

Improving Variational Autoencoders with Density Gap-based Regularization

Jianfei Zhang^{1,2} Jun Bai^{1,2} Chenghua Lin³ Yanmeng Wang⁴ Wenge Rong^{1,2}

¹State Key Laboratory of Software Development Environment, Beihang University, China

²School of Computer Science and Engineering, Beihang University, China

³Department of Computer Science, University of Sheffield, United Kingdom

⁴Ping An Technology, China

{zhangjf, bai_jun, w.rong}@buaa.edu.cn

c.lin@sheffield.ac.uk, wangyanmeng219@pingan.com.cn

Catalogue

1. Background:

- a. Variational Autoencoder —the design and theory of VAEs
- b. Posterior Collapse and Hole Problem —the two problems in VAEs we intend to solve
- c. Existing methods —existing methods to solve the two problems

2. Methodology

- a. Regularization on the aggregated posterior distribution —the theoretical support and previous methods
- b. Density Gap-based regularization —the proposed PDF-oriented regularization method
- c. Marginal regularization for more Mutual Information —regularization over marginal distributions
- d. Aggregation size for ablation

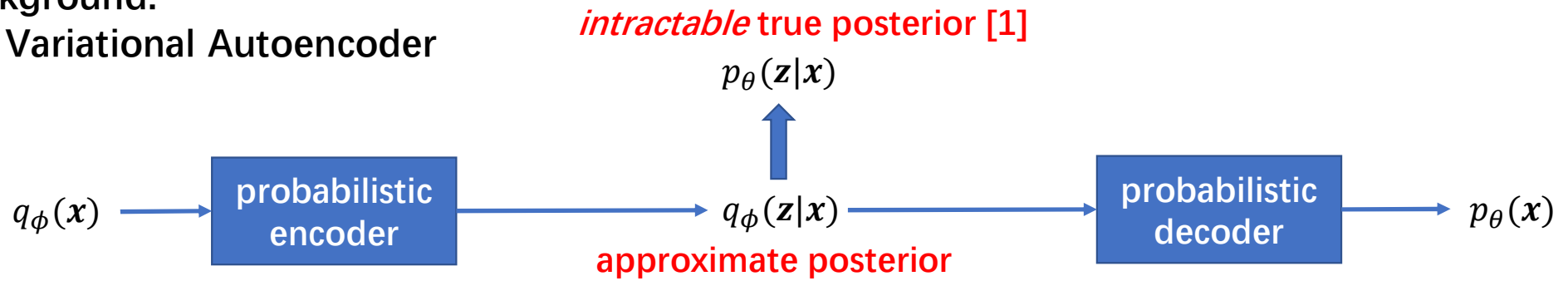
3. Experiment

- a. Language modeling
- b. Visualization of the posterior
- c. Interpolation study

Improving Variational Autoencoders with Density Gap-based Regularization

1. Background:

a. Variational Autoencoder



$$\begin{aligned}\mathcal{L}_{ELBO}(\theta, \phi, \mathbf{x}) &= \underbrace{E_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})]}_{\text{conditional log-likelihood}} - \underbrace{D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}))}_{\text{KL regularization}} \\ &= E_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x}) + \log p_\theta(\mathbf{z})] \\ &= E_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})] \\ &= E_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{z}|\mathbf{x}) - \log q_\phi(\mathbf{z}|\mathbf{x}) + \log p_\theta(\mathbf{x})] \\ &= E_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x})] - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x})) \\ &\leq E_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x})] \\ &= \log p_\theta(\mathbf{x})\end{aligned}$$

where,

$q_\phi(\mathbf{x})$: the data distribution, described by the dataset and received by the encoder ϕ

