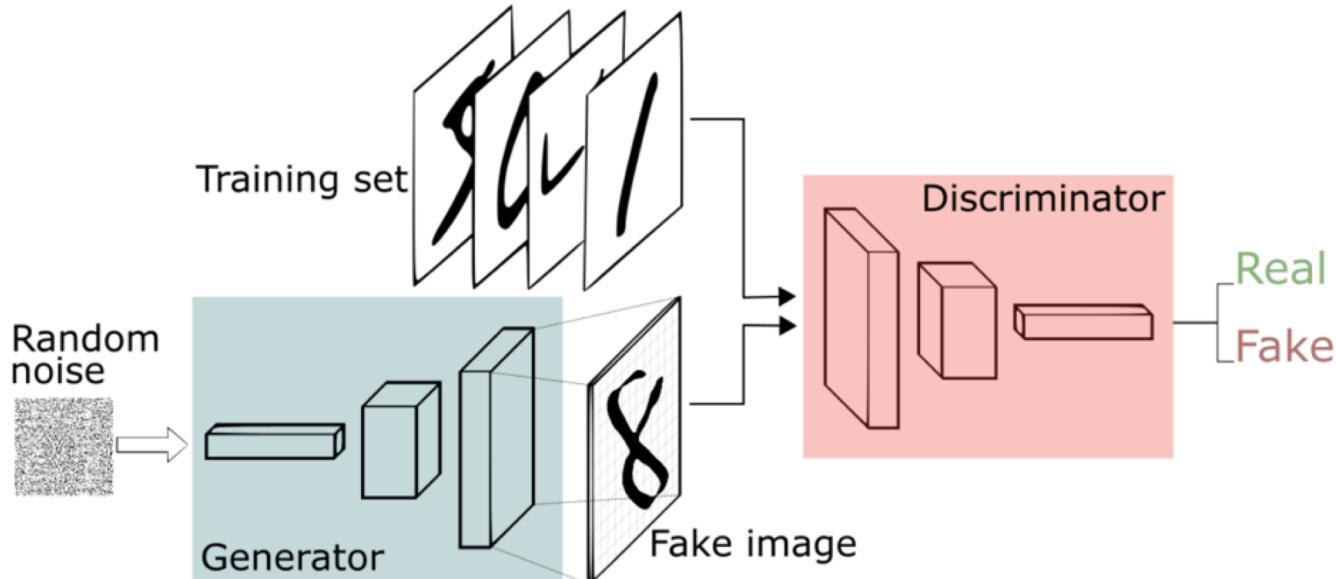


# VAE+GAN

Monday, Oct 8, 2018

Jianguo Zhang

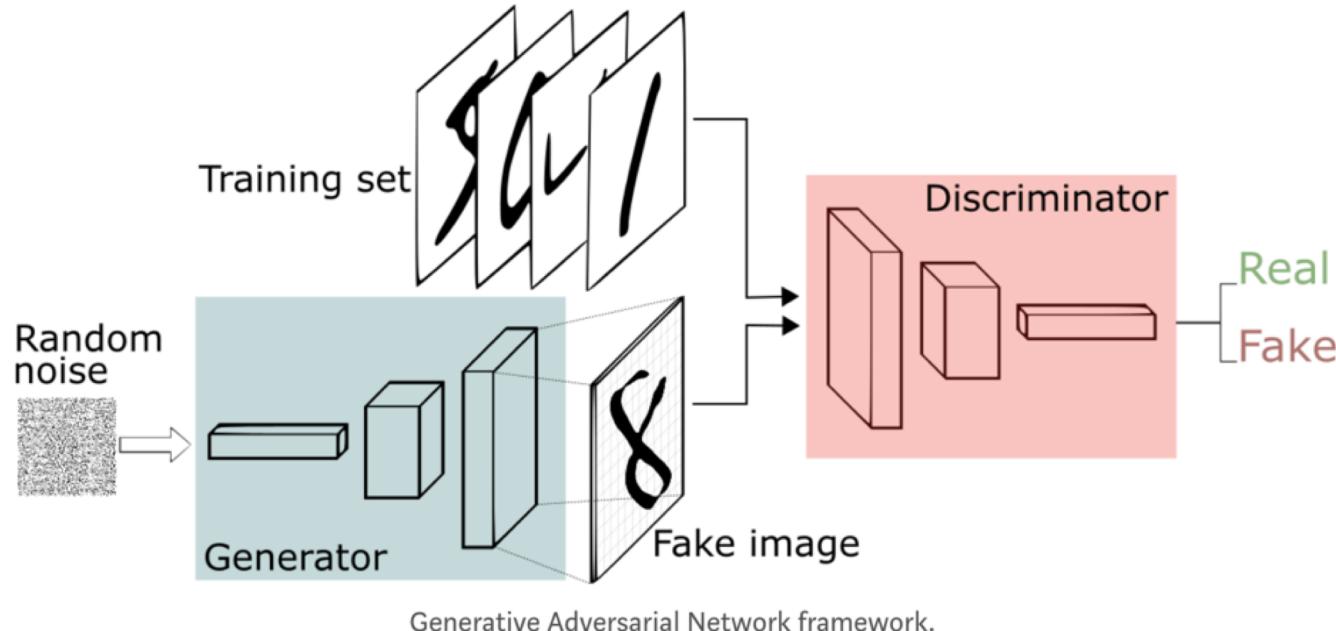
# Generative Adversarial Network (GAN)



Generative Adversarial Network framework.

