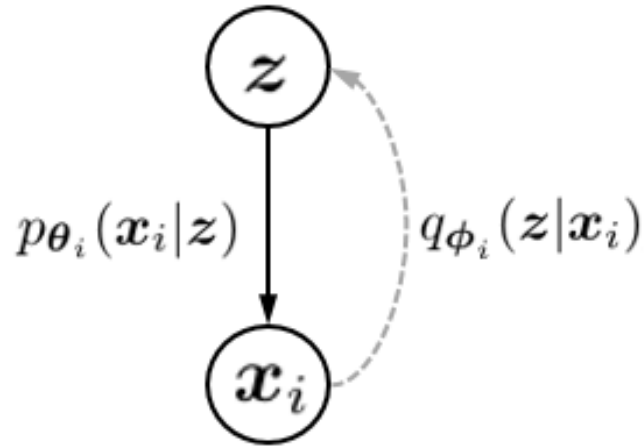
[2] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks. ArXiv, 2017

[3] Pu Y, Dai S, Gan Z, et al. JointGAN: Multi-Domain Joint Distribution Learning with Generative Adversarial Nets ICML. 2018.

[4] Choi Y, Choi M, Kim M, et al. StarGAN: Unified Generative Adversarial Networks for Multi-Domain Image-to-Image Translation. CVPR 2018.

# Solution 3: ALI Ensemble

---

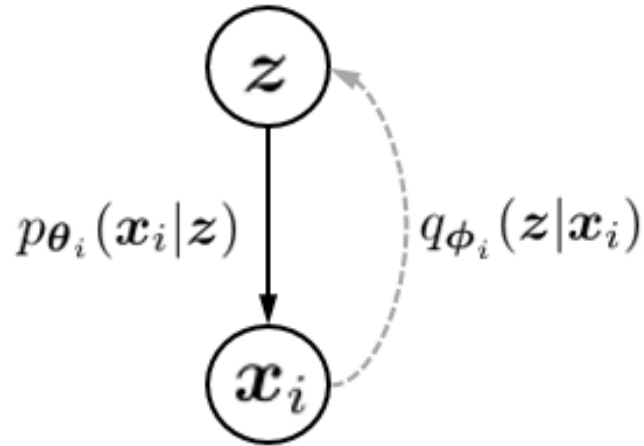


**Adversarially Learned Inference (ALI)**  
Model [5]

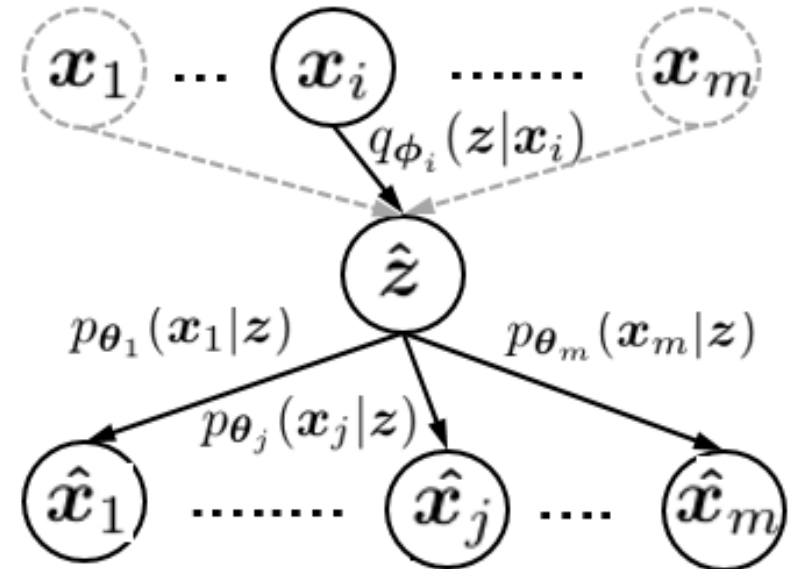
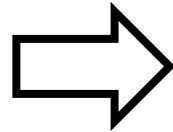
$$\min_{\theta_i, \phi_i} \max_{\omega_i} \mathcal{L}_{\text{ALI}}^{(i)}(\theta_i, \phi_i, \omega_i) =$$
$$\mathbb{E}_{\mathbf{x}_i \sim p(\mathbf{x}_i), \hat{\mathbf{z}} \sim q_{\phi_i}(\hat{\mathbf{z}}|\mathbf{x}_i)} [\log f_{\omega_i}(\mathbf{x}_i, \hat{\mathbf{z}})]$$
$$+ \mathbb{E}_{\hat{\mathbf{x}}_i \sim p_{\theta_i}(\hat{\mathbf{x}}_i|\mathbf{z}), \mathbf{z} \sim q(\mathbf{z})} [\log (1 - f_{\omega_i}(\hat{\mathbf{x}}_i, \mathbf{z}))]$$