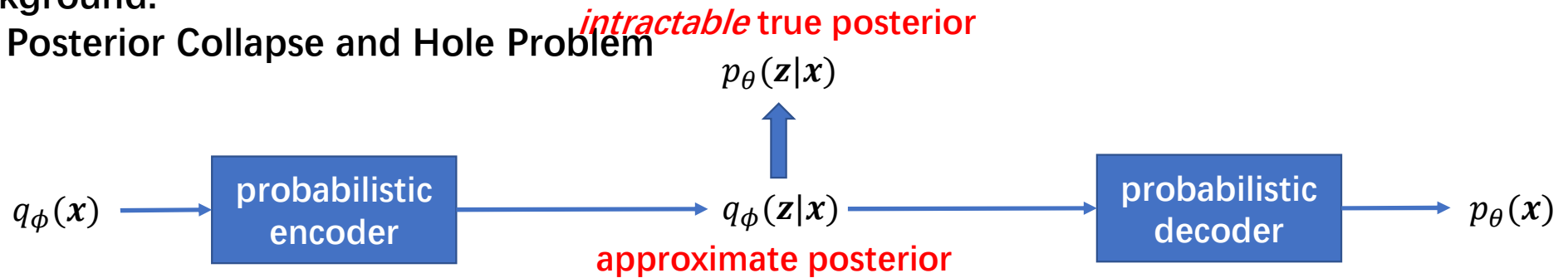
$p_\theta(\mathbf{z})$: the prior distribution of latent variable \mathbf{z} in decoder θ

$p_\theta(\mathbf{x})$: the generative data distribution by decoder θ (or the generative likelihood)

Improving Variational Autoencoders with Density Gap-based Regularization

1. Background:

b. Posterior Collapse and Hole Problem



$$\mathcal{L}_{ELBO}(\theta, \phi, x) = E_{q_\phi(\mathbf{z}|x)}[\log p_\theta(x|\mathbf{z})] - D_{KL}(q_\phi(\mathbf{z}|x)||p_\theta(\mathbf{z}))$$

Posterior Collapse:

$$\forall x D_{KL}(q_\phi(\mathbf{z}|x)||p_\theta(\mathbf{z})) \approx 0$$

$$\rightarrow \forall x p_\theta(\mathbf{z}|x) \approx q_\phi(\mathbf{z}|x) \approx p_\theta(\mathbf{z})$$

