

# A Surprisingly Effective Fix for Deep Latent Variable Modeling of Text

Bohan Li<sup>\*1</sup>, Junxian He<sup>\*1</sup>, Graham Neubig<sup>1</sup>, Taylor Berg-Kirkpatrick<sup>2</sup>, Yiming Yang<sup>1</sup>

## Core idea:

- Pretraining encoder using AE objective
- Initializing VAE with pertained AE and training with FB objective

## Results:

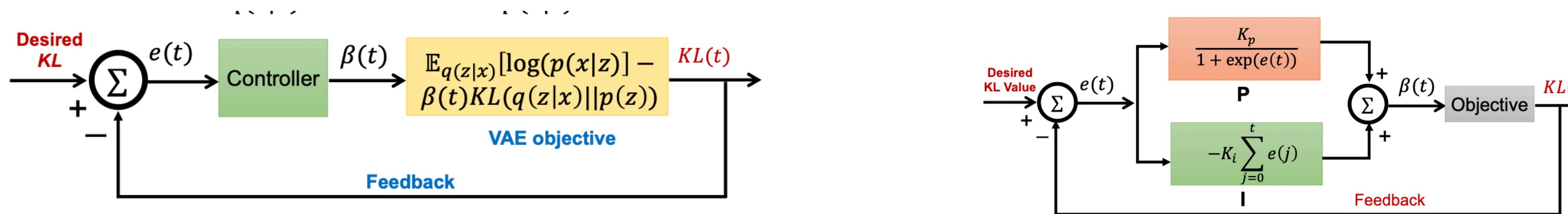
	Yahoo				
LSTM-LM	60.75	-	-	-	-
VAE	61.52	329.10	0	0.00	329.10
+ anneal	61.21	328.80	0	0.00	328.80
+ cyclic	66.93	333.80	4	2.83	336.63
+ aggressive	59.77	322.70	15	5.70	328.40
+ FBP ( $\lambda = 9$ )	62.59	322.91	6	9.08	331.99
+ FBP ( $\lambda = 7$ )	62.76	324.66	5	7.03	331.69
+ FBP ( $\lambda = 5$ )	62.78	326.26	3	5.07	331.32
+ FBP ( $\lambda = 3$ )	62.88	328.13	2	3.06	331.19
Ours ( $\lambda = 6$ )	59.23	317.39	32	12.09	329.48
Ours ( $\lambda = 8$ )	<b>59.51</b>	<b>315.31</b>	<b>32</b>	15.02	330.33
Ours ( $\lambda = 9$ )	59.60	315.09	32	15.49	330.58

# ControVAE: Controllable Variational Autoencoder (ICML 2020)

Huajie Shao et al

## Core ideas:

- Enhance  $\beta$ -VAE by controlling  $\beta$  with a new non-linear PI controller,  $\beta(t)$ . The controller is not learnable, with  $K_p$ ,  $K_i$ , Desired KL being hyperparameters.



- Intuition: Comparing error between Desired KL and Sampled KL (from inference output); if Sampled KL is too small,  $\beta(t)$  will be small so KL-divergence is encouraged to grow, and vice versa.

# ControlVAE: Controllable Variational Autoencoder (ICML 2020)

Huajie Shao et al

## Core ideas:

- PI Algorithm

---

### Algorithm 1 PI algorithm.

---

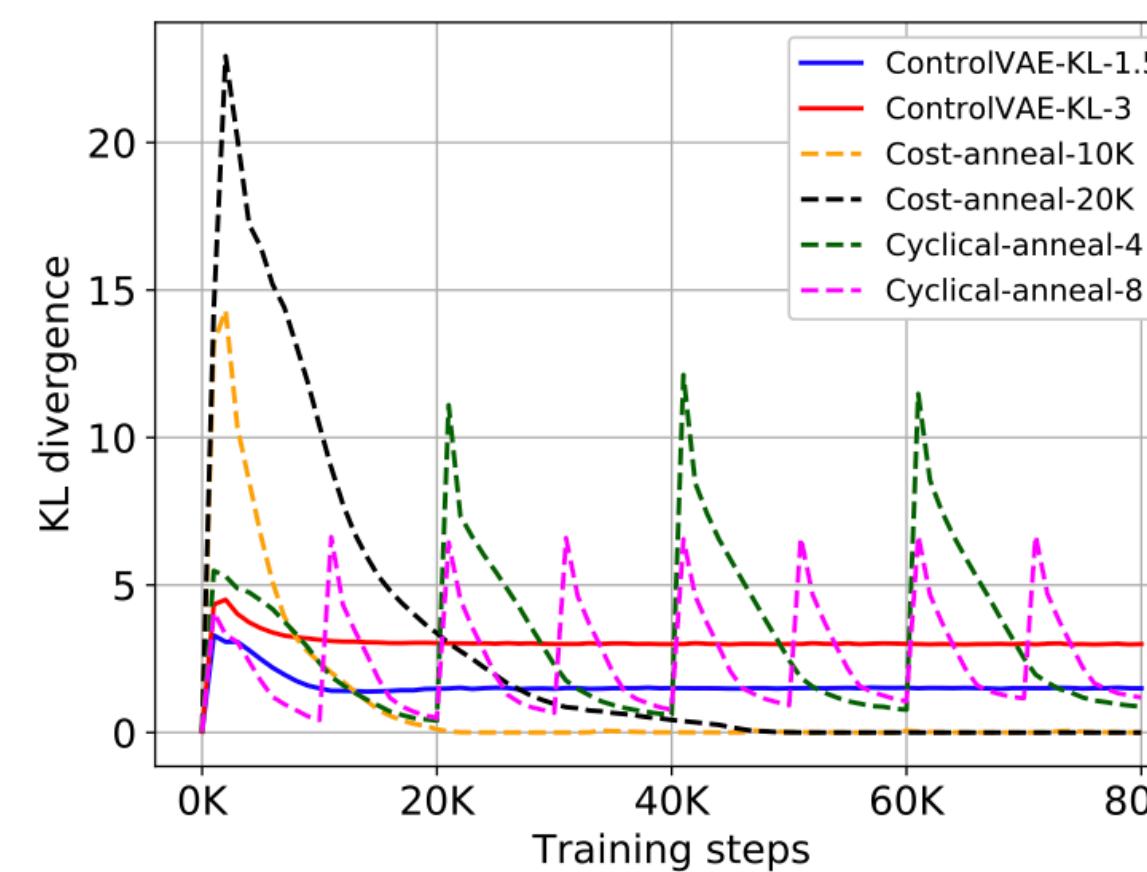
```
1: Input: desired KL  $v_{kl}$ , coefficients  $K_p, K_i$ , max/min value  
    $\beta_{max}, \beta_{min}$ , iterations  $N$   
2: Output: hyperparameter  $\beta(t)$  at training step  $t$   
3: Initialization:  $I(0) = 0, \beta(0) = 0$   
4: for  $t = 1$  to  $N$  do  
5:   Sample KL-divergence,  $\hat{v}_{kl}(t)$   
6:    $e(t) \leftarrow v_{kl} - \hat{v}_{kl}(t)$   
7:    $P(t) \leftarrow \frac{K_p}{1+\exp(e(t))}$   
8:   if  $\beta_{min} \leq \beta(t-1) \leq \beta_{max}$  then  
9:      $I(t) \leftarrow I(t-1) - K_i e(t)$   
10:    else  
11:       $I(t) = I(t-1)$  // Anti-windup  
12:    end if  
13:     $\beta(t) = P(t) + I(t) + \beta_{min}$   
14:    if  $\beta(t) > \beta_{max}$  then  
15:       $\beta(t) = \beta_{max}$   
16:    end if  
17:    if  $\beta(t) < \beta_{min}$  then  
18:       $\beta(t) = \beta_{min}$   
19:    end if  
20:  Return  $\beta(t)$   
21: end for
```