**Lemma 1** ((Dumoulin et al., 2016)). *The optimal generation, inference and critic nets w.r.t.,  $\{\theta_i^*, \phi_i^*, \omega_i^*\}$  ( $\forall i \in [m]$ ) refer to a saddle point in Eq.2  $\iff p_{\theta_i^*}(\mathbf{x}_i|\mathbf{z})q(\mathbf{z}) = q_{\phi_i^*}(\mathbf{z}|\mathbf{x}_i)p_i(\mathbf{x}_i)$ .*

# Solution 3: ALI Ensemble



ALI ensemble  
across  $m$  domains

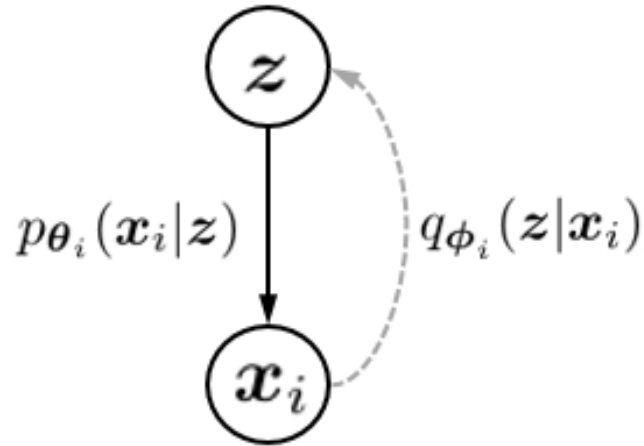


Advantages: (1). Linear-parameter scalability as  $m$  increases.

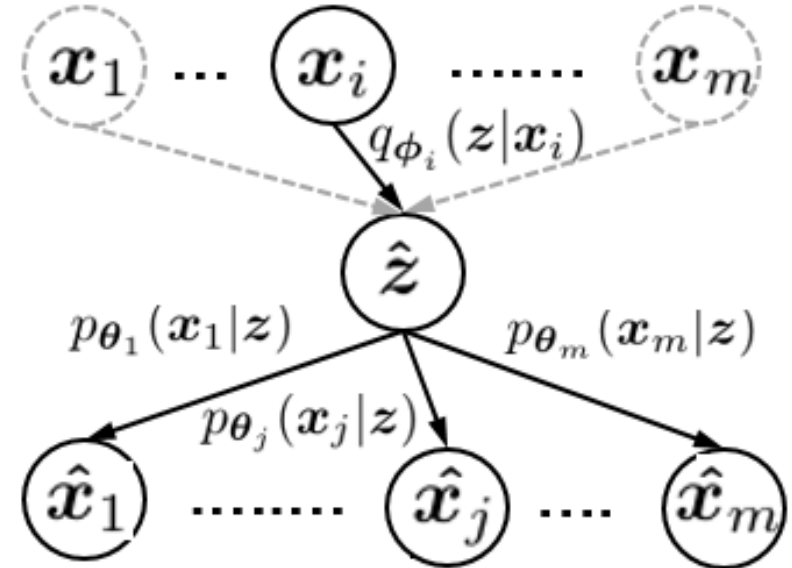
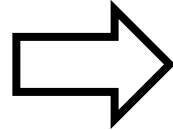
(2). Generative model capability: **Proposition 1.** *Given a pair of domains  $\forall i, j \in [m], i \neq j$ , their well-trained ALIs (in Lemma.1) construct a cross-domain transfer process  $p_{\Phi, \Theta}(\hat{x}_j | x_i)$  that satisfies*

$$p_{\Phi^*, \Theta^*}(\hat{x}_j) = \int p_{\Phi^*, \Theta^*}(\hat{x}_j | x_i) p_i(x_i) dx_i = p_j(\hat{x}_j)$$

# Solution 3: ALI Ensemble



ALI ensemble  
across  $m$  domains



Advantages: (1). Linear-parameter scalability as  $m$  increases.

(2). Generative model capability:

**Marginal matching do not imply joint density matching !**

**Proposition 1.** Given a pair of domains  $\forall i, j \in [m], i \neq j$ , their well-trained ALIs (in Lemma.1) construct a cross-domain transfer process  $p_{\Phi, \Theta}(\hat{x}_j | x_i)$  that satisfies

$$p_{\Phi^*, \Theta^*}(\hat{x}_j) = \int p_{\Phi^*, \Theta^*}(\hat{x}_j | x_i) p_i(x_i) dx_i = p_j(\hat{x}_j)$$

# JDM Criteria

---

In supervised learning, data are drawn from  $p(\mathbf{x}_1, \dots, \mathbf{x}_m)$ , each of them presenting as  $m$ -tuple. So the criterion can be written as

$$\min_{\Phi, \Theta} - \mathbb{E}_p \left[ \log p_{\Phi, \Theta}(\{\mathbf{x}_i\}_{i=1}^m) \right]$$