- **Discriminator tries to correctly **distinguish** the true data and the fake data**
- **Generator tries to generate high-quality data to **fool** discriminator**
- **Ideally**, when D cannot distinguish the true and generated data, G can keeps generate realistic images

# Generative Adversarial Network (GAN)



**Probability that a real sample is classified as Real**

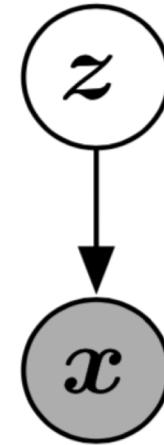
$$\min_G \max_D V(D, G) = E_{x \sim p_{data}(x)} [\log D(x)] + E_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$

**Probability that a fake sample is classified as Fake**

# Generative Adversarial Network (GAN)

## Attractiveness

- Generator Networks  $x = G(z; \theta^{(G)})$
- It is only required that,  $G$  is differentiable.
- So, having training data  $x \sim p_{\text{data}}(x)$   
what we want is a model that can draw samples  
 $x \sim p_{\text{model}}(x)$ , where  $p_{\text{model}} \approx p_{\text{data}}$
- **Don't write a formula for  $p_{\text{data}}(x)$ , just learn to draw sample directly.**



# Generative Adversarial Network (GAN)

**Drawbacks-1: Gradient Vanishing**

**Drawbacks-2: Training Instability**

**Drawbacks-3: Mode Collapse**

(Generate Undesirable Modes)

(Generate a Single Mode, lack of variety)

# Generative Adversarial Network (GAN)

## Drawbacks-1: Gradient Vanishing

using the original form of the objective of  $\mathbf{G}$

$$\mathbb{E}_{z \sim p(z)} [\log(1 - D(g_{\theta}(z)))]$$

will result in gradient vanishing issue of  $\mathbf{D}$  for  $\mathbf{G}$  because *intuitively*, at the very early phase of training,  $\mathbf{D}$  is very easy to be confident in detecting  $\mathbf{G}$ , so  $\mathbf{D}$  will output almost always 0

# Generative Adversarial Network (GAN)

## Drawbacks-1: Gradient Vanishing

using the original form of the objective of  $\mathbf{G}$

$$\mathbb{E}_{z \sim p(z)} [\log(1 - D(g_{\theta}(z)))]$$

will result in gradient vanishing issue of  $\mathbf{D}$  for  $\mathbf{G}$   
because *theoretically*, when  $\mathbf{D}$  is *optimal*,  
minimizing the loss is equal to minimizing the *JS divergence* (Arjovsky & Bottou, 2017)

# Generative Adversarial Network (GAN)

## Drawbacks-1: Gradient Vanishing

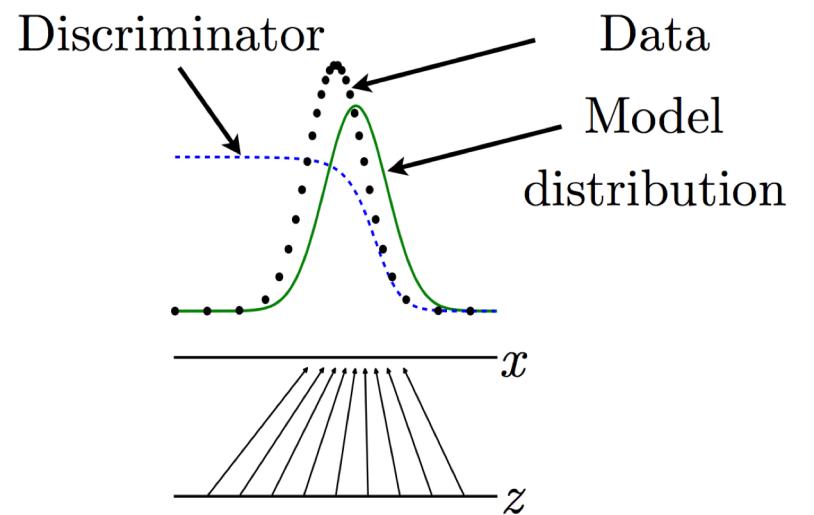
$$L(D, g_\theta) = \mathbb{E}_{x \sim \mathbb{P}_r} [\log D(x)] + \mathbb{E}_{x \sim \mathbb{P}_g} [\log(1 - D(x))]$$

### Derivative

$$\frac{\delta}{\delta D(\mathbf{x})} J^{(D)} = 0.$$

$$D^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\text{model}}(\mathbf{x})}.$$

$$L(D^*, g_\theta) = 2JSD(\mathbb{P}_r \parallel \mathbb{P}_g) - 2 \log 2$$



# Generative Adversarial Network (GAN)

## Drawbacks-1: Gradient Vanishing

$$L(D^*, g_\theta) = 2JSD(\mathbb{P}_r \parallel \mathbb{P}_g) - 2\log 2$$

JSD

$$D_{JS}(p|q) = D_{JS}(q|p) = \frac{1}{2}D_{KL}(p|r) + \frac{1}{2}D_{KL}(q|r)$$

$$r = \frac{1}{2}(p + q)$$

The JS divergence for the two distributions  $P_r$  and  $P_g$  is (almost) always log2 because  $P_r$  and  $P_g$  hardly can overlap (Arjovsky & Bottou, 2017, Theorem 2.1~2.3)