---

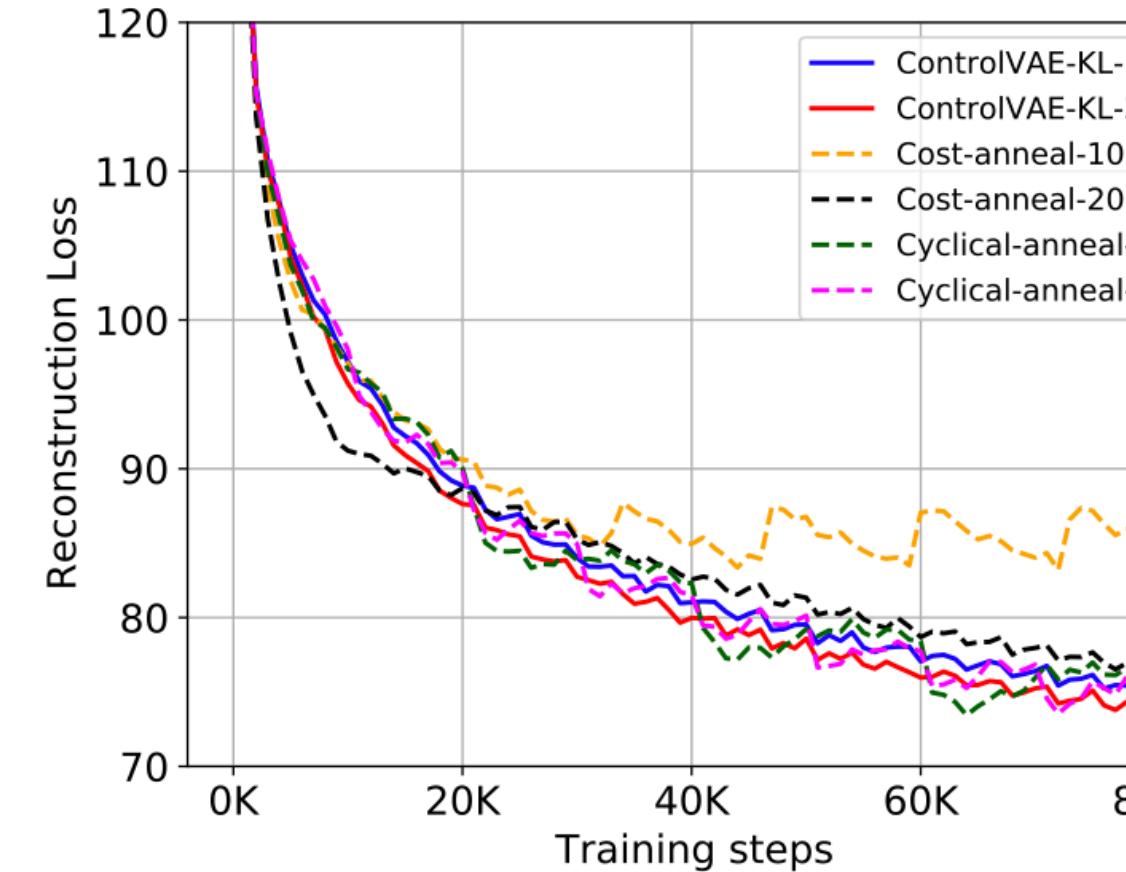
# ControlVAE: Controllable Variational Autoencoder (ICML 2020)

Huajie Shao et al

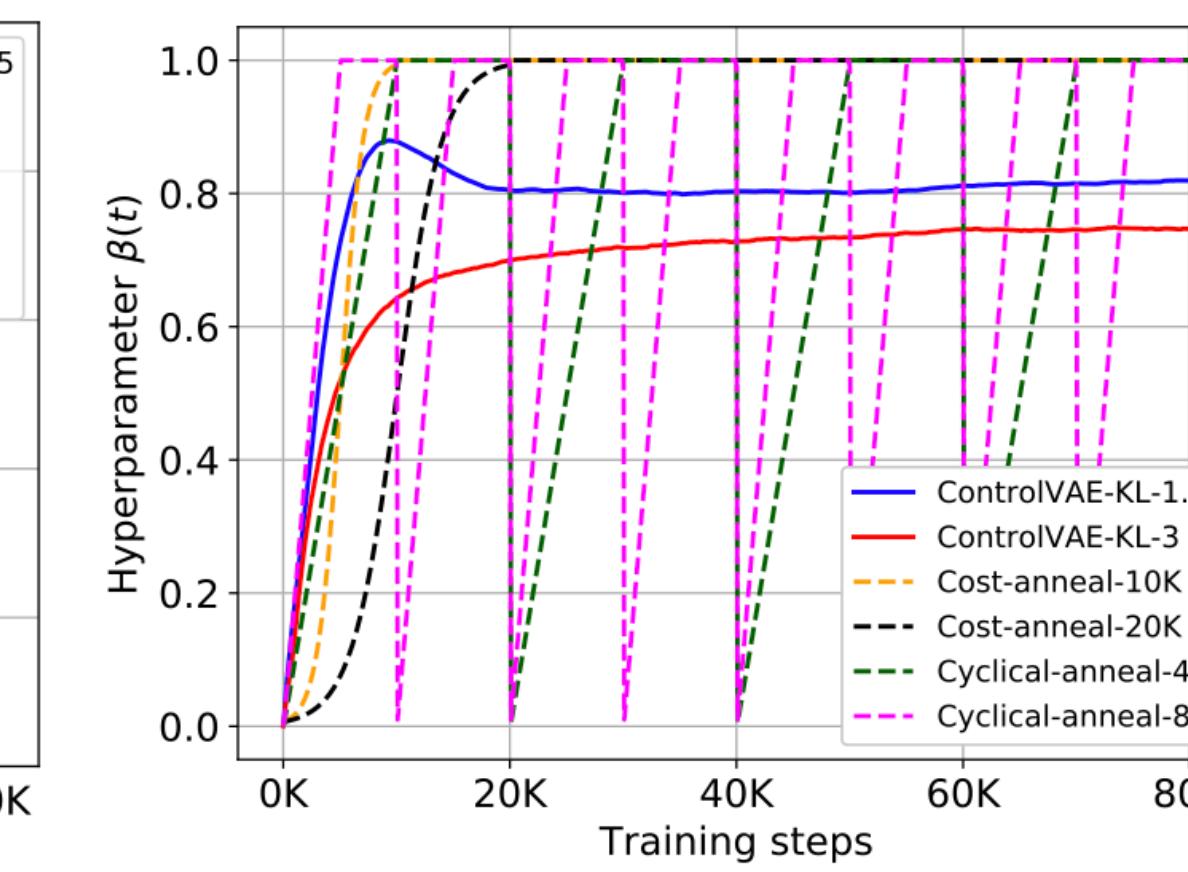
## Results:



(a) KL divergence



(b) Reconstruction loss



(c)  $\beta(t)$

Methods/metric	Dis-1	Dis-2	self-BLEU-2	self-BLEU-3	PPL
ControlVAE-KL-35	<b>6.27K</b> $\pm$ 41	<b>95.86K</b> $\pm$ 1.02K	<b>0.663</b> $\pm$ 0.012	<b>0.447</b> $\pm$ 0.013	<b>8.81</b> $\pm$ 0.05
ControlVAE-KL-25	6.10K $\pm$ 60	83.15K $\pm$ 4.00K	0.698 $\pm$ 0.006	0.495 $\pm$ 0.014	12.47 $\pm$ 0.07
Cost anneal-KL-17	5.71K $\pm$ 87	69.60K $\pm$ 1.53K	0.721 $\pm$ 0.010	0.536 $\pm$ 0.008	16.82 $\pm$ 0.11
Cyclical (KL = 21.5)	5.79K $\pm$ 81	71.63K $\pm$ 2.04K	0.710 $\pm$ 0.007	0.524 $\pm$ 0.008	17.81 $\pm$ 0.33

# Implicit Deep Latent Variable Models for Text Generation (EMNLP 2019)

Le Fang<sup>†</sup>, Chunyuan Li<sup>§</sup>, Jianfeng Gao<sup>§</sup>, Wen Dong<sup>†</sup>, Changyou Chen<sup>†</sup>

## Core ideas:

- Attributing posterior collapse to the restrictive Gaussian assumption, and advocate more flexible sample-based posterior representation.
- Proposed iVAE:
  - Instead of assuming posterior as Gaussian, use sampling mechanism for  $q_\phi(z|x)$ .  $z_{x,i} = G_\phi(x, \varepsilon_i)$ ,  $\varepsilon_i \sim q(\varepsilon)$
  - How to evaluate KL on this? Use dual form of  $\text{KL}(q_\phi(z|x) \parallel p(z))$ :

$$\begin{aligned} & \text{KL}(q_\phi(z|x) \parallel p(z)) \\ &= \max_\nu \mathbb{E}_{z \sim q_\phi(z|x)} \nu_\psi(x, z) - \mathbb{E}_{z \sim p(z)} \exp(\nu_\psi(x, z)), \end{aligned} \tag{7}$$

- New iVAE objective:

$$\begin{aligned} \mathcal{L}_{\text{iVAE}} &= \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{z \sim q_\phi(z|x)} \log p_\theta(x|z) \\ &\quad - \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{z \sim q_\phi(z|x)} \nu_\psi(x, z) \\ &\quad + \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{z \sim p(z)} \exp(\nu_\psi(x, z)), \end{aligned} \tag{8}$$

- Here,  $\nu_\psi$  is a MLP taking in  $(x, z)$ .