In unsupervised learning, no access is provided to draw  $m$ -tuple from  $p(\mathbf{x}_1, \dots, \mathbf{x}_m)$ . Extending the observation from [6], the criterion is to minimize the conditional entropy

$$\begin{aligned} \min_{\Phi, \Theta} H(\mathbf{x}_i | \{\hat{\mathbf{x}}_j\}_{j \in [m] \& j \neq i}) \\ = -\mathbb{E}_{p_{\Phi, \Theta}} \left[ \log p_{\Phi, \Theta}(\mathbf{x}_i | \{\hat{\mathbf{x}}_j\}_{j \in [m] \& j \neq i}) \right] \end{aligned}$$



# Multivariate Mutual Information

---

Mutual Information:  $I(\mathbf{x}; \mathbf{y}) = I(\mathbf{y}; \mathbf{x}) := H(\mathbf{y}) - H(\mathbf{y}|\mathbf{x})$

Multivariate Mutual Information:  $I(\mathbf{y}_1; \dots; \mathbf{y}_n)$   
 $:= I(\mathbf{y}_1; \dots; \mathbf{y}_{n-1}) - I(\mathbf{y}_1; \dots; \mathbf{y}_{n-1} | \mathbf{y}_n)$

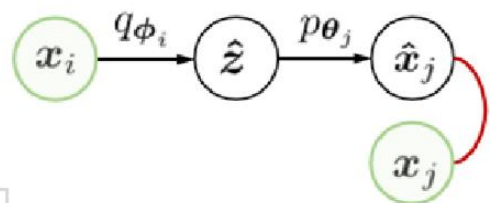
Our method aims to achieve:  $\min_{\Phi, \Theta} - \sum_{i, j \in [m], i \neq j} I_{\Phi, \Theta}(\mathbf{x}_i; \mathbf{x}_j; \mathbf{z})$

# Multivariate Mutual Information Leads To JDM

**Observation 1.** Given empirical draws from  $p_i$  ( $\forall i \in [m]$ ), in supervised learning,

$$\begin{aligned}
 & -I_{\Phi, \Theta}(\mathbf{x}_i; \mathbf{x}_j; \hat{\mathbf{z}}) \leq H_{\Phi, \Theta}(\mathbf{x}_i | \mathbf{x}_j) \\
 & \leq \mathbb{E}_{\mathbf{x}_i, \mathbf{x}_j \sim p_{i,j}} - \left[ \log \int p_{\theta_i}(\mathbf{x}_i | \hat{\mathbf{z}}) q_{\phi_j}(\hat{\mathbf{z}} | \mathbf{x}_j) d\hat{\mathbf{z}} \right] \triangleq \mathcal{L}_{\Phi, \Theta}^{\text{con}}(\mathbf{x}_i, \mathbf{x}_j)
 \end{aligned} \tag{10}$$

where  $p_{i,j} = p(\mathbf{x}_i, \mathbf{x}_j)$ .



$$\mathcal{L}_{\Phi, \Theta}^{\text{con}}(\mathbf{x}_i, \mathbf{x}_j)$$

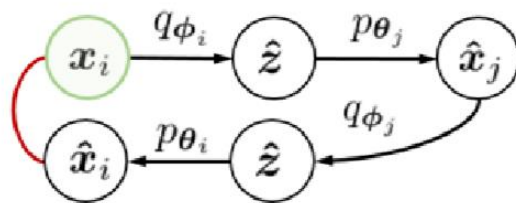
(a).Observation.1

**Observation 2.** Given empirical draws from  $p_i$  ( $\forall i \in [m]$ ), in unsupervised learning,

$$\begin{aligned}
 & -I_{\Phi, \Theta}(\mathbf{x}_i; \hat{\mathbf{x}}_j; \hat{\mathbf{z}}) \leq H_{\Phi, \Theta}(\mathbf{x}_i | \hat{\mathbf{x}}_j) \\
 & \leq \mathbb{E}_{\mathbf{x}_i, \hat{\mathbf{x}}_j \sim p_{\theta_j, \phi_i}} - \left[ \log \int p_{\theta_i}(\mathbf{x}_i | \hat{\mathbf{z}}) q_{\phi_j}(\hat{\mathbf{z}} | \hat{\mathbf{x}}_j) d\hat{\mathbf{z}} \right] \triangleq \mathcal{L}_{\Phi, \Theta}^{\text{cycle}}(\mathbf{x}_i, \hat{\mathbf{x}}_j)
 \end{aligned} \tag{11}$$

where  $p_{\theta_j, \phi_i} = p(\mathbf{x}_i) \int_{\hat{\mathbf{z}}} p_{\theta_j}(\hat{\mathbf{x}}_j | \hat{\mathbf{z}}) q_{\phi_i}(\hat{\mathbf{z}} | \mathbf{x}_i) d\hat{\mathbf{z}}$ .