# Generative Adversarial Network (GAN)

## Drawbacks-2: Training Instability

It is heuristically motivated that generator can still learn even when discriminator successfully rejects all generator samples, but not theoretically guaranteed

using the alternative form of the objective of  $\mathbf{G}$

$$\mathbb{E}_{z \sim p(z)} [-\log D(g_\theta(z))]$$

will result in gradient unstable issue and mode missing problem because *theoretically*, when  $\mathbf{D}$  is optimal, minimizing the loss is equal to **minimizing** the *KL divergence* meanwhile **maximizing** the *JS divergence* (Arjovsky & Bottou, 2017, Theorem 2.5):

$$KL(\mathbb{P}_{g_\theta} \parallel \mathbb{P}_r) - 2JSD(\mathbb{P}_{g_\theta} \parallel \mathbb{P}_r)$$

# Generative Adversarial Network (GAN)

## Drawbacks-3: Mode Collapse

minimizing the *KL divergence* only is biased:

$$KL(\mathbb{P}_{g_\theta} \parallel \mathbb{P}_r) - 2JSD(\mathbb{P}_{g_\theta} \parallel \mathbb{P}_r)]$$

because *KL divergence* is asymmetric, and thus it is not equally treated when  $\mathbf{G}$  generates an unreal sample and when  $\mathbf{G}$  fails to generate real sample

Therefore,  $\mathbf{G}$  will generate too many few-mode (less diverse) but real samples , a safer strategy

# Generative Adversarial Network (GAN)

Wasserstein GANs (Arjovsky et al., 2017)

Wasserstein-1 Distance (Earth-Mover Distance):

$$W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|]$$

- *Why is it superior to KL and JS divergence?*

# Generative Adversarial Network (GAN)

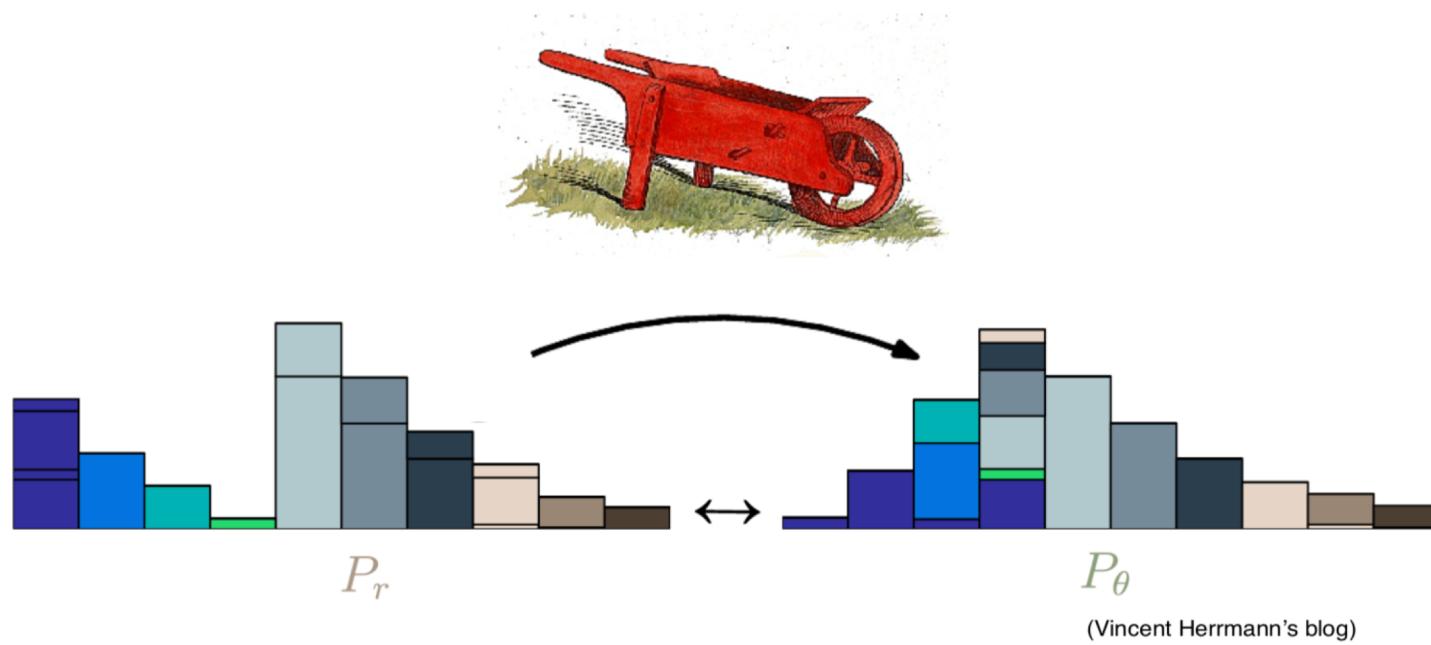
- Wasserstein-1 Distance (Earth-Mover Distance):

$$W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|]$$

where  $\Pi(\mathbb{P}_r, \mathbb{P}_g)$  denotes the set of all joint distributions  $\gamma(x, y)$  whose marginals are respectively  $\mathbb{P}_r$  and  $\mathbb{P}_g$ . Intuitively,  $\gamma(x, y)$  indicates how much “mass” must be transported from  $x$  to  $y$  in order to transform the distributions  $\mathbb{P}_r$  into the distribution  $\mathbb{P}_g$ . The EM distance then is the “cost” of the optimal transport plan.