# Implicit Deep Latent Variable Models for Text Generation (EMNLP 2019)

Le Fang<sup>†</sup>, Chunyuan Li<sup>§</sup>, Jianfeng Gao<sup>§</sup>, Wen Dong<sup>†</sup>, Changyou Chen<sup>†</sup>

## Core ideas:

- Training Scheme:

- Sample a mini-batch of  $\mathbf{x}_i \sim \mathcal{D}$ ,  $\epsilon_i \sim q(\epsilon)$ , and generate  $\mathbf{z}_{\mathbf{x}_i, \epsilon_i} = G(\mathbf{x}_i, \epsilon_i; \phi)$ ; Sample a mini-batch of  $\mathbf{z}_i \sim p(\mathbf{z})$ .
- Update  $\psi$  in  $\nu_\psi(\mathbf{x}, \mathbf{z})$  to maximize

$$\sum_i \nu_\psi(\mathbf{x}_i, \mathbf{z}_{\mathbf{x}_i, \epsilon_i}) - \sum_i \exp(\nu_\psi(\mathbf{x}_i, \mathbf{z}_i)) \quad (9)$$

- Update parameters  $\{\phi, \theta\}$  to maximize

$$\sum_i \log p_\theta(\mathbf{x}_i | \mathbf{z}_{\mathbf{x}_i, \epsilon_i}) - \sum_i \nu_\psi(\mathbf{x}_i, \mathbf{z}_{\mathbf{x}_i, \epsilon_i}) \quad (10)$$

# Implicit Deep Latent Variable Models for Text Generation (EMNLP 2019)

Le Fang<sup>†</sup>, Chunyuan Li<sup>§</sup>, Jianfeng Gao<sup>§</sup>, Wen Dong<sup>†</sup>, Changyou Chen<sup>†</sup>

## Core ideas:

- Proposed Mutual Information Regularized iVAE
  - Replace  $-\text{KL}(q_\phi(z|x) \parallel p(z))$  with  $-\text{KL}(q_\phi(z) \parallel p(z))$ , where  $q_\phi(z) = \int q(x)q_\phi(z|x)dx$ , estimated by ancestral sampling in practice.
  - This objective also maximize the mutual information  $I(x,z)$ , as they claimed.
  - The new objective:

$$\begin{aligned} & \text{KL}(q_\phi(z) \parallel p(z)) \\ &= \max_{\nu} \mathbb{E}_{z \sim q_\phi(z)} \nu_\psi(z) - \mathbb{E}_{z \sim p(z)} \exp(\nu_\psi(z)). \end{aligned} \quad (13)$$

$$\begin{aligned} \mathcal{L}_{\text{iVAE}_{\text{MI}}} &= \mathbb{E}_{x \sim D} \mathbb{E}_{z \sim q_\phi(z|x)} \log p_\theta(x|z) \\ &\quad - \mathbb{E}_{z \sim q_\phi(z)} \nu_\psi(z) + \mathbb{E}_{z \sim p(z)} \exp(\nu_\psi(z)), \end{aligned} \quad (14)$$

- The only difference with iVAE is that  $\nu_\psi$  is a MLP taking in only  $z$ .

# Implicit Deep Latent Variable Models for Text Generation (EMNLP 2019)

Le Fang<sup>†</sup>, Chunyuan Li<sup>§</sup>, Jianfeng Gao<sup>§</sup>, Wen Dong<sup>†</sup>, Changyou Chen<sup>†</sup>

## Results:

Methods	-ELBO↓	PPL↓	KL↑	MI↑	AU↑
Dataset: PTB					
VAE	102.6	108.26	1.08	0.8	2
$\beta(0.5)$ -VAE	104.5	117.92	<b>7.50</b>	3.1	5
SA-VAE	102.6	107.71	1.23	0.7	2
Cyc-VAE	103.1	110.50	3.48	1.8	5
iVAE	<b>87.6</b>	<b>54.46</b>	6.32	<b>3.5</b>	<b>32</b>
iVAE <sub>MI</sub>	<b>87.2</b>	<b>53.44</b>	<b>12.51</b>	<b>12.2</b>	<b>32</b>
Dataset: Yahoo					
VAE	328.6	61.21	0.0	0.0	0
$\beta(0.4)$ -VAE	328.7	61.29	6.3	2.8	8
SA-VAE	327.2	60.15	5.2	2.7	10
Lag-VAE	326.7	59.77	5.7	2.9	15
iVAE	<b>309.5</b>	<b>48.22</b>	<b>8.0</b>	<b>4.4</b>	<b>32</b>
iVAE <sub>MI</sub>	<b>309.1</b>	<b>47.93</b>	<b>11.4</b>	<b>10.7</b>	<b>32</b>
Dataset: Yelp					
VAE	357.9	40.56	0.0	0.0	0
$\beta(0.4)$ -VAE	358.2	40.69	4.2	2.0	4
SA-VAE	355.9	39.73	2.8	1.7	8
Lag-VAE	355.9	39.73	3.8	2.4	11
iVAE	<b>348.2</b>	<b>36.70</b>	<b>7.6</b>	<b>4.6</b>	<b>32</b>
iVAE <sub>MI</sub>	<b>348.7</b>	<b>36.88</b>	<b>11.6</b>	<b>11.0</b>	<b>32</b>

# FlowPrior: Learning Expressive Priors for Latent Variable Sentence Models (NAACL 2021)

Xiaoan Ding, Kevin Gimpel

## Core ideas:

- Using normalizing flows (NF) for the prior distribution. In this paper, using *real-valued non-volume preserving* (real NVP) transformations ([Dinh et al., 2016](#))
  - $z_L = f_L \circ f_{L-1} \circ \dots \circ f_1(z_0)$ ,  $z_0 \sim N(0, I)$ ,  $z_L$  is the sentence latent variable.

$$z_0 = f_1^{-1} \circ \dots \circ f_{L-1}^{-1} \circ f_L^{-1}(z_L) \quad (5)$$

$$\log p_\psi(z_L) = \log p_0(z_0) - \sum_{l=1}^L \log |\det\left(\frac{\partial f_l(z_{l-1})}{\partial z_{l-1}}\right)| \quad (6)$$

- Inspired by Real-NVP, the choice of  $f$  is an affine transformation, which is  $f: z_{i-1} \rightarrow z_i$

$$\begin{aligned} z_i^{(1:d)} &= z_{i-1}^{(1:d)} \\ z_i^{(d+1:D)} &= z_i^{(d+1:D)} \odot \exp(s(z_{i-1}^{(1:d)})) + t(z_{i-1}^{(1:d)}) \end{aligned}$$

- 1) easily invertible
- 2) its Jacobian determinant is easy to compute. (Do not require inverse/Jacobian of  $s$  and  $t$ )
- So we can model  $s$  and  $t$  using NNs, denoted by  $p_\psi(z)$

$$\mathbf{J} = \begin{bmatrix} \mathbb{I}_d & \mathbf{0}_{d \times (D-d)} \\ \frac{\partial \mathbf{y}_{d+1:D}}{\partial \mathbf{x}_{1:d}} & \text{diag}(\exp(s(\mathbf{x}_{1:d}))) \end{bmatrix}$$

# FlowPrior: Learning Expressive Priors for Latent Variable Sentence Models (NAACL 2021)

Xiaoan Ding, Kevin Gimpel

## Core ideas:

- New Objective: Using importance weighting + Monte Carlo for KL

$$\begin{aligned}\mathcal{L}(\theta, \phi, \psi; x) &= \log \frac{1}{N} \sum_{i=1}^N \frac{p_\theta(x|z^{(i)})p_\psi(z^{(i)})}{q_\phi(z^{(i)}|x)} \\ &+ \frac{1}{N} \sum_{i=1}^N \log p_\theta(x|z^{(i)}) - \text{KL}_{\phi,\psi}(x, \{z^{(i)}\}_{i=1}^N) \\ \text{s.t. } z^{(i)} &\sim q_\phi(z|x) \quad (7)\end{aligned}$$

# FlowPrior: Learning Expressive Priors for Latent Variable Sentence Models (NAACL 2021)

Xiaoan Ding, Kevin Gimpel

## Core ideas:

- Training Scheme

1. Draw  $N$  samples  $z_L^{(1)}, z_L^{(2)}, \dots, z_L^{(N)}$  from the inference network using the reparameterization trick.
2. Perform the inverse transformation to get the image of each point under the base distribution:  
 $z_0^{(1)}, z_0^{(2)}, \dots, z_0^{(N)}$ .
3. Compute the exact log likelihood of the sample prior with change of variable theorem (Eq. 6).
4. Compute and backpropagate the loss (Eq. 7).