$\mathcal{L}_{\Phi, \Theta}^{\text{cycle}}(\mathbf{x}_i, \hat{\mathbf{x}}_j)$  is constructed as illustrated in Fig.2.b.



$$\mathcal{L}_{\Phi, \Theta}^{\text{cycle}}(\mathbf{x}_i, \hat{\mathbf{x}}_j)$$

(b).Observation.2

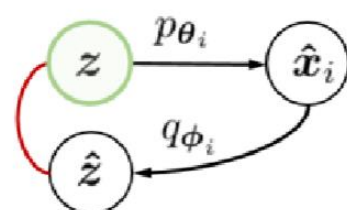
**Observation 3.** Given empirical draws from  $q(z)$ ,

$$-I_{\Phi, \Theta}(\hat{\mathbf{x}}_i; \hat{\mathbf{x}}_j; z) \leq H_{\Phi, \Theta}(z | \hat{\mathbf{x}}_i) + H_{\Phi, \Theta}(\hat{\mathbf{x}}_j | \hat{\mathbf{x}}_i) \tag{12}$$

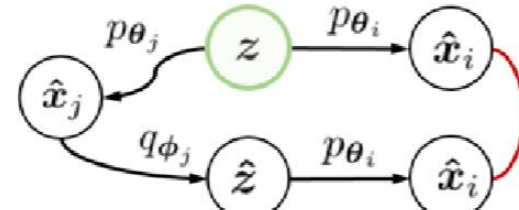
$$H_{\Phi, \Theta}(z | \hat{\mathbf{x}}_i) = \mathbb{E}_{\hat{\mathbf{x}}_i \sim p_{\theta_i}, z \sim q(z)} - \log q_{\phi_i}(z | \hat{\mathbf{x}}_i) \triangleq \mathcal{L}_{\Phi, \Theta}^{\text{cycle}}(z, \hat{\mathbf{x}}_i)$$

$$H_{\Phi, \Theta}(\hat{\mathbf{x}}_j | \hat{\mathbf{x}}_i) = \mathbb{E}_{\hat{\mathbf{x}}_i \sim p_{\theta_i}, \hat{\mathbf{x}}_j \sim p_{\theta_j}} - \left[ \log \int_z p_{\theta_j}(\hat{\mathbf{x}}_j | z) q_{\phi_i}(z | \hat{\mathbf{x}}_i) dz \right]$$

$$\triangleq \mathcal{L}_{\Phi, \Theta}^{\text{cycle}}(z, \hat{\mathbf{x}}_i, \hat{\mathbf{x}}_j)$$

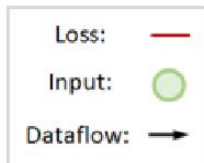


$$\mathcal{L}_{\Phi, \Theta}^{\text{cycle}}(z, \hat{\mathbf{x}}_i)$$



$$\mathcal{L}_{\Phi, \Theta}^{\text{cycle}}(z, \hat{\mathbf{x}}_i, \hat{\mathbf{x}}_j)$$

(c).Observation.3



# Multivariate Mutual Information Leads To JDM

---

$$\begin{aligned}\mathcal{R}_{\text{SL}}(\Theta, \Phi) &= \sum_{i,j \in [m], i \neq j} \mathcal{L}_{\Phi, \Theta}^{\text{con}}(\mathbf{x}_i, \mathbf{x}_j) + \mathcal{L}_{\Phi, \Theta}^{\text{cycle}}(z, \hat{\mathbf{x}}_i) \\ &\quad + \mathcal{L}_{\Phi, \Theta}^{\text{cycle}}(z, \hat{\mathbf{x}}_i, \hat{\mathbf{x}}_j) \\ \mathcal{R}_{\text{UL}}(\Theta, \Phi) &= \sum_{i,j \in [m], i \neq j} \mathcal{L}_{\Phi, \Theta}^{\text{cycle}}(\mathbf{x}_i, \hat{\mathbf{x}}_j) + \mathcal{L}_{\Phi, \Theta}^{\text{cycle}}(z, \hat{\mathbf{x}}_i) \\ &\quad + \mathcal{L}_{\Phi, \Theta}^{\text{cycle}}(z, \hat{\mathbf{x}}_i, \hat{\mathbf{x}}_j)\end{aligned}$$

**Theorem 1.** *Suppose that true and parameterized domain marginal distributions maintain a high likelihood to domain variables,  $\mathcal{R}_{\text{SL}} \rightarrow 0$  leads to the optima in  $\min_{\Phi, \Theta} -\mathbb{E}_p[\log p_{\Phi, \Theta}(\{\mathbf{x}_i\}_{i=1}^m)]$*

