

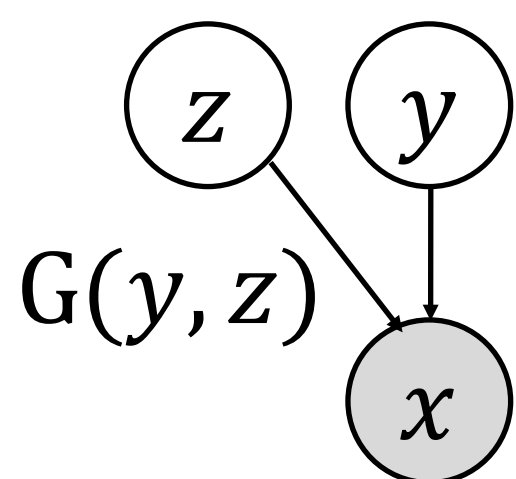
## Problem

Semi-supervised conditional generative modeling

- Conditional generative models are quite useful
  - Generate data samples with designated semantics
  - Synthetic data help supervised training of downstream tasks
- Challenges: labels are scarce
  - How to accurately capture the conditions during generating process?
  - How to separate semantics of interest from other factors of variations?
- Problem: ensure *controllability* and *disentanglability*
  - *Controllability*: the ability to conditionally generate data strictly following the designated semantics
  - *Disentanglability*: the ability to disentangle the modeled semantic of interest from other factors

## Intuition

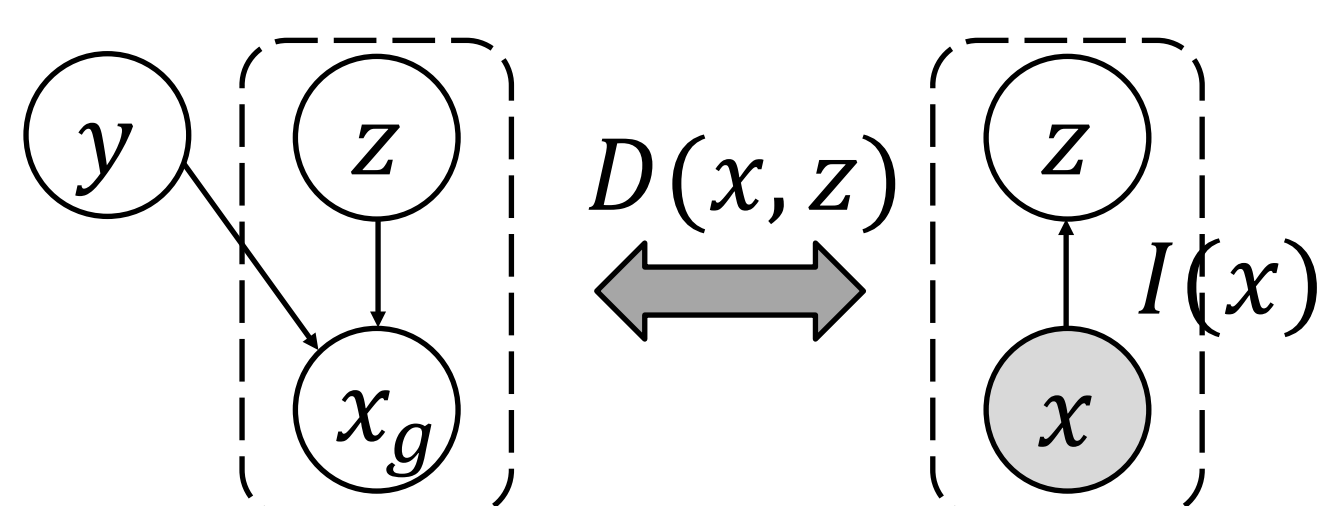
- Hidden space shall be structured as
  - Semantic of our interest  $y$
  - Other factors of variations  $z$
- Hence our goal: learn a generated model  $p_g(x|y, z)$  with
  - *Controllability*: semantics of our interest are fully captured by  $y$
  - *Disentanglability*:  $y$  and  $z$  are not cluttered as much as possible
  - However, directly learning  $p(x, y, z)$  is difficult



Characterizing  $p(x, z)$  and  $p(x, y)$  instead of  $p(x, y, z)$

## Model

- Step 1: Learn joint distribution  $p(x, z)$ 
  - Introduce an inference network  $I(x): x \rightarrow z$
  - Estimate  $p(z|x)$  via adversarial learned inference