# Generative Adversarial Network (GAN)

- Wasserstein-1 Distance (Earth-Mover Distance):



- It is continuous everywhere and
- differentiable almost everywhere.
- *And most importantly, it can reflect the distance of two distributions even if they do not overlap, and thus can provide meaningful gradients*

# Wasserstein GANs

- This new value function of WGAN gives rise to the additional requirement that the discriminator must lie within in the space of 1-Lipschitz functions:

$$\min_G \max_{D \in \mathcal{D}} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r} [D(\mathbf{x})] - \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g} [D(\tilde{\mathbf{x}})]$$

- In other words,  $\mathcal{D}$  is the set of 1-Lipschitz functions
  - Lipschitz continuity

# Wasserstein GANs

## Lipschitz Continuity

- real-value function:  $f: R \rightarrow R$
- positive constant:  $K$

$$|f(x_1) - f(x_2)| \leq K|x_1 - x_2|$$

- In other words, a Lipschitz continuous function has bounded first derivative. Intuitively, the slope of a KK-Lipschitz function never exceeds KK, for a more general definition of slope.

$$d_Y(f(x_1), f(x_2)) \leq K d_X(x_1, x_2)$$

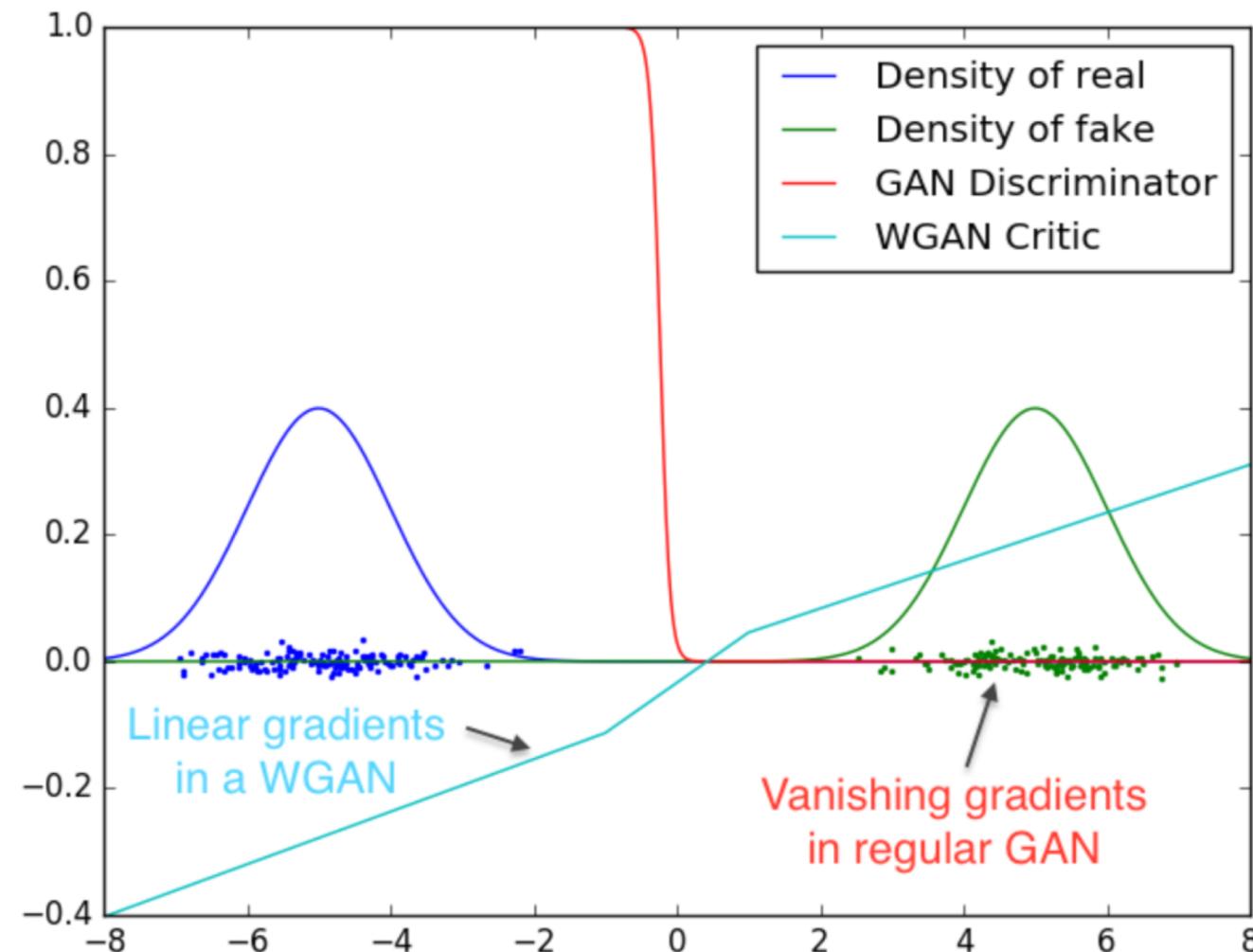
# Wasserstein GANs

- This new value function of WGAN gives rise to the additional requirement that the discriminator must lie within in the space of 1-Lipschitz functions:

$$\min_G \max_{D \in \mathcal{D}} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r} [D(\mathbf{x})] - \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g} [D(\tilde{\mathbf{x}})]$$