*$\mathcal{R}_{\text{UL}} \rightarrow 0$  leads to the optima in  $\min_{\Phi, \Theta} H(\mathbf{x}_i | \{\hat{\mathbf{x}}_j\}_{j \in [m] \& j \neq i})$*

# Adversarial Ensemble Learning

---

$$\begin{aligned}
 \min_{\Theta, \Phi} \max_{\Omega} (1 - \gamma) \sum_{i=1}^m \mathcal{L}_{\text{ALI}}^{(i)} + \gamma \sum_{i=1}^m \mathcal{L}_{\text{DMAE}}^{(i)} + \beta \mathcal{R}_{\text{SL}} / \mathcal{R}_{\text{UL}} \\
 \text{s.t. } \mathcal{L}_{\text{DMAE}}^{(i)}(\Phi, \Theta, \Omega) = \mathbb{E}_{\mathbf{x}_i, \hat{\mathbf{z}} \sim q_{\phi_i}(\mathbf{x}_i, \hat{\mathbf{z}})} [\log f_{\omega_i}(\mathbf{x}_i, \hat{\mathbf{z}})] \\
 + \sum_{j=1}^m \pi_j \left( \mathbb{E}_{\hat{\mathbf{x}}_i \sim p_{\theta_i}(\hat{\mathbf{x}}_i | \mathbf{z}), \mathbf{z} \sim q_{\phi_j}} [\log (1 - f_{\omega_i}(\hat{\mathbf{x}}_i, \mathbf{z}))] \right)
 \end{aligned}$$

where  $\mathcal{R}_{\text{SL}} / \mathcal{R}_{\text{UL}}$  are switched by supervised/unsupervised learning and  $\beta > 0$  denotes the loss-balance factor.

**Proposition 2.** *The optimum of the generation, inference and critic networks in*

$$\min_{\Theta, \Phi} \max_{\Omega} (1 - \gamma) \sum_{i=1}^m \mathcal{L}_{\text{ALI}}^{(i)} + \gamma \sum_{i=1}^m \mathcal{L}_{\text{MALI}}^{(i)}$$

*refer to their saddle points in Lemma.1 if and only if  $\forall i \in [m]$ , there exist  $p_{\theta_i^*}(\mathbf{x} | \mathbf{z}) q(\mathbf{z}) = q_{\phi_i^*}(\mathbf{z} | \mathbf{x}) p(\mathbf{x})$ .*