# FlowPrior: Learning Expressive Priors for Latent Variable Sentence Models (NAACL 2021)

Xiaoan Ding, Kevin Gimpel

## Results:

Model	PPL( $\downarrow$ )	Recon( $\downarrow$ )	KL	AU( $\uparrow$ )	MI( $\uparrow$ )
VAE	101.40	101.28	0.00	0	0.00
Cyc-VAE	107.73	101.17	2.01	5	1.24
Lag-VAE	100.25	100.41	1.04	3	0.79
VAE + FB	101.56	99.84	4.46	32	0.90
Pre-VAE + FB	96.35	94.52	8.15	32	6.30
MoG-VAE	98.22	100.54	0.00	0	0.00
MoG-VAE + FB	97.50	99.44	2.35	32	0.68
Vamp-VAE	98.27	100.56	0.00	0	0.00
Vamp-VAE + FB	97.83	99.53	2.31	32	0.72
FlowPrior	94.72	98.46	3.28	2	2.25
FlowPrior + FB	93.58	99.20	7.21	31	2.83

Table 1: Language modeling results on PTB dataset.

# Adversarial Poincaré Variational Autoencoder (APo-VAE) (NAACL 2021)

Shuyang Dai · Zhe Gan · Yu Cheng · Chenyang Tao · Lawrence Carin · Jingjing Liu

## Core ideas:

- Natural languages have latent hierarchy, but prior distribution from Euclidean space couldn't capture it, so resorting to Riemannian Geometry.
- Proposing a prior based on Poincaré Ball Model in Hyperbolic space, with the following advantages:

$$\mathbb{B}_c^n := \{z \in \mathbb{R}^n \mid c\|z\|^2 < 1\}$$

- It implies significant representation ability.
- It's a well defined compact metric space, with well defined algebraic operations that allows back-propagation.

$$z \oplus_c z' := \frac{(1 + 2c\langle z, z' \rangle + c\|z'\|^2)z + (1 - c\|z\|^2)z'}{1 + 2c\langle z, z' \rangle + c^2\|z\|^2\|z'\|^2}. \quad (4)$$

$$\exp_{\mu}^c(u) := \mu \oplus_c (\tanh(\sqrt{c}\frac{\lambda_{\mu}^c\|u\|}{2})\frac{u}{\sqrt{c}\|u\|}),$$

$$\log_{\mu}^c(y) := \frac{2}{\sqrt{c}\lambda_{\mu}^c} \tanh^{-1}(\sqrt{c}\|\kappa_{\mu,y}\|) \frac{\kappa_{\mu,y}}{\|\kappa_{\mu,y}\|}, \quad (5)$$

$$P_{0 \rightarrow \mu}^c(v) = \log_{\mu}^c(\mu \oplus_c \exp_0^c(v)) = \frac{\lambda_0^c}{\lambda_{\mu}^c} v. \quad (6)$$

where  $\kappa_{\mu,y} := (-\mu) \oplus_c y$ .

# Adversarial Poincaré Variational Autoencoder (APo-VAE) (NAACL 2021)

Shuyang Dai · Zhe Gan · Yu Cheng · Chenyang Tao · Lawrence Carin · Jingjing Liu

## Core ideas:

- How does the defined compact metric space contribute to the formulation of a prior?
  - Poincaré ball based normal prior  $\mathbf{z} \sim N_{B^n_c}(\mu, \Sigma)$ :

$$\mathbf{z} = \exp_{\boldsymbol{\mu}}^c \left( \frac{\lambda_0^c}{\lambda_{\boldsymbol{\mu}}^c} \mathbf{v} \right), \mathbf{v} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}). \quad (7)$$

- Inspired by iVAE, we can enhance  $\mathbf{v} := G(\mathbf{x}, \xi; \phi_1)$  and  $\boldsymbol{\mu} := F(\mathbf{x}; \phi_2)$ , as outputs of encoder ( $\phi$ ). Then we get  $\mathbf{z}$  from the above formula.
- Same as iVAE, they use dual form to evaluate  $KL(q_\phi(z|x) \parallel p(z))$ :

$$\begin{aligned} D_{KL}(q_\phi(z|x) \parallel p(z)) &= \max_{\psi} && (9) \\ &\left\{ \mathbb{E}_{\mathbf{z} \sim q_\phi(z|x)} \nu_\psi(\mathbf{x}, \mathbf{z}) - \mathbb{E}_{\mathbf{z} \sim p(z)} \exp \nu_\psi(\mathbf{x}, \mathbf{z}) \right\}, \end{aligned}$$

- Same as iVAE,  $\nu_\psi$  is optimized by the following:

$$\begin{aligned} \mathcal{L}_1 &= \mathbb{E}_{\mathbf{x} \sim \mathbf{X}} [\mathbb{E}_{\mathbf{z} \sim q_\phi(z|x)} \nu_\psi(\mathbf{x}, \mathbf{z}) \\ &\quad - \mathbb{E}_{\mathbf{z} \sim p(z)} \exp \nu_\psi(\mathbf{x}, \mathbf{z})], \quad (10) \end{aligned}$$

# Adversarial Poincaré Variational Autoencoder (APo-VAE) (NAACL 2021)

Shuyang Dai · Zhe Gan · Yu Cheng · Chenyang Tao · Lawrence Carin · Jingjing Liu

## Core ideas:

- Different from iVAE, here the real prior is not assumed Gaussian. It's estimated by sampling scheme used in VampPrior (Tomczak and Welling, 2018)

$$p_{\delta}(z) = \frac{1}{K} \sum_{k=1}^K q_{\phi}(z|s_k), \quad (12)$$

- $\delta := \{s_k\}_{k=1}^K$  is now learnable pseudo inputs. Replacing  $p(z)$  with  $p_{\delta}(z)$  seeks to match the aggregated posterior  $q(z) = \sum q_{\phi}(z|x_i) / N$ .

# Adversarial Poincaré Variational Autoencoder (APo-VAE) (NAACL 2021)

Shuyang Dai · Zhe Gan · Yu Cheng · Chenyang Tao · Lawrence Carin · Jingjing Liu

## Core ideas:

- For the geometry-aware decoder, they use a deterministic (and learnable) hyperbolic linear function  $f$  to extract feature, then pass to LSTM decoder. ( $a$ ,  $b$ , together with the LSTM, are trainable parameters of decoder  $\theta$ )

$$f_{\mathbf{a}, \mathbf{b}}^c(\mathbf{z}) = \text{sign}(\langle \mathbf{a}, \log_b^c(\mathbf{z}) \rangle_b) \|\mathbf{a}\|_b d_c^{\mathbb{B}}(\mathbf{z}, H_{\mathbf{a}, \mathbf{b}}^c), \quad (8)$$

where  $H_{\mathbf{a}, \mathbf{b}}^c = \{\mathbf{z} \in \mathbb{B}_c^n | \langle \mathbf{a}, \log_b^c(\mathbf{z}) \rangle_b = 0\}$ ,

$$d_c^{\mathbb{B}}(\mathbf{z}, H_{\mathbf{a}, \mathbf{b}}^c) = \frac{1}{\sqrt{c}} \sinh^{-1} \left( \frac{2\sqrt{c} |\langle \kappa_{\mathbf{b}, \mathbf{z}}, \mathbf{a} \rangle|}{(1 - c \|\kappa_{\mathbf{b}, \mathbf{z}}\|^2) \|\mathbf{a}\|} \right)$$

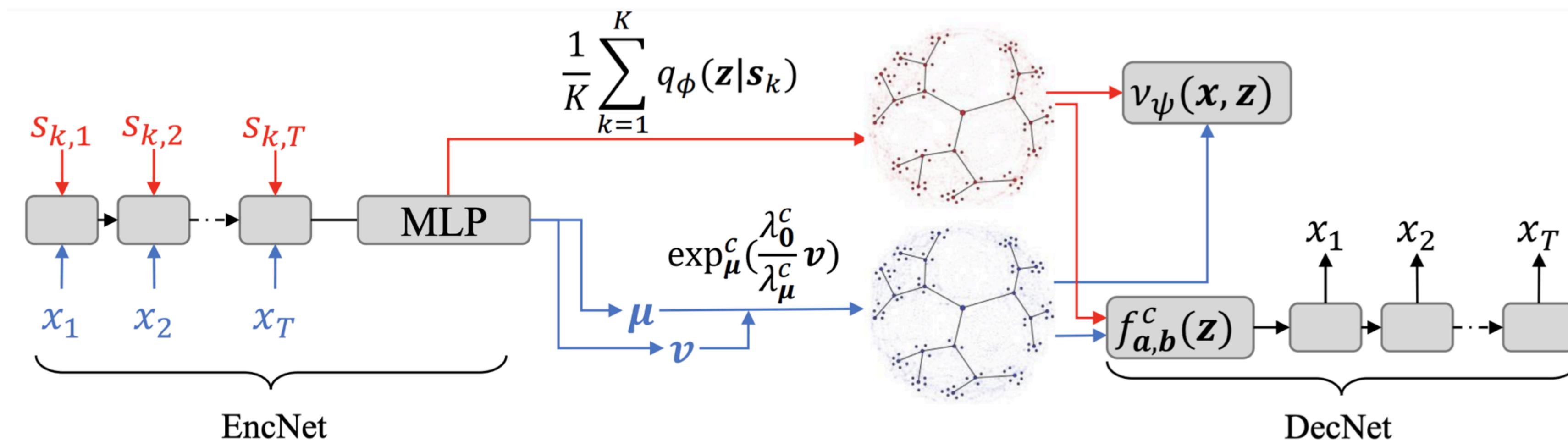
- $\theta$  and  $\phi$  are optimized by the following objective

$$\begin{aligned} \mathcal{L}_2 = \mathbb{E}_{\mathbf{x} \sim \mathbf{X}} \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z} | \mathbf{x})} [ & \log p_{\theta}(\mathbf{x} | \mathbf{z}) \\ & - \nu_{\psi}(\mathbf{x}, \mathbf{z}) ]. \end{aligned} \quad (11)$$