- To satisfy this requirement, WGAN enforces the weights of  $D$  lie within a compact space  $[-c, c]$  by applying **weight clipping**

# Wasserstein GANs



# Improving Variational Encoder-Decoders in Dialogue Generation

## VAE

$$\begin{aligned}-\log p_\theta(x) &\leq -\log p_\theta(x) + \text{KL}(q_\phi(z|x)||p_\theta(z|x)) \\&= -\mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)] + \text{KL}(q_\phi(z|x)||p(z))\end{aligned}$$

The weaker the decoder family is, the more it will be biased to utilizing latent variables.

A more flexible prior distribution  $p_\theta(z)$  will also increase the chance as it provides more possibilities for the utilisation .

## CVAE

$$-\mathbb{E}_{q_\phi(z|x,c)}[\log p_\theta(x|c,z)] + \text{KL}(q_\phi(z|x,c)||p_\theta(z|c))$$

CVAE can be rewrite as

$$-\log \int_z p_\theta(z|c)p_\theta(x|z,c)dz + \text{KL}(q_\phi(z|x,c)||p_\theta(z|x,c))$$

Let's first take a look at the first item,  $\log \int_z p_\theta(z|c)p_\theta(x|z,c)dz = \log p_\theta(x|c)$ . When the family of  $p_\theta(x|z,c)$  is complex enough and includes the real distribution of  $x$ , the optimal value of this item is  $p(x|c)$  and the reliance on  $z$  is not necessary. However, reliance on  $z$  provides the model with a chance of taking advantage of  $z$ 's distribution and reduces the complexity requirement for the distribution family  $p_\theta(x|z,c)$ . For

CVAE can be rewrite as

$$-\log \int_z p_\theta(z|c)p_\theta(x|z, c)dz + \text{KL}(q_\phi(z|x, c)||p_\theta(z|x, c))$$

**KL divergence**

$$p_\theta(z|x, c) = \frac{p_\theta(x|z, c)p_\theta(z|c)}{p_\theta(x|c)}$$

modelled by  $q_\phi(z|x)$  (Hinton et al. 1995). However, when  $x$  represents sentences with variable length, the value of  $p_\theta(x|z)$  vanishes greatly when the length grows, which makes the adjusting task much more difficult. This implies the second item will always prefer ignoring the latent variables, so long as the approximated posterior is not powerful enough to perfectly match the real posterior. The weaker the approximating

The weaker the decoder family is, the more it will be biased to utilizing latent variables.

A more flexible prior distribution  $p_\theta(z)$  will also increase the chance as it provides more possibilities for the utilisation .

The weaker the approximating posterior distribution family is, the more it will be biased to ignoring latent variables

# Proposed Method

## Original: Adversarial Encoder-Decoder (AED)

$p_\theta(z|c)$  are implicitly defined by passing context-dependent Gaussian random variables  $\epsilon$  through multi-layer perceptions

## New: Replacing GAN with VAE

$$-\mathbb{E}_{q_\phi(\epsilon|c, \tilde{z})} p_\theta(z|c, \epsilon) + KL(q_\phi(\epsilon|c, \tilde{z}) || p_\theta(\epsilon|c))$$

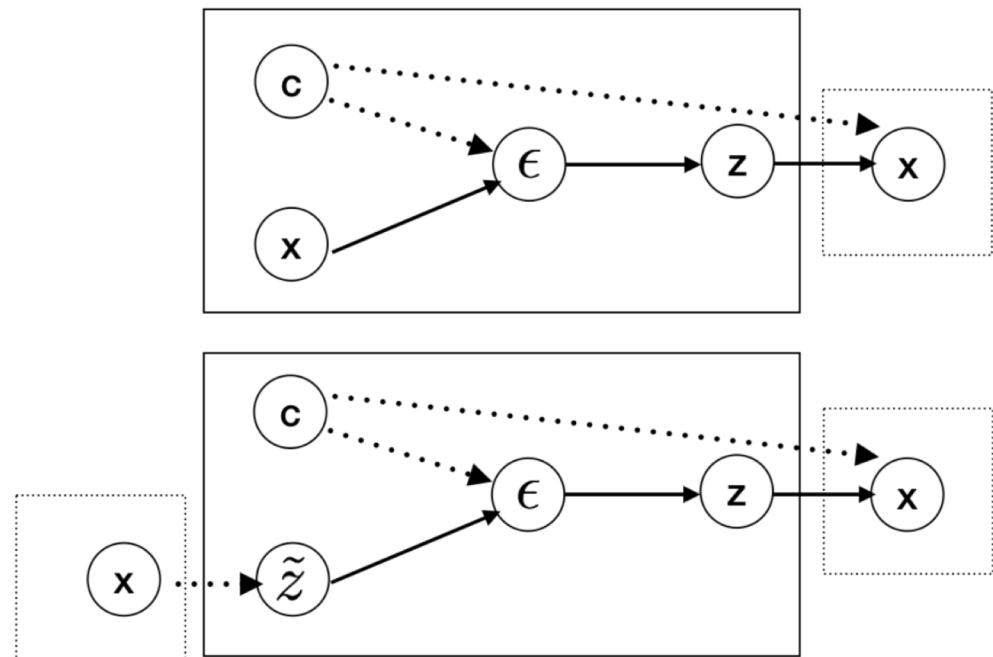


(approximated posterior)