i.e., the latent variable \mathbf{z} contains **little information** of x

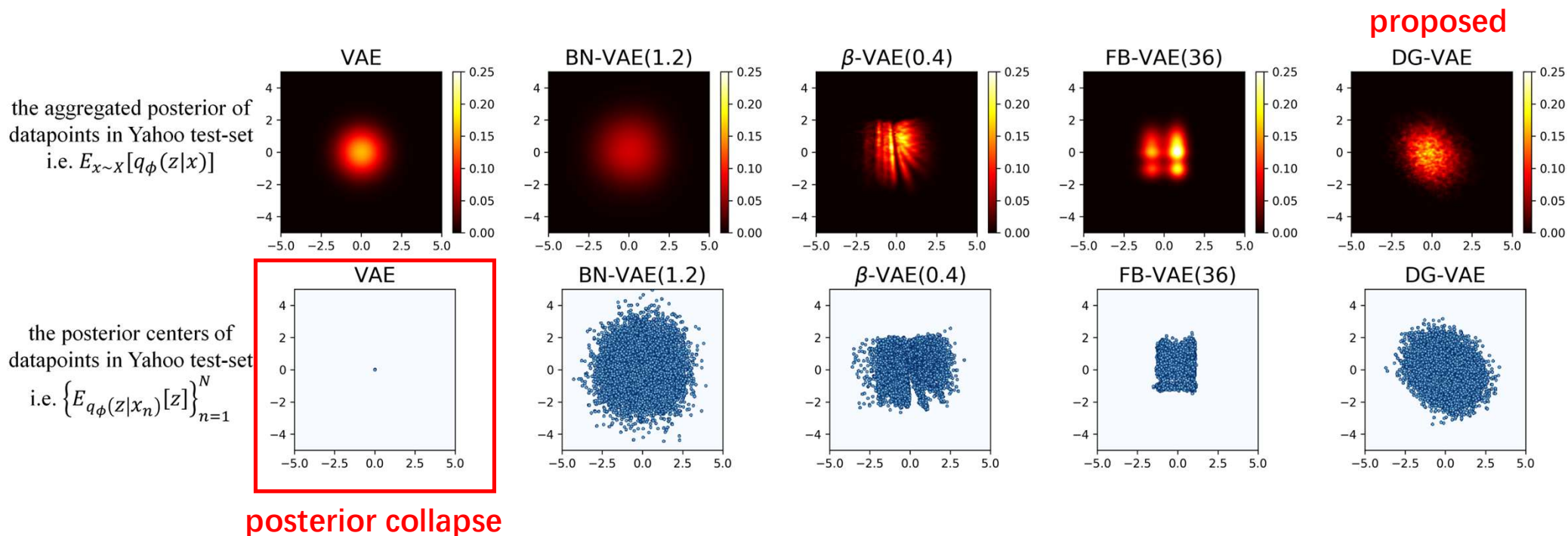
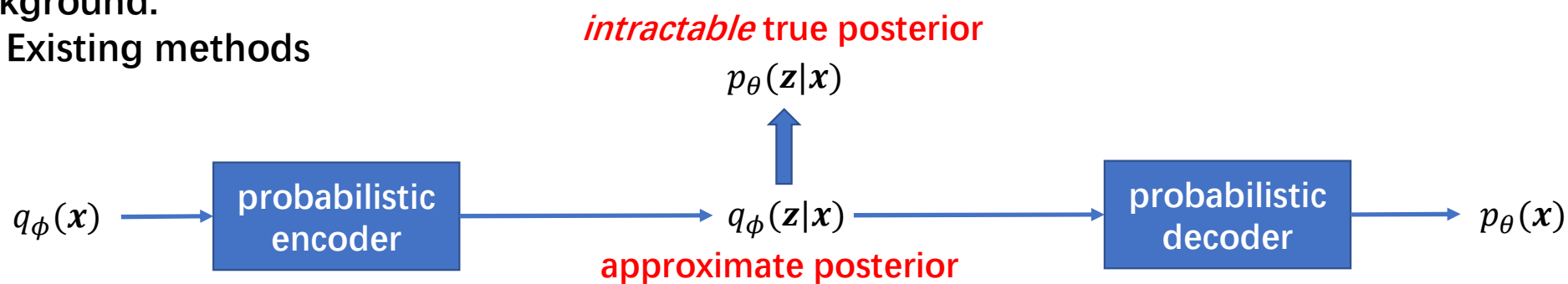
$$\rightarrow \forall x p_\theta(x|\mathbf{z}) = \frac{p_\theta(x,\mathbf{z})}{p_\theta(\mathbf{z})} \approx \frac{p_\theta(x,\mathbf{z})}{p_\theta(\mathbf{z}|x)} = p_\theta(x)$$

i.e., the decoder θ becomes **insensitive** to \mathbf{z}

i.e., the decoder **degenerates** to an unconditional language model (for NLG)

Improving Variational Autoencoders with Density Gap-based Regularization

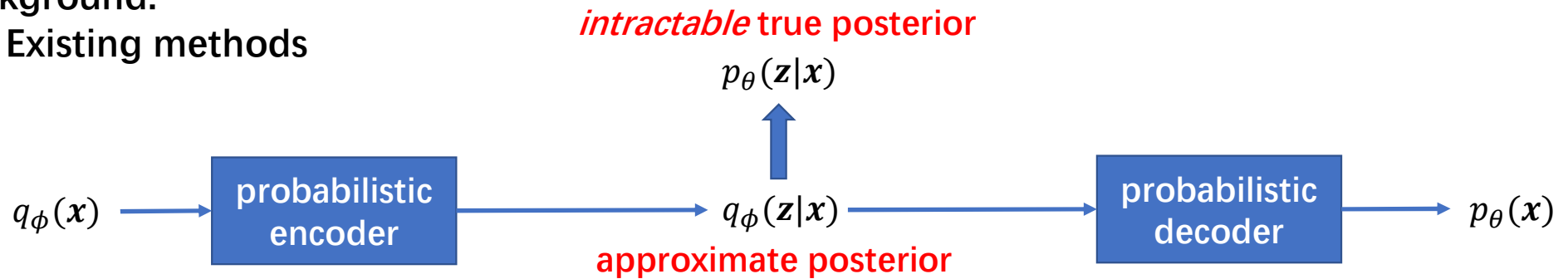
1. Background:
c. Existing methods



Improving Variational Autoencoders with Density Gap-based Regularization

1. Background:

c. Existing methods



$$\mathcal{L}_{ELBO}(\theta, \phi, \mathbf{x}) = E_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}))$$