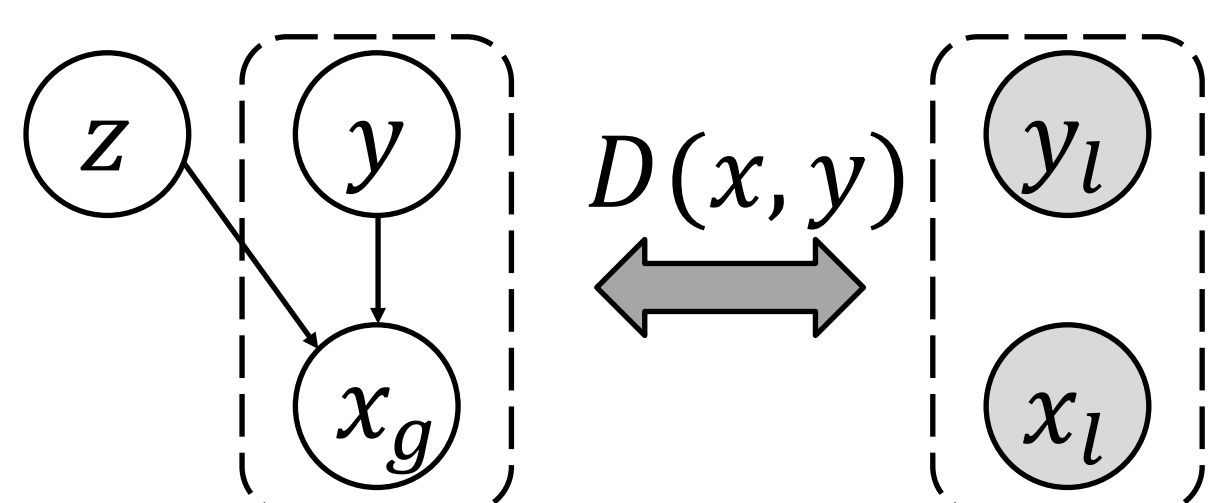


$$p_g(x, z) = p(z)p_g(x|z)$$

$$\updownarrow \text{Match}$$

$$p_i(x, z) = p(x)p_g(z|x)$$

- Step 2: Learn joint distribution  $p(x, y)$ 
  - Estimate  $p(y|x)$  via adversarial learned inference

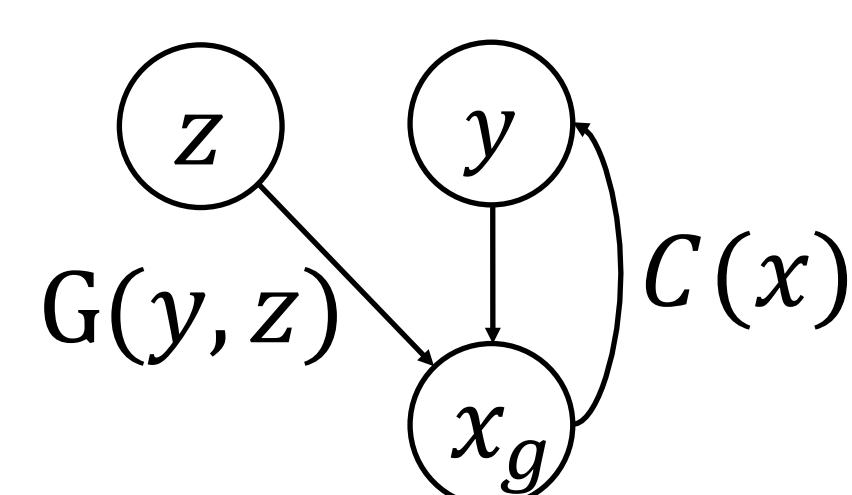


$$p_g(x, y) = p(y)p_g(x|y)$$

$$\updownarrow \text{Match}$$

$$p_i(x, y) = p(x)p_g(y|x)$$

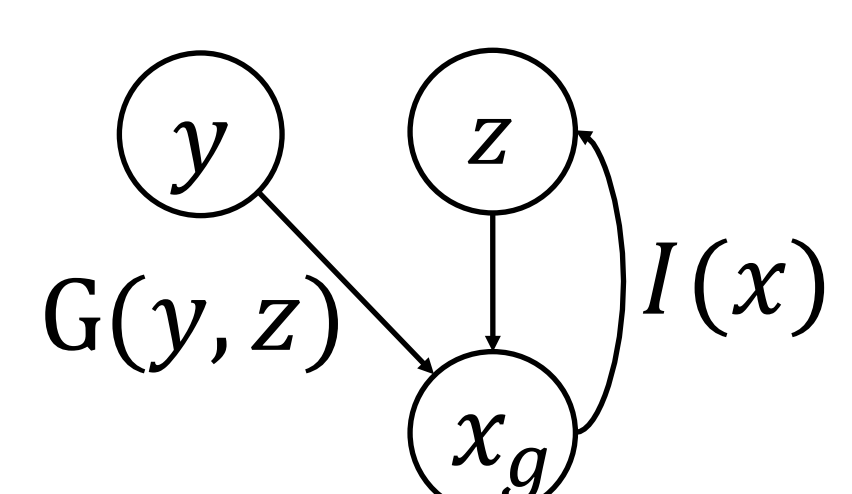
- Step 3: Enforce  $y$  to capture all semantic of interest
  - Therefore, enhance the controllability of the generator
  - Introduce an inference network  $C(x): x \rightarrow y$
  - Minimize reconstruction error  $\mathcal{R}_y$