# Adversarial Poincaré Variational Autoencoder (APo-VAE) (NAACL 2021)

Shuyang Dai<sup>1</sup> Zhe Gan<sup>2</sup> Yu Cheng<sup>3</sup> Chenyang Tao<sup>1</sup> Lawrence Carin<sup>1</sup> Jingjing Liu<sup>1</sup>

**Core ideas:**



# Adversarial Poincaré Variational Autoencoder (APo-VAE) (NAACL 2021)

Shuyang Dai<sup>1</sup> Zhe Gan<sup>2</sup> Yu Cheng<sup>3</sup> Chenyang Tao<sup>4</sup> Lawrence Carin<sup>5</sup> Jingjing Liu<sup>1</sup>

## Core ideas:

- Training Procedure

---

### Algorithm 1 Training procedure of APo-VAE.

---

- 1: **Input:** Data samples  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N$ , Poincaré curvature  $c$ , and number of pseudo-input  $K$ .
- 2: Initialize  $\theta$ ,  $\phi$ ,  $\psi$ , and  $\delta$ .
- 3: **for**  $iter$  from 1 to  $max\_iter$  **do**
- 4:   Sample a mini-batch  $\{\mathbf{x}_m\}_{m=1}^M$  from  $\mathbf{X}$  of size  $M$ .
- 5:   **# Sampling in the Hyperbolic Space.**
- 6:   Obtain  $\mu_m$  and  $v_m$  from  $\text{EncNet}_\phi(\mathbf{x}_m)$ .
- 7:   Move  $v_m$  to  $u_m = P_{0 \rightarrow \mu_m}^c(v_m)$  by (6).
- 8:   Map  $u_m$  to  $z_m = \exp_{\mu_m}^c(u_m)$  by (5).
- 9:   **# Update the dual function and the pseudo-input.**
- 10:   Sample  $\tilde{z}_m$  by (12).
- 11:   Update  $\psi$  and  $\delta$  by gradient ascent on (10)
- 12:   **# Update the encoder and decoder networks.**
- 13:   Update  $\theta$  and  $\phi$  by gradient ascent on (11).
- 14: **end for**

# Adversarial Poincaré Variational Autoencoder (APo-VAE) (NAACL 2021)

Shuyang Dai<sup>1</sup> Zhe Gan<sup>2</sup> Yu Cheng<sup>3</sup> Chenyang Tao<sup>4</sup> Lawrence Carin<sup>5</sup> Jingjing Liu<sup>6</sup>

## Results:

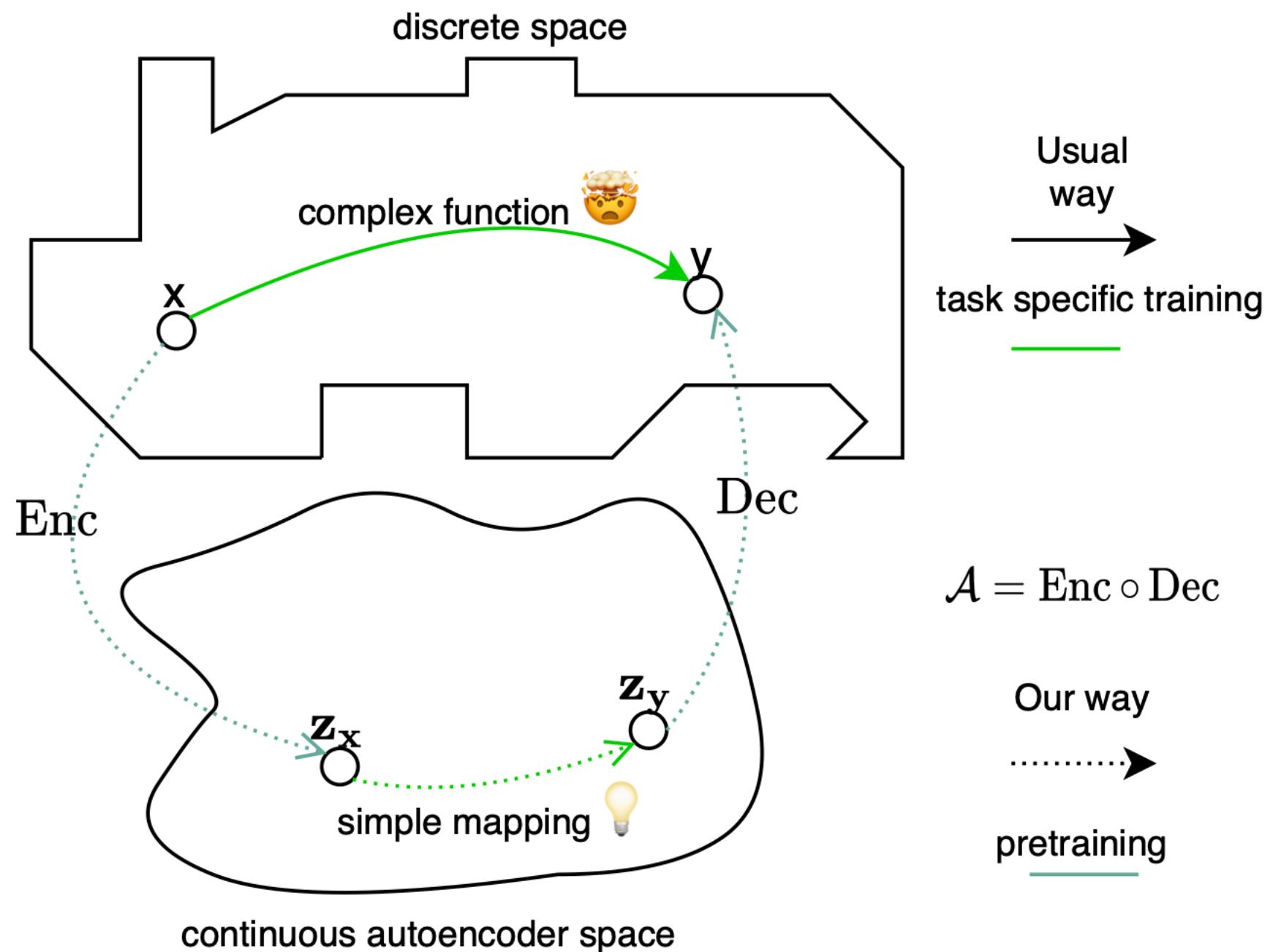
Model	-ELBO	PPL	KL	MI	AU
	PTB				
VAE	102.6	108.26	1.1	0.8	2
$\beta$ -VAE	104.5	117.92	7.5	3.1	5
SA-VAE	102.6	107.71	1.2	0.7	2
vMF-VAE	95.8	93.70	2.9	3.2	21
$\mathcal{P}$ -VAE	91.4	76.13	4.5	2.9	23
iVAE	87.2	53.44	<b>12.5</b>	<b>12.2</b>	<b>32</b>
APo-VAE	87.2	53.32	8.4	4.8	<b>32</b>
APo-VAE+VP	<b>87.0</b>	<b>53.02</b>	8.9	4.5	<b>32</b>
Yahoo					
VAE	328.6	61.21	0.0	0.0	0
$\beta$ -VAE	328.7	61.29	6.3	2.8	8
SA-VAE	327.2	60.15	5.2	2.9	10
LAG-VAE	326.7	59.77	5.7	2.9	15
vMF-VAE	318.5	53.92	6.3	3.7	23
$\mathcal{P}$ -VAE	313.4	50.57	7.2	3.3	27
iVAE	309.1	47.93	<b>11.4</b>	<b>10.7</b>	<b>32</b>
APo-VAE	286.2	47.00	6.9	4.1	<b>32</b>
APo-VAE+VP	<b>285.6</b>	<b>46.61</b>	8.1	4.9	<b>32</b>
Yelp					
VAE	357.9	40.56	0.0	0.0	0
$\beta$ -VAE	358.2	40.69	4.2	2.0	4
SA-VAE	357.8	40.51	2.8	1.7	8
LAG-VAE	355.9	39.73	3.8	2.4	11
vMF-VAE	356.2	51.03	4.1	3.9	13
$\mathcal{P}$ -VAE	355.4	50.64	4.3	4.8	19
iVAE	348.7	36.88	11.6	<b>11.0</b>	<b>32</b>
APo-VAE	319.7	34.10	12.1	7.5	<b>32</b>
APo-VAE+VP	<b>316.4</b>	<b>32.91</b>	<b>12.7</b>	6.2	<b>32</b>

# Plug and Play Autoencoders for Conditional Text Generation

Florian Mai et al

## Core ideas:

- Learning mappings from continuous space is easier than from discrete space.