Posterior Collapse:

$$\forall \mathbf{x} D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z})) \approx 0$$

→ training strategy:

Cyclic-VAEs (cyclic annealing schedule); AE pretraining;

→ semantic learning of \mathbf{z} :

Skip-VAE (skip connection on \mathbf{z}); BOW-VAEs (Bag-of-Word loss term on \mathbf{z});

→ hard restriction on $q_\phi(\mathbf{z}|\mathbf{x})$:

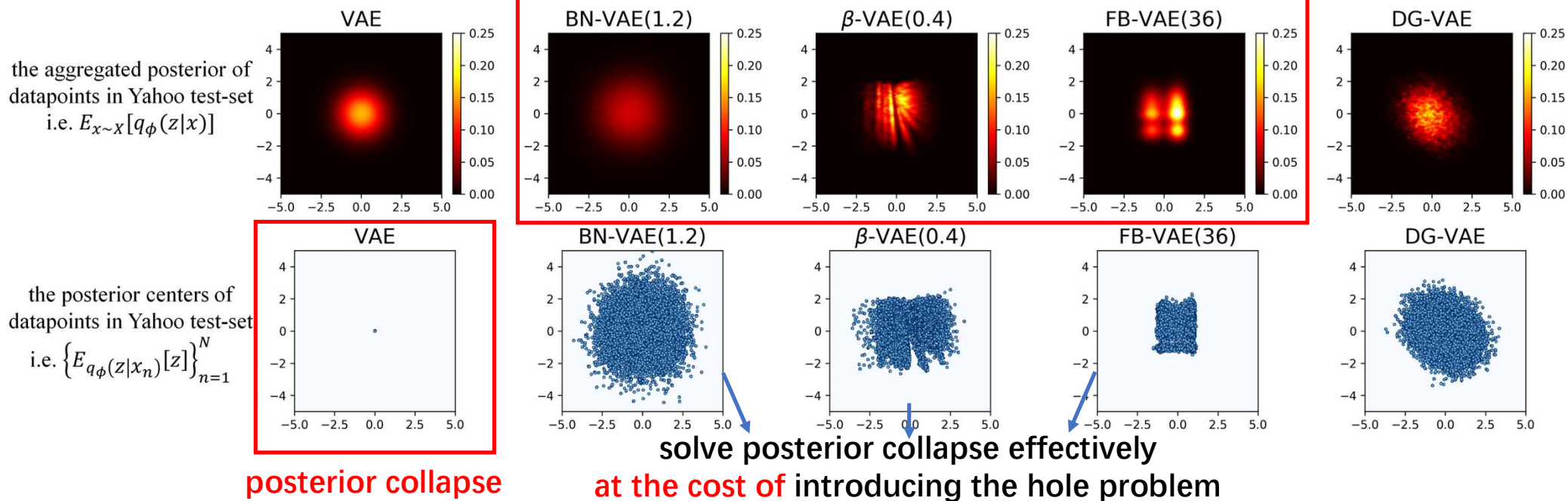
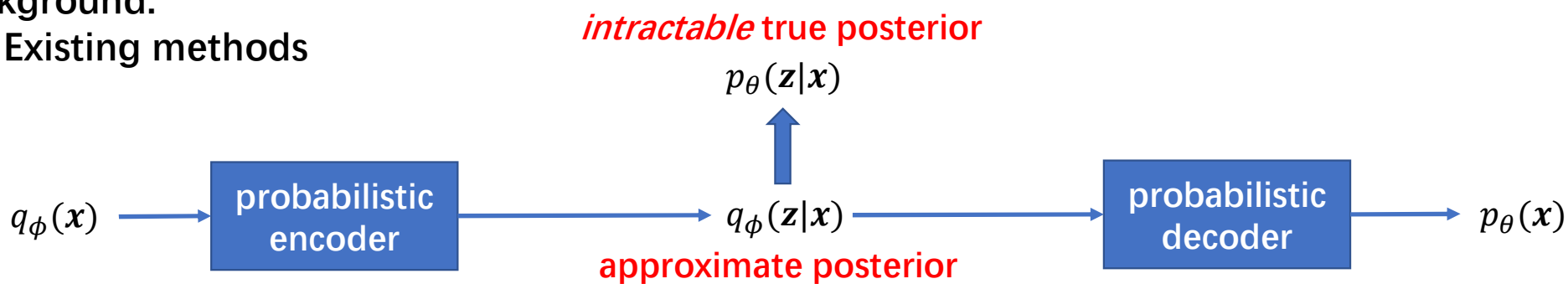
BN-VAEs (BN layer on $q_\phi(\mathbf{z}|\mathbf{x})$); vMF-VAEs (vMF distributions for $q_\phi(\mathbf{z}|\mathbf{x})$ and $p_\theta(\mathbf{z})$);