$$\min_{C, G} \mathcal{R}_y = -\mathbb{E}_{(x, y) \sim p(x, y)} [\log p_c(y|x)]$$

$$-\mathbb{E}_{(x, y) \sim p_g(x, y)} [\log p_c(y|x)]$$

- Step 4: Enforce  $z$  to capture other factors of variations
  - Therefore, enhance the disentanglability of the generator
  - Reuse the inference network  $I(x): x \rightarrow z$
  - Minimize reconstruction error  $\mathcal{R}_z$



$$\min_{I, G} \mathcal{R}_z = -\mathbb{E}_{(x, z) \sim p_g(x, z)} [\log p_i(z|x)]$$

## Training

- Key training techniques
  - Augment labeled dataset with  $(x_c, y_c) \sim p_c(x, y)$
  - Mix  $(x_c, y_c) \sim p_c(x, y)$ ,  $(x_g, y_g) \sim p_g(x, y)$  with labeled data with appropriate mixing proportion

**Algorithm 1** Training Structured Generative Adversarial Networks (SGAN).

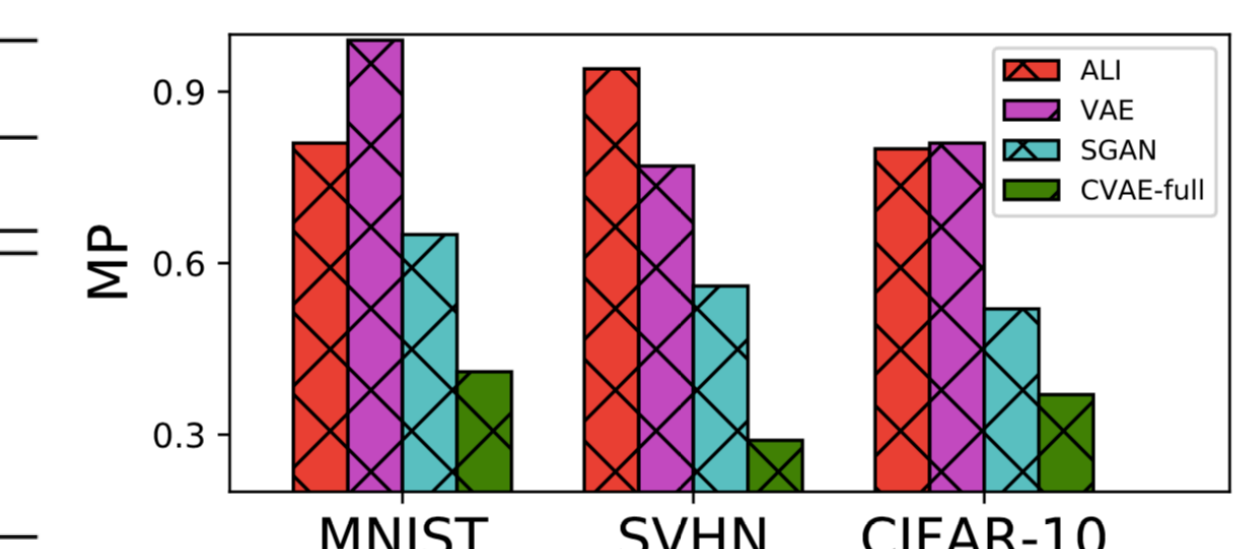
```

1: Pretrain  $C$  by minimizing the first term of Eq. 4 w.r.t.  $C$  using  $X_l$ .
2: repeat
3:   Sample a batch of  $x: x_u \sim p(x)$ .
4:   Sample batches of pairs  $(x, y): (x_l, y_l) \sim p(x, y), (x_g, y_g) \sim p_g(x, y), (x_c, y_c) \sim p_c(x, y)$ .
5:   Obtain a batch  $(x_m, y_m)$  by mixing data from  $(x_l, y_l), (x_g, y_g), (x_c, y_c)$  with proper mixing portion.
6:   for  $k = 1 \rightarrow K$  do
7:     Train  $D_{xz}$  by maximizing the first term of  $\mathcal{L}_{xz}$  using  $x_u$  and the second using  $x_g$ .
8:     Train  $D_{xy}$  by maximizing the first term of  $\mathcal{L}_{xy}$  using  $(x_m, y_m)$  and the second using  $(x_g, y_g)$ .
9:   end for
10:  Train  $I$  by minimizing  $\mathcal{L}_{xz}$  using  $x_u$  and  $\mathcal{R}_z$  using  $x_g$ .
11:  Train  $C$  by minimizing  $\mathcal{R}_y$  using  $(x_m, y_m)$  (see text).
12:  Train  $G$  by minimizing  $\mathcal{L}_{xy} + \mathcal{L}_{xz} + \mathcal{R}_y + \mathcal{R}_z$  using  $(x_g, y_g)$ .
13: until convergence.
  