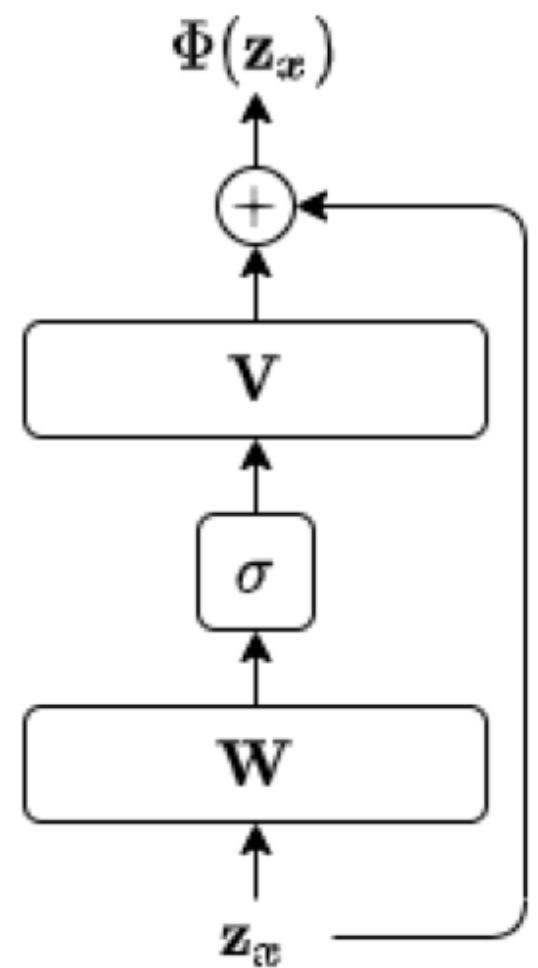


# Plug and Play Autoencoders for Conditional Text Generation

Florian Mai et al

## Core ideas:

- Encoder and Decoder are pretrained and fixed during training.
- Adding a mapping function  $\Phi$ , parametrized by a OffsetNet
  - Different from DAAE, where operations in latent space is just adding and subtracting. Here the operation is conditioned on the input, and during inference, using FGIM to fulfill style transfer.



(b) OffsetNet

# Plug and Play Autoencoders for Conditional Text Generation

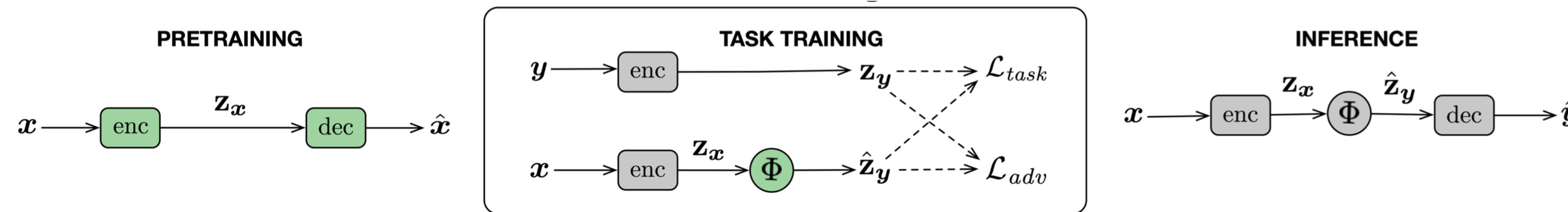
Florian Mai et al

## Core ideas:

- Adversarial learning to ensure the manifold of the mapping  $\Phi$  resembles the original latent space.

$$\mathcal{L}_{adv}(\Phi(\mathbf{z}_{\mathbf{x}_i}); \boldsymbol{\theta}) = -\log(\text{disc}(\Phi(\mathbf{z}_{\mathbf{x}_i}); \boldsymbol{\theta})) \quad (11)$$

- For supervised tasks:



$$\mathcal{L} = \mathcal{L}_{task} + \lambda_{adv} \mathcal{L}_{adv}. \quad (1)$$

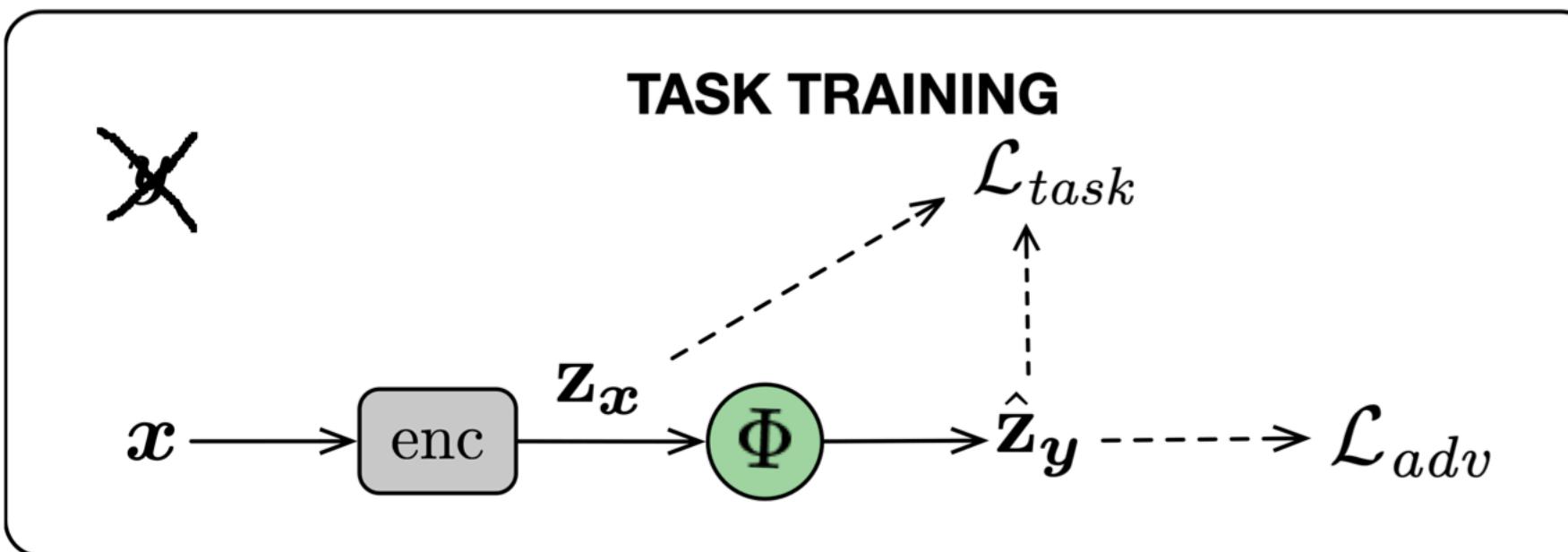
$$\mathcal{L}_{task} = \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{emb}(\Phi(\mathbf{z}_{\mathbf{x}_i}; \boldsymbol{\theta}), \mathbf{z}_{\mathbf{y}_i}). \quad (8)$$

# Plug and Play Autoencoders for Conditional Text Generation

Florian Mai et al

## Core ideas:

- For unsupervised tasks:



$$\mathcal{L}_{task}(\hat{\mathbf{z}}_y, \mathbf{z}_x) = \lambda_{sty} \mathcal{L}_{sty}(\hat{\mathbf{z}}_y) + (1 - \lambda_{sty}) \mathcal{L}_{cont}(\hat{\mathbf{z}}_y, \mathbf{z}_x)$$

$$\mathcal{L}_{sty}(\Phi(\mathbf{z}_{x_i}; \theta), \mathbf{z}_{x_i}) = -\log(c(\Phi(\mathbf{z}_{x_i}; \theta))).$$

Here  $c$  is a pretrained classifier for style. They freeze  $c$  when training and encourage  $\Phi$  to produce outputs of the target attribute ( $y=1$ ).