→ weakening $D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}))$ in $\mathcal{L}_{ELBO}(\theta, \phi, \mathbf{x})$:

β -VAEs (smaller weight of $D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}))$ in $\mathcal{L}_{ELBO}(\theta, \phi, \mathbf{x})$);

FB-VAEs (hinge loss of $D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}))$ in $\mathcal{L}_{ELBO}(\theta, \phi, \mathbf{x})$);

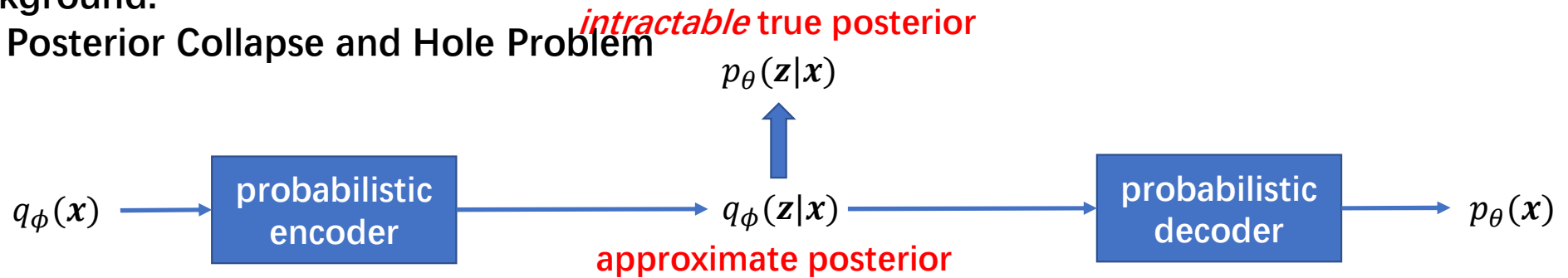
1. Background:
c. Existing methods



Improving Variational Autoencoders with Density Gap-based Regularization

1. Background:

b. Posterior Collapse and Hole Problem



$$\mathcal{L}_{ELBO}(\theta, \phi, x) = E_{q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}(q_\phi(z|x) || p_\theta(z))$$

Hole Problem:

$$q_\phi(z) \neq p_\theta(z)$$

where, $q_\phi(z) = E_{q_\phi(x)}[q_\phi(z|x)]$: the aggregated approximate posterior distribution

→ $\exists z q_\phi(z) \neq p_\theta(z)$

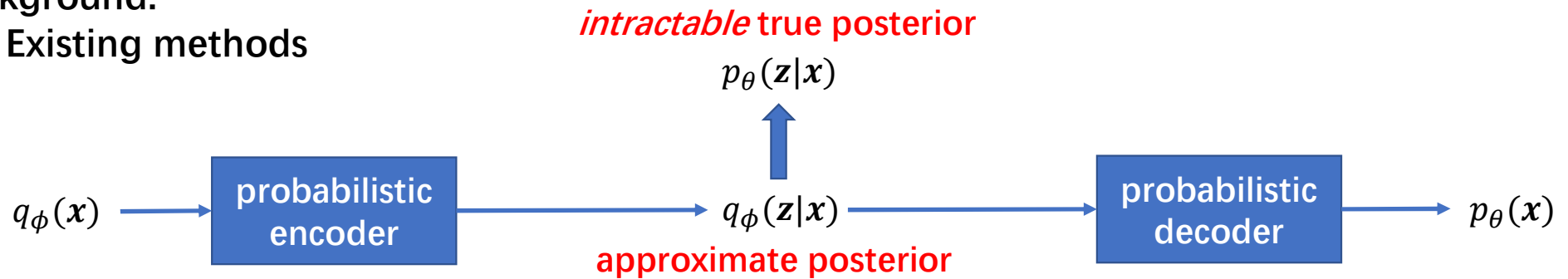
i.e. there exist areas (named as holes) with mismatch between density in $q_\phi(z)$ and $p_\theta(z)$