```

## Results

- Improved controllability and disentanglability
  - Evaluate controllability: generate samples with designated semantics, classify the samples using gold classifiers
  - Evaluate disentanglability: mutual predictability measure (MP)

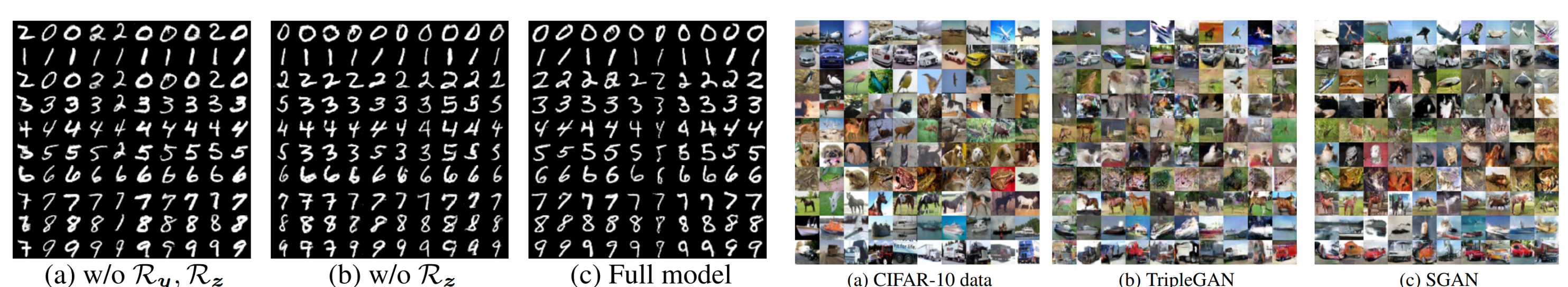
Model	# labeled samples		
	$n = 20$	$n = 50$	$n = 100$
CVAE-semi	33.05	10.72	5.66
TripleGAN	3.06	1.80	1.29
SGAN	<b>1.68</b>	<b>1.23</b>	<b>0.93</b>



- Better results on semi-supervised classification
  - State-of-the-art results across multiple standard datasets
  - More advantages at low-shot settings

Method	MNIST		SVHN	CIFAR-10
	$n = 20$	$n = 50$	$n = 100$	$n = 1000$
Ladder [22]	-	-	<b>0.89(±0.50)</b>	-
VAE [12]	-	-	3.33(±0.14)	36.02(±0.10)
CatGAN [28]	-	-	1.39(±0.28)	-
ALI [5]	-	-	-	7.3
ImprovedGAN [27]	16.77(±4.52)	2.21(±1.36)	0.93(±0.07)	8.11(±1.3)
TripleGAN [15]	5.40(±6.53)	1.59(±0.69)	0.92(±0.58)	5.83(±0.20)
SGAN	<b>4.0(±4.14)</b>	<b>1.29(±0.47)</b>	<b>0.89(±0.11)</b>	<b>5.73(±0.12)</b>

- Controllable generation
  - Ablation studies reveal that  $\mathcal{R}_y$  and  $\mathcal{R}_z$  help align the semantics
- Visual quality
  - Report an inception score 6.91(±0.07), higher than that of TripleGAN and Improved-GAN w/o minibatch discrimination



- SGAN enables more interesting applications
  - Image progression: generate images with interpolated  $z$  – SGAN generalizes instead of memorizing data
  - Style transfer: infer  $z$  given an image, generate a new image with the same  $z$  but different semantic of interest  $y$