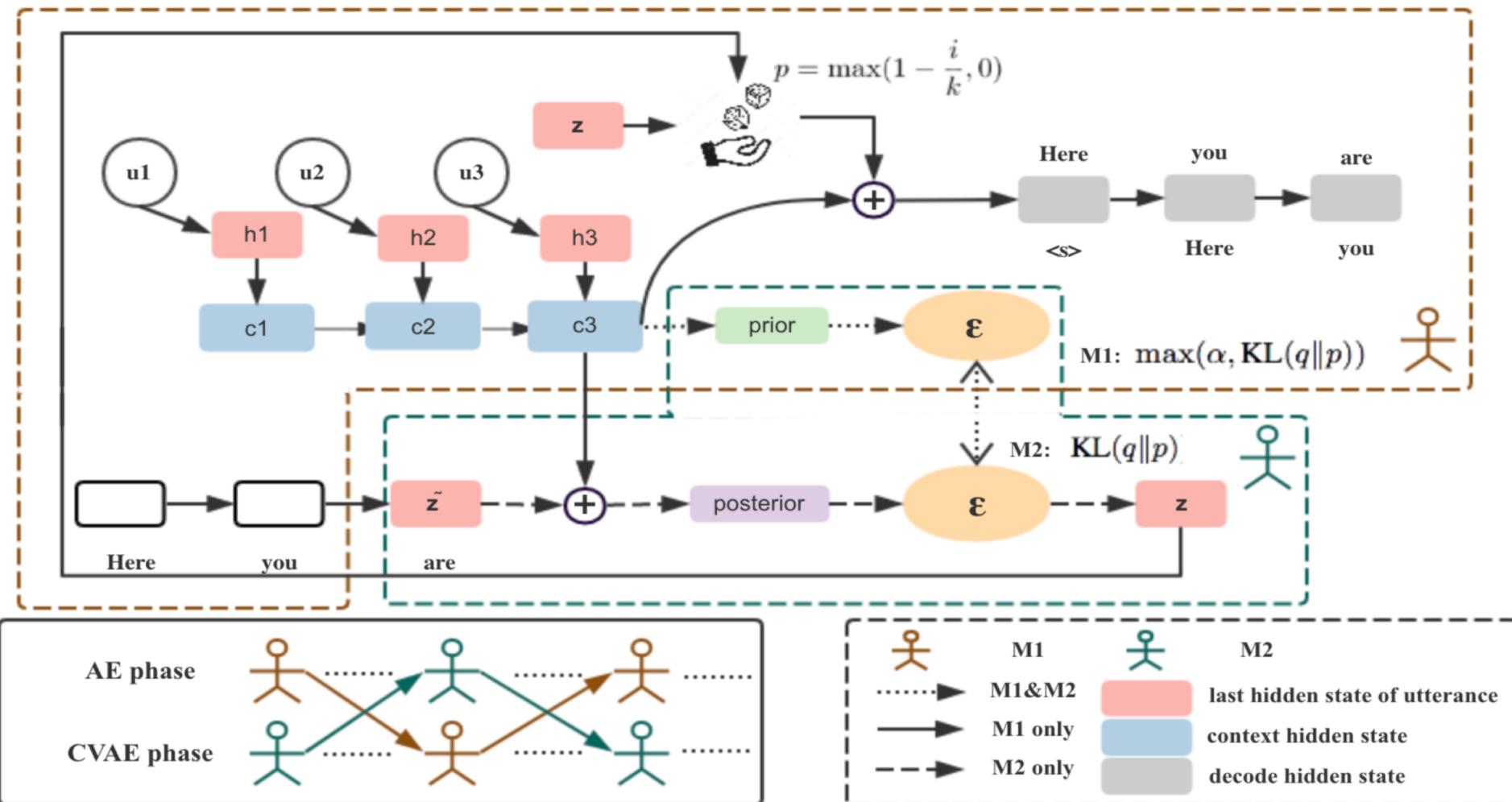
**real hidden vector  $\tilde{z}$  or the noisy  $z$**

The accuracy of the CVAE objective relies on the matching degree of  $q_\phi(\epsilon | c, \tilde{z})$  and  $p_\theta(\epsilon | c, \tilde{z})$

**Original and New have the same global optimum**



# Proposed Method



$$-\mathbb{E}_{q_\phi(\epsilon|c, \tilde{z})} p_\theta(z|c, \epsilon) + KL(q_\phi(\epsilon|c, \tilde{z})||p_\theta(\epsilon|c))$$