# Experiments

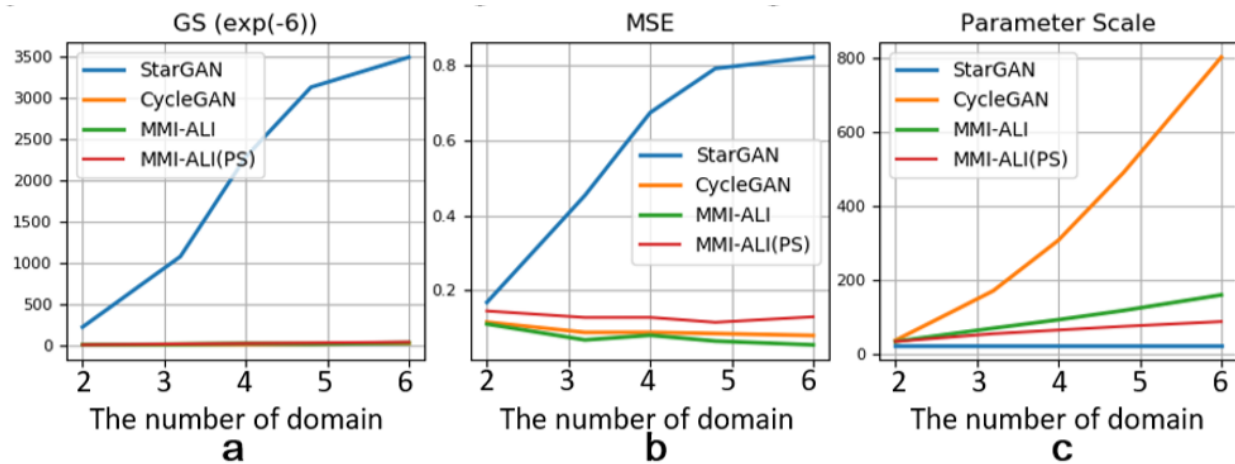