Empirically, inferences located in such areas are observed to **perform low-quality generation**, e.g., obscure and corrupted images, or sentences against commonsense.

Improving Variational Autoencoders with Density Gap-based Regularization

1. Background:

c. Existing methods



$$\mathcal{L}_{ELBO}(\theta, \phi, x) = E_{q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}(q_\phi(z|x) || p_\theta(z))$$

Hole Problem:

$$q_\phi(z) \neq p_\theta(z)$$

For image generation:

→ ascribed to the limited expressivity of $p_\theta(z)$ ($p_\theta(z) = N(\mathbf{0}, \mathbf{I})$ by default)

→ tackled by increasing the flexibility of $p_\theta(z)$ through:

hierarchical priors, energy-based models, a mixture of encoders, etc.

For text generation:

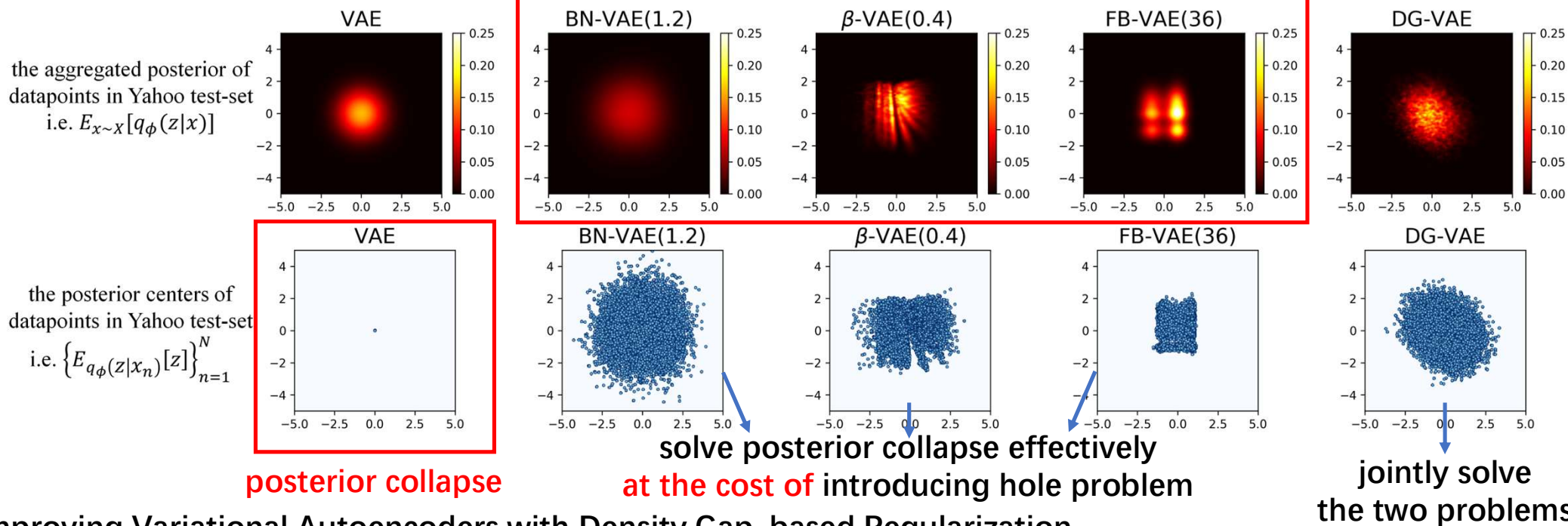