# Training Process

## CVAE phase

$$\min_{\phi} \text{KL}(q_{\phi}(\epsilon|\tilde{z}, c) \| p_{\phi}(\epsilon|c)) + \frac{1}{2} \mathbb{E}_{q_{\phi}(\epsilon|\tilde{z}, c)} \|g_{\phi}(\epsilon) - \tilde{z}\|_2^2;$$

$$\tilde{z} = f_{\theta}(x)$$

A sample  $\tilde{z}$  is obtained from the AE

## AE phase

$$\min_{\theta} \max(\alpha, \text{KL}(q_{\phi}(\epsilon|\tilde{z}, c) \| p_{\phi}(\epsilon|c)))$$

$$- \mathbb{E}_{q_{\phi}(z|\tilde{z}, c)} [\log(p_{\theta}(x|z, c))];$$

$$\tilde{z} = f_{\theta}(x), z = (1-p)g_{\phi}(\epsilon) + p\tilde{z}$$

The corresponding latent variable  $z$  is sampled from the posterior distribution  $q_{\phi}(z|\tilde{z}, c)$  provided by the CVAE part

The output of the CVAE part are latent variables, which can represent a much broader distribution family than mean-field Gaussian

# DialogWAE: Multimodal Response Generation with Conditional Wasserstein Auto-Encoder

# Data-driven Conversation

## **Safe Response Problem**

**Fail to generate meaningful, diverse on-topic responses**

## **the Posterior Collapse Problem**

**The decoder learns to ignore the latent variable and degrades to a vanilla RNN**

# GAN-based Methods

**However, training with REINFORCE has been observed to be unstable due to the high variance of the gradient estimate**

**A new approach injects diversity in the word level rather than the level of the whole responses**

Our main contributions are two-fold: (1) A novel GAN-based model for neural dialogue modeling, which employs GAN to generate samples of latent variables. (2) A Gaussian mixture prior network to sample random noise from a multimodal prior distribution. **To the best of our knowledge**, the proposed DialogWAE is the first GAN conversation model that exploits multimodal latent structures.

of the whole responses. DialogWAE differs from exiting GAN conversation models in that it shapes the distribution of responses in a high level latent space rather than direct tokens and does not rely on RL where the gradient variances are large.

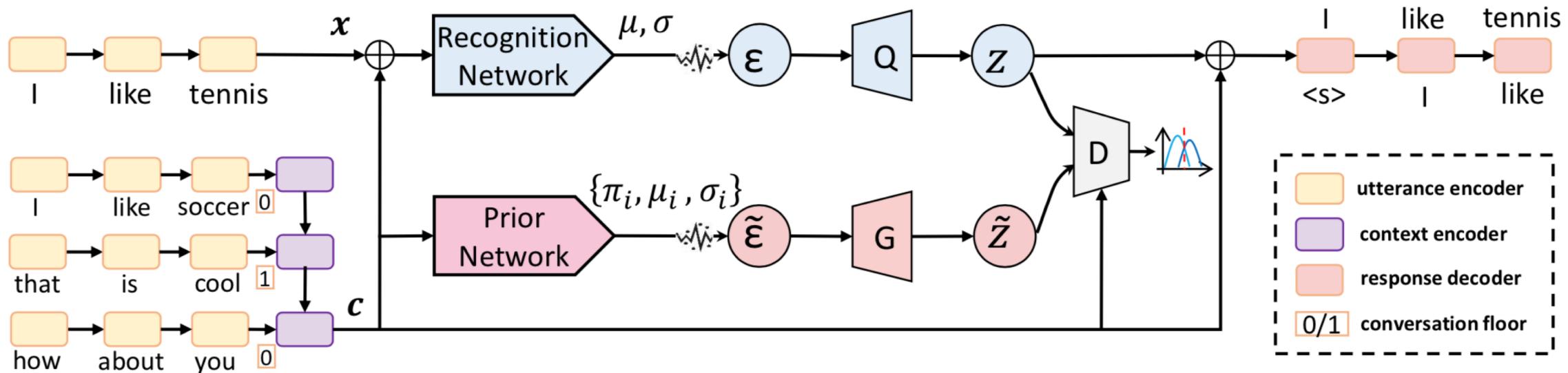


Figure 1: Architecture of DialogWAE

Let  $d=[u_1, \dots, u_k]$  denote a dialogue of  $k$  utterances where  $u_i=[w_1, \dots, w_{|u_i|}]$  represents an utterance and  $w_n$  denotes the  $n$ -th word in  $u_i$ . Let  $c=[u_1, \dots, u_{k-1}]$  denote a dialogue context, the  $k-1$  historical utterances, and  $x=u_k$  be a response which means the next utterance. Our goal is to estimate the conditional distribution  $p_\theta(x|c)$ .