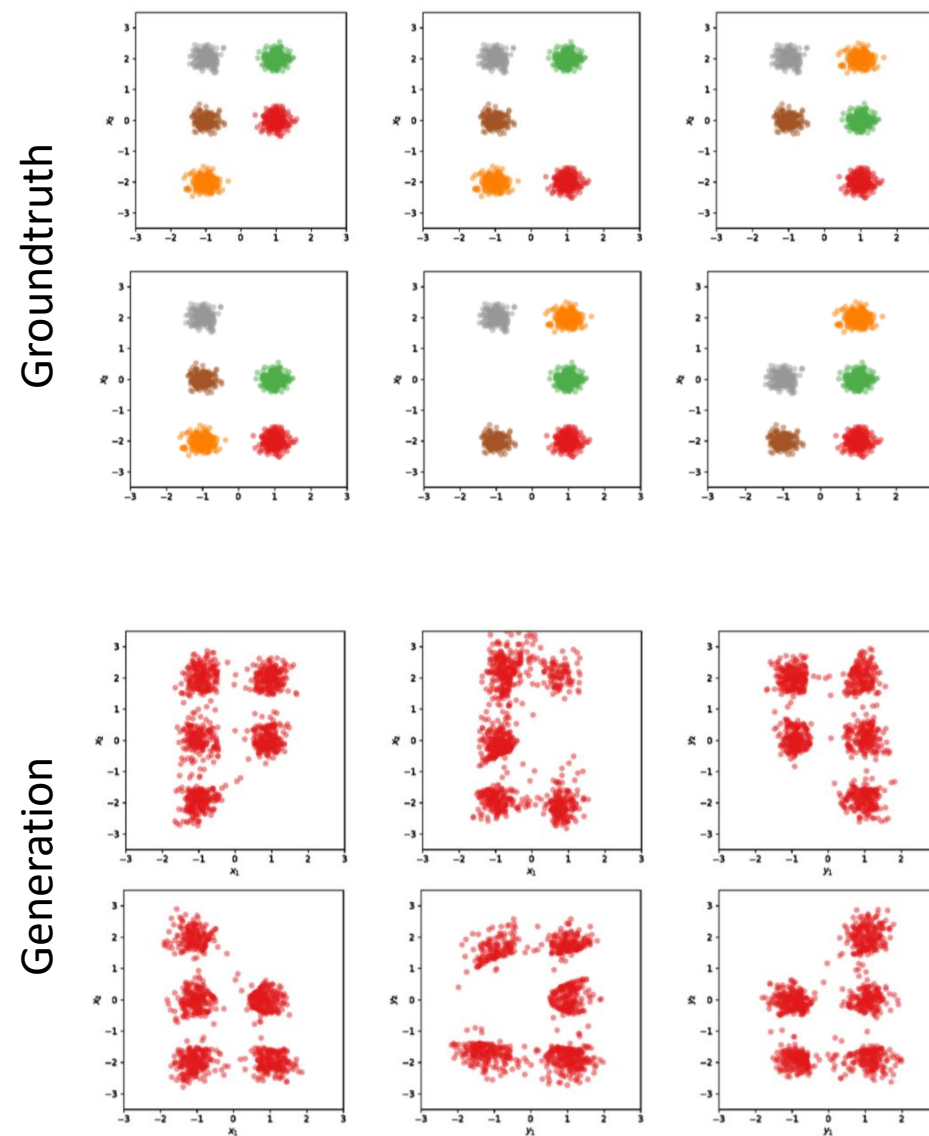
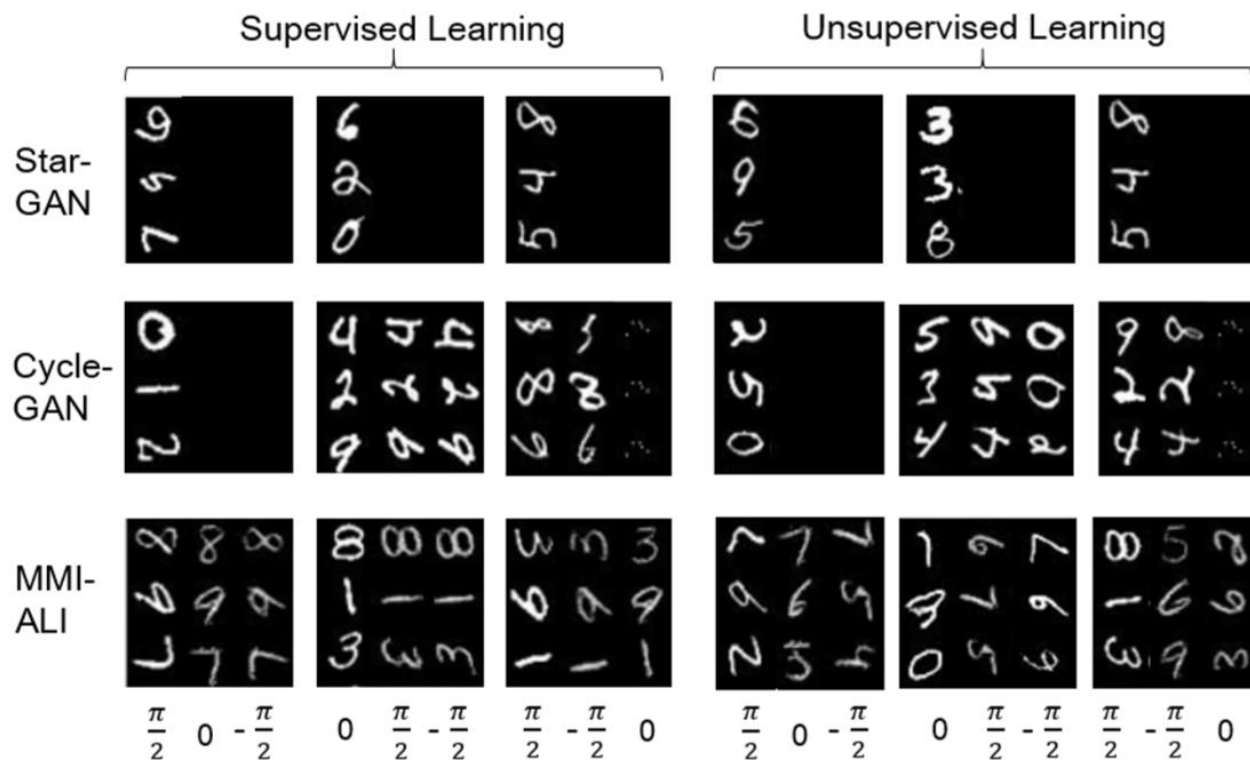