# Plug and Play Autoencoders for Conditional Text Generation

Florian Mai et al

## Core ideas:

- For unsupervised tasks, how to better manipulate latent space (i.e to invert sentiment)? Use FGIM until the confidence of the sentiment classifier is greater than a threshold.

$$\hat{\hat{\mathbf{z}}}_{\mathbf{x}_i} = \hat{\mathbf{z}}_{\mathbf{x}_i} + \omega \nabla_{\hat{\mathbf{z}}_{\mathbf{x}_i}} \mathcal{L}(\hat{\mathbf{z}}_{\mathbf{x}_i}, \mathbf{z}_{\mathbf{x_i}}),$$

# Plug and Play Autoencoders for Conditional Text Generation

Florian Mai et al

## Results:

Model	BLEU	SARI	Time
S2S-Scratch	3.6	15.6	3.7×
S2S-Pretrain	5.4	16.2	3.7×
S2S-MLP	10.5	17.7	3.7×
S2S-Freeze	23.3	22.4	2.2×
Emb2Emb	<b>34.7</b>	<b>25.4</b>	1.0×

Table 1: Text simplification performance of model variants of end2end training on the test set. “Time” is wall time of one training epoch, relative to our model, Emb2Emb.

Model	Acc.	s-BLEU	+Time
Shen et al.	96.8	6.5	0.5×
FGIM	94.9	10.8	70.0×
Emb2Emb + FGIM	93.1	18.1	2820.0×
Emb2Emb	87.1	22.1	1.0×

Table 2: Self-BLEU (“s-BLEU”) on the Yelp sentiment transfer test set for the configurations in Figure 5 with highest transfer accuracy (“Acc.”). “+Time” reports the inference-time slowdown factor due to each model’s additional computation (relative to our method).

# Wasserstein Auto-Encoders

Ilya Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Scho'lkopf

## Core ideas:

- Minimizing  $W_c(P_X, P_G)$  between the true (but unknown) data distribution  $P_X$  and a latent variable model  $P_G$

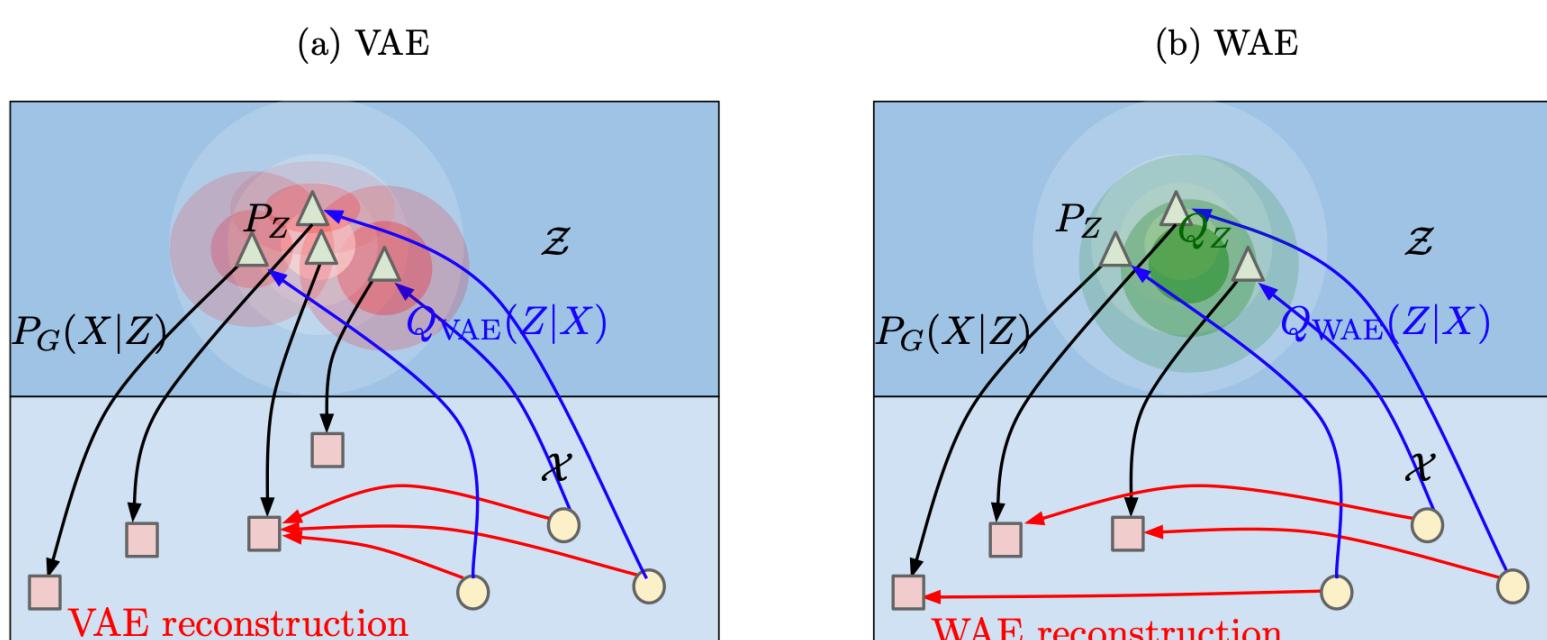
$$W_c(P_X, P_G) := \inf_{\Gamma \in \mathcal{P}(X \sim P_X, Y \sim P_G)} \mathbb{E}_{(X, Y) \sim \Gamma} [c(X, Y)],$$

**Theorem 1.** For  $P_G$  as defined above with deterministic  $P_G(X|Z)$  and any function  $G: \mathcal{Z} \rightarrow \mathcal{X}$

$$\inf_{\Gamma \in \mathcal{P}(X \sim P_X, Y \sim P_G)} \mathbb{E}_{(X, Y) \sim \Gamma} [c(X, Y)] = \inf_{Q: Q_Z = P_Z} \mathbb{E}_{P_X} \mathbb{E}_{Q(Z|X)} [c(X, G(Z))],$$

where  $Q_Z$  is the marginal distribution of  $Z$  when  $X \sim P_X$  and  $Z \sim Q(Z|X)$ .

- Enforcing aggregated posterior ( $Q_Z := \int Q(Z|X)dP_X$ ) to match the Prior  $P_z$



$$D_{\text{WAE}}(P_X, P_G) := \inf_{Q(Z|X) \in \mathcal{Q}} \mathbb{E}_{P_X} \mathbb{E}_{Q(Z|X)} [c(X, G(Z))] + \lambda \cdot \mathcal{D}_Z(Q_Z, P_Z),$$

# Wasserstein Auto-Encoders

Ilya Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Schölkopf

## Core ideas:

- Two proposed penalties for  $D_z$ 
  - GAN-based  $D_z$

$$\mathcal{D}_Z(Q_Z, P_Z) = D_{\text{JS}}(Q_Z, P_Z)$$

---

**Algorithm 1** Wasserstein Auto-Encoder  
with GAN-based penalty (WAE-GAN).

---

**Require:** Regularization coefficient  $\lambda > 0$ .

Initialize the parameters of the encoder  $Q_\phi$ ,  
decoder  $G_\theta$ , and latent discriminator  $D_\gamma$ .

**while**  $(\phi, \theta)$  not converged **do**

    Sample  $\{x_1, \dots, x_n\}$  from the training set

    Sample  $\{z_1, \dots, z_n\}$  from the prior  $P_Z$

    Sample  $\tilde{z}_i$  from  $Q_\phi(Z|x_i)$  for  $i = 1, \dots, n$

    Update  $D_\gamma$  by ascending:

$$\frac{\lambda}{n} \sum_{i=1}^n \log D_\gamma(z_i) + \log(1 - D_\gamma(\tilde{z}_i))$$

    Update  $Q_\phi$  and  $G_\theta$  by descending:

$$\frac{1}{n} \sum_{i=1}^n c(x_i, G_\theta(\tilde{z}_i)) - \lambda \cdot \log D_\gamma(\tilde{z}_i)$$

**end while**

# Wasserstein Auto-Encoders

Ilya Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Scho'lkopf

## Core ideas:

- Two proposed penalties for  $D_z$ 
  - MMD-based  $D_z$

$$\text{MMD}_k(P_Z, Q_Z) = \left\| \int_{\mathcal{Z}} k(z, \cdot) dP_Z(z) - \int_{\mathcal{Z}} k(z, \cdot) dQ_Z(z) \right\|_{\mathcal{H}_k},$$

---

**Algorithm 2** Wasserstein Auto-Encoder  
with MMD-based penalty (WAE-MMD).

---

**Require:** Regularization coefficient  $\lambda > 0$ ,

characteristic positive-definite kernel  $k$ .

Initialize the parameters of the encoder  $Q_\phi$ ,  
decoder  $G_\theta$ , and latent discriminator  $D_\gamma$ .

**while**  $(\phi, \theta)$  not converged **do**

    Sample  $\{x_1, \dots, x_n\}$  from the training set

    Sample  $\{z_1, \dots, z_n\}$  from the prior  $P_Z$

    Sample  $\tilde{z}_i$  from  $Q_\phi(Z|x_i)$  for  $i = 1, \dots, n$

    Update  $Q_\phi$  and  $G_\theta$  by descending:

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n c(x_i, G_\theta(\tilde{z}_i)) + \frac{\lambda}{n(n-1)} \sum_{\ell \neq j} k(z_\ell, z_j) \\ & + \frac{\lambda}{n(n-1)} \sum_{\ell \neq j} k(\tilde{z}_\ell, \tilde{z}_j) - \frac{2\lambda}{n^2} \sum_{\ell, j} k(z_\ell, \tilde{z}_j) \end{aligned}$$

**end while**

# Wasserstein Auto-Encoders

Ilya Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Scho'lkopf

## Core ideas:

- Difference with WGAN
  - Claiming the intractability of infimum, WGAN uses dual form of 1-Wasserstein Distance ( $W_1$ ), while WAE could use any cost function  $c$  ( $W_c$ )

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

$$W(\mathbb{P}_r, \mathbb{P}_\theta) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim \mathbb{P}_r} [f(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta} [f(x)]$$

---

**Algorithm 1** WGAN, our proposed algorithm. All experiments in the paper used the default values  $\alpha = 0.00005$ ,  $c = 0.01$ ,  $m = 64$ ,  $n_{\text{critic}} = 5$ .

---

**Require:** :  $\alpha$ , the learning rate.  $c$ , the clipping parameter.  $m$ , the batch size.  
 $n_{\text{critic}}$ , the number of iterations of the critic per generator iteration.

**Require:** :  $w_0$ , initial critic parameters.  $\theta_0$ , initial generator's parameters.

```
1: while  $\theta$  has not converged do
2:   for  $t = 0, \dots, n_{\text{critic}}$  do
3:     Sample  $\{x^{(i)}\}_{i=1}^m \sim \mathbb{P}_r$  a batch from the real data.
4:     Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
5:      $g_w \leftarrow \nabla_w [\frac{1}{m} \sum_{i=1}^m f_w(x^{(i)}) - \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)}))]$ 
6:      $w \leftarrow w + \alpha \cdot \text{RMSProp}(w, g_w)$ 
7:      $w \leftarrow \text{clip}(w, -c, c)$ 
8:   end for
9:   Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
10:   $g_\theta \leftarrow -\nabla_\theta \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)}))$ 
11:   $\theta \leftarrow \theta - \alpha \cdot \text{RMSProp}(\theta, g_\theta)$ 
12: end while
```

# Wasserstein Auto-Encoders

Ilya Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Schölkopf

## Core ideas:

- Relation to other models:
  - AAE: When  $c$  equals to L2 norm, WAE-GAN becomes AAE

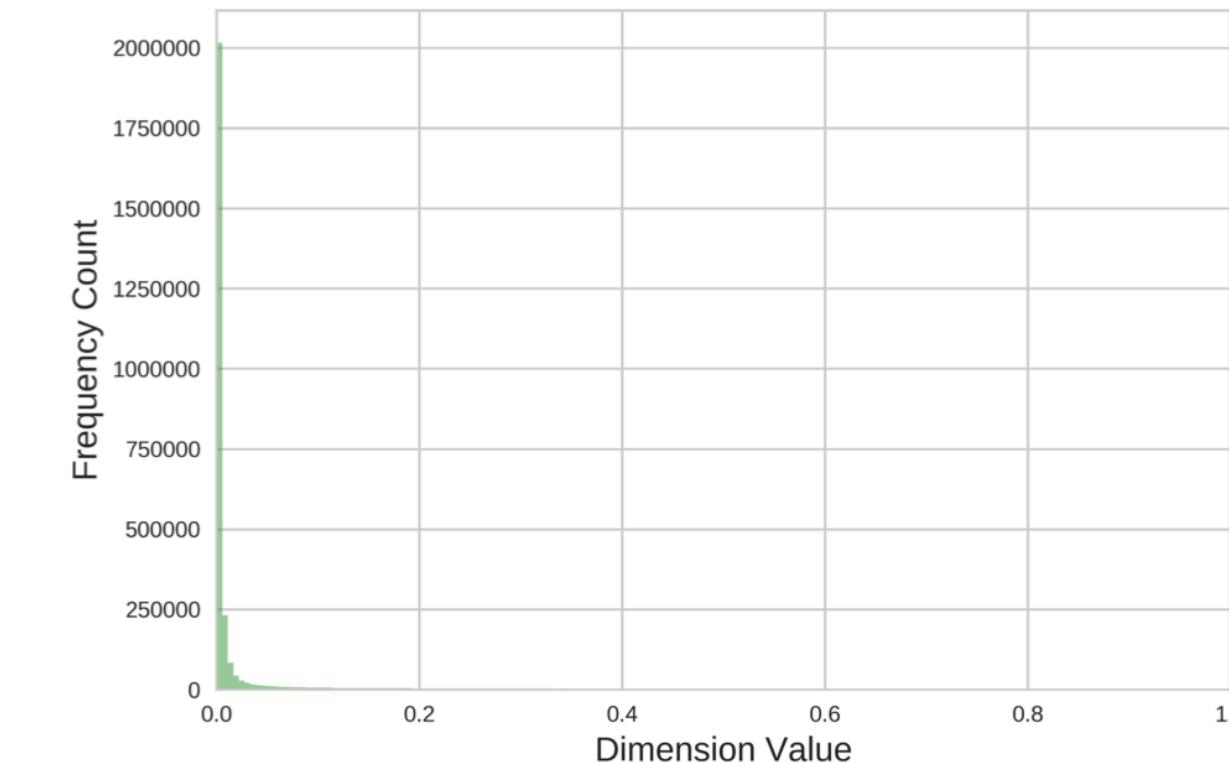
# Stochastic Wasserstein Autoencoder for Probabilistic Sentence Generation (NAACL 2019)

Hareesh Bahuleyan, Lili Mou, Hao Zhou, Olga Vechtomova

## Core ideas:

- Theoretically and empirically, the encoded Gaussian of the original WAE would be Dirac-delta ( $\sigma \rightarrow 0$ )

**Theorem 1.** Suppose we have a Gaussian family  $\mathcal{N}(\mu, \text{diag } \sigma^2)$ , where  $\mu$  and  $\sigma$  are parameters. The covariance is diagonal, meaning that the variables are independent. If the gradient of  $\sigma$  completely comes from sample gradient and  $\sigma$  is small at the beginning of training, then the Gaussian converges to a Dirac delta function with stochastic gradient descent, i.e.,  $\sigma \rightarrow 0$ . (See Appendix A for the proof.) □



(a)  $\lambda_{\text{KL}} = 0$

- Add a stochastic term to original WAE-MMD, yielding a relaxed optimization of WAE loss with a constraint on  $\sigma$

$$\begin{aligned} J &= J_{\text{rec}} + \lambda_{\text{WAE}} \cdot \widehat{\text{MMD}} \\ &+ \lambda_{\text{KL}} \sum_n \text{KL} \left( \mathcal{N}(\boldsymbol{\mu}_{\text{post}}^{(n)}, \text{diag}(\boldsymbol{\sigma}_{\text{post}}^{(n)})^2) \middle\| \mathcal{N}(\boldsymbol{\mu}_{\text{post}}^{(n)}, \mathbf{I}) \right) \end{aligned} \quad (5)$$

# Stochastic Wasserstein Autoencoder for Probabilistic Sentence Generation (NAACL 2019)

Hareesh Bahuleyan, Lili Mou, Hao Zhou, Olga Vechtomova

## Results:

- For generation, WAE keeps continuity and smoothness as VAE, while have a much higher reconstruction performance. WAE-S further encourages the stochasticity.

	<b>BLEU<math>\uparrow</math></b>	<b>PPL<math>\downarrow</math></b>	<b>UniKL<math>\downarrow</math></b>	<b>Entropy</b>	<b>AvgLen</b>
<b>Corpus</b>	-	-	-	$\rightarrow 5.65$	$\rightarrow 9.6$
<b>DAE</b>	<b>86.35</b>	146.2	0.178	6.23	11.0
<b>VAE (KL-annealed)</b>	43.18	79.4	0.081	5.04	8.8
<b>WAE-D</b> $\lambda_{\text{WAE}} = 3$	86.03	113.8	0.071	<b>5.59</b>	10.0
<b>WAE-D</b> $\lambda_{\text{WAE}} = 10$	84.29	104.9	0.073	5.57	9.9
<b>WAE-S</b> $\lambda_{\text{KL}} = 0.0$	75.66	115.2	0.069	<b>5.61</b>	9.9
<b>WAE-S</b> $\lambda_{\text{KL}} = 0.01$	82.01	84.9	<b>0.058</b>	5.26	<b>9.4</b>
<b>WAE-S</b> $\lambda_{\text{KL}} = 0.1$	47.63	<b>62.5</b>	0.150	4.65	8.7

Table 1: Results of SNLI-style sentence generation, where WAE is compared with DAE and VAE. **D** and **S** refer to the deterministic and stochastic encoders, respectively.  $\uparrow/\downarrow$ The larger/lower, the better. For **Entropy** and **AvgLen**, the closer to corpus statistics, the better (indicated by the  $\rightarrow$  arrow).