Log-likelihood of a response

$$p_\theta(x|c) = \int_z p(x|c, z)p(z|c)dz.$$

Evidence Lower Bound (ELBO)

$$\begin{aligned} \log p_\theta(x|c) &= \log \int_z p(x|c, z)p(z|c)dz \\ &\geq \ell(x, c) = \mathbf{E}_{z \sim q_\phi(z|x, c)}[\log p_\psi(x|c, z)] - \text{KL}(q_\phi(z|x, c)||p(z|c)), \end{aligned}$$

approximate posterior sample  $z \sim q_\phi(z|c, x)$

*prior network and recognition network*

$$\tilde{z} = G_\theta(\tilde{\epsilon}), \quad \tilde{\epsilon} \sim \mathcal{N}(\epsilon; \tilde{\mu}, \tilde{\sigma}^2 I), \quad \begin{bmatrix} \tilde{\mu} \\ \log \tilde{\sigma}^2 \end{bmatrix} = \tilde{W}f_\theta(c) + \tilde{b}$$

$$z = Q_\phi(\epsilon), \quad \epsilon \sim \mathcal{N}(\epsilon; \mu, \sigma^2 I), \quad \begin{bmatrix} \mu \\ \log \sigma^2 \end{bmatrix} = Wg_\phi(\begin{bmatrix} x \\ c \end{bmatrix}) + b,$$

Goal

$$\min_{\theta, \phi, \psi} -E_{q_\phi(z|x, c)} \log p_\psi(x|z, c) + W(q_\phi(z|x, c)||p_\theta(z|c)),$$

Reconstruction loss

$$\mathcal{L}_{rec} = -E_{z=Q(\epsilon), \epsilon \sim \text{RecNet}(x, c)} \log p_\psi(x|c, z)$$

Discriminator loss

$$\mathcal{L}_{disc} = E_{\epsilon \sim \text{RecNet}(x, c)}[D(Q(\epsilon), c)] - E_{\tilde{\epsilon} \sim \text{PriNet}(c)}[D(G(\tilde{\epsilon}), c)]$$

situations [Sato *et al.*, 2017], topics and sentiments. A random noise with normal distribution could restrict the generator to output a latent space with a single dominant mode due to the single-modal nature of Gaussian distribution. Consequently, the generated responses could follow simple prototypes.

Mixture of Gaussian distributions

$$\text{GMM}(\{\pi_k, \mu_k, \sigma_k^2 I\}_{k=1}^K)$$

$$p(\epsilon|c) = \sum_{k=1}^K v_k \mathcal{N}(\epsilon; \mu_k, \sigma_k^2 I),$$

$$\pi_k = \frac{\exp(e_k)}{\sum_{i=1}^K \exp(e_i)}, \text{ where } \begin{bmatrix} e_k \\ \mu_k \\ \log \sigma_k^2 \end{bmatrix} = W f_\theta(c) + b$$

Gumbel-Softmax re-parametrization

$$v_k = \frac{\exp((e_k + g_k)/\tau)}{\sum_{i=1}^K \exp((e_i + g_i)/\tau)},$$

$$g_i = -\log(-\log(u_i)), u_i \sim U(0, 1)$$

---

**Algorithm 1:** DialogWAE Training (UEnc: utterance encoder; CEnc: context encoder; RecNet: recognition network; PriNet: prior network; Dec: decoder) K=3,  $n_{\text{disc}}=5$  in all experiments

---

**In:** a dialog corpus  $\mathcal{D}=\{(c_i, x_i)\}_{i=1}^{|\mathcal{D}|}$ , the number of prior modes  $K$ , discriminator iterations  $n_{\text{critic}}$

1 Initialize  $\{\theta_{\text{UEnc}}, \theta_{\text{CEnc}}, \theta_{\text{PriNet}}, \theta_{\text{RecNet}}, \theta_Q, \theta_G, \theta_D, \theta_{\text{Dec}}\}$

2 **while** *not convergence* **do**

3     Initialize  $\mathcal{D}$

4     **while**  $\mathcal{D}$  has unsampled batches **do**

5         Sample a mini-batch of  $N$  instances  $\{(x_n, c_n)\}_{n=1}^N$  from  $\mathcal{D}$

6         Get the representations of context and response  $x_n = \text{UEnc}(x_n)$ ,  $c_n = \text{CEnc}(c_n)$

7         Sample  $\epsilon_n$  from  $\text{RecNet}(x_n, c_n)$  according to Equation 4

8         Sample  $\hat{\epsilon}_n$  from  $\text{PriNet}(c_n, K)$  according to Equation 8–10

9         Generate  $z_n = Q(\epsilon_n)$ ,  $\tilde{z}_n = G(\hat{\epsilon}_n)$

10         Update  $\{\theta_Q, \theta_G, \theta_{\text{PriNet}}, \theta_{\text{RecNet}}\}$  by gradient ascent on discriminator loss

11              $\mathcal{L}_{\text{disc}} = \frac{1}{N} \sum_{n=1}^N D(z_n, c_n) - \frac{1}{N} \sum_{n=1}^N D(\tilde{z}_n, c_n)$

12         **for**  $i \in \{1, \dots, n_{\text{critic}}\}$  **do**

13             | Repeat 5–9

14             | Update  $\theta_D$  by gradient descent on the discriminator loss  $\mathcal{L}_{\text{disc}}$  with gradient penalty

15         **end**

16         Update  $\{\theta_{\text{UEnc}}, \theta_{\text{CEnc}}, \theta_{\text{RecNet}}, \theta_Q, \theta_{\text{Dec}}\}$  by gradient descent on the reconstruction loss

17              $\mathcal{L}_{\text{rec}} = -\frac{1}{N} \sum_{n=1}^N \log p(x_n | z_n, c_n)$

18     **end**

19 **end**