Figure 4. Transfer evaluations with 2~6 synthetic domains: (a). Geometric Score (GS, lower is better); (b). Mean Square Error (MSE, lower is better); (c). Parameter Scale (lower is better).

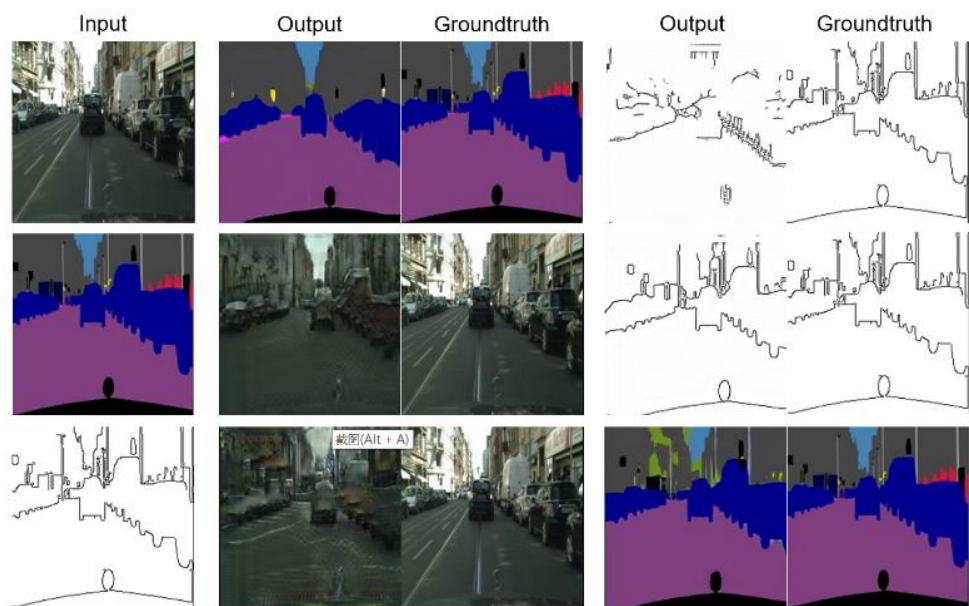


# Experiments

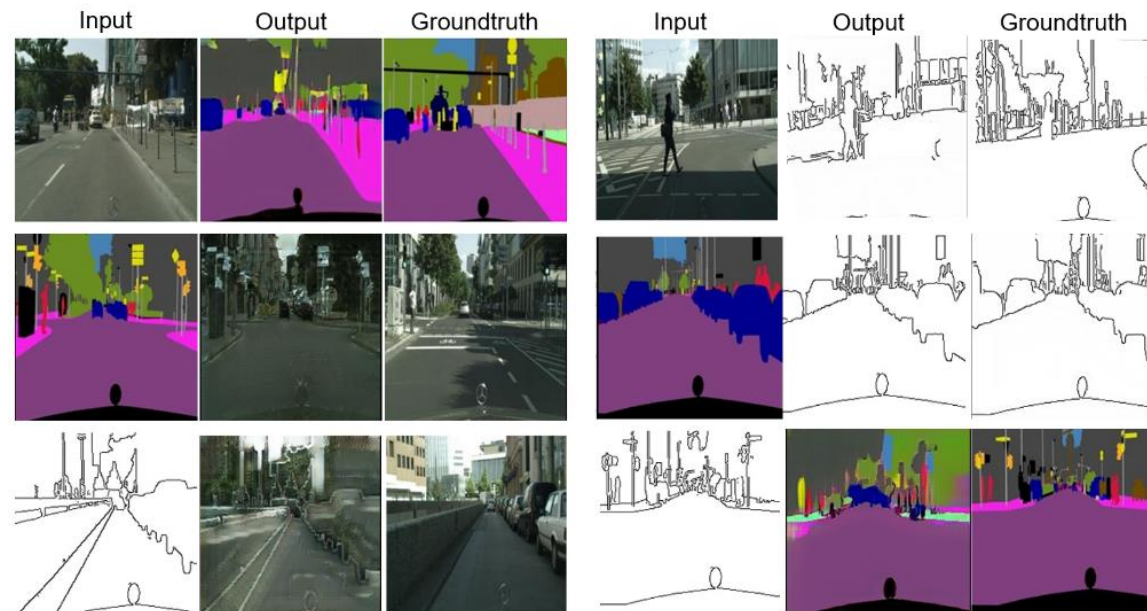


	$-\frac{\pi}{2} \rightarrow 0$	$\frac{\pi}{2} \rightarrow 0$	$0 \rightarrow \frac{\pi}{2}$	$-\frac{\pi}{2} \rightarrow \frac{\pi}{2}$	$-\frac{\pi}{2} \rightarrow 0$	$\frac{\pi}{2} \rightarrow -\frac{\pi}{2}$
StarGAN	$1.00 \pm 0.00$	$1.00 \pm 0.00$	$1.00 \pm 0.00$	$1.00 \pm 0.00$	$1.00 \pm 0.00$	$1.00 \pm 0.00$
CycleGAN	$8.34 \pm 0.12$	$6.13 \pm 0.37$	$2.25 \pm 0.02$	$2.38 \pm 0.06$	$1.71 \pm 0.03$	$1.04 \pm 0.01$
MMI-ALI (wo MMI)	$7.48 \pm 0.05$	$6.06 \pm 0.10$	$3.19 \pm 0.04$	$2.90 \pm 0.05$	$2.73 \pm 0.09$	$2.47 \pm 0.12$
MMI-ALI ( $\gamma = 0$ )	$8.34 \pm 0.20$	$8.27 \pm 0.10$	<b><math>3.26 \pm 0.06</math></b>	$3.06 \pm 0.10$	$3.15 \pm 0.11$	$2.92 \pm 0.12$
MMI-ALI	<b><math>8.99 \pm 0.06</math></b>	<b><math>9.01 \pm 0.00</math></b>	$2.95 \pm 0.08$	<b><math>3.86 \pm 0.12</math></b>	<b><math>3.31 \pm 0.12</math></b>	<b><math>3.08 \pm 0.05</math></b>

# Experiments



Supervised Learning



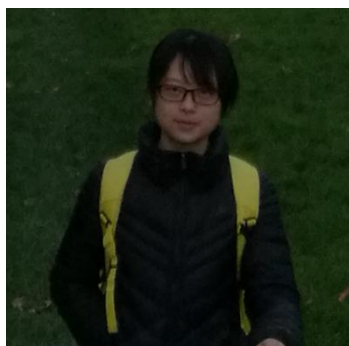
Unsupervised Learning

		R→Seg	Seg→R	R→Ske	Ske→R	Seg→Ske	Ske→Seg
Unsuper	ST	405.16	372.59	385.08	388.97	357.19	417.39
	CG	224.04	<b>213.43</b>	164.65	<b>222.24</b>	<b>60.20</b>	<b>144.07</b>
	Ours	<b>202.93</b>	254.41	<b>150.98</b>	246.04	101.30	192.13
Super	ST	382.90	440.53	419.11	383.72	400.70	299.82
	CG	<b>217.28</b>	260.41	<b>171.04</b>	<b>223.43</b>	65.18	228.61
	Ours	250.48	<b>246.01</b>	196.06	229.45	<b>55.76</b>	<b>143.20</b>

# Collaborators



Zhanfu Yang



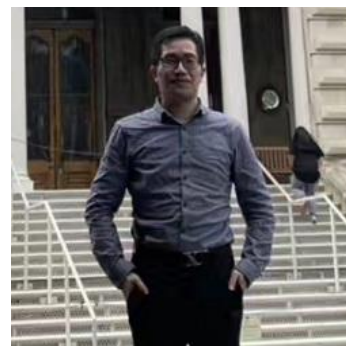
Xiaoxi Wang



Xiaodan Liang



Xiaopeng Yan



Guanbin Li



Liang Lin

Thank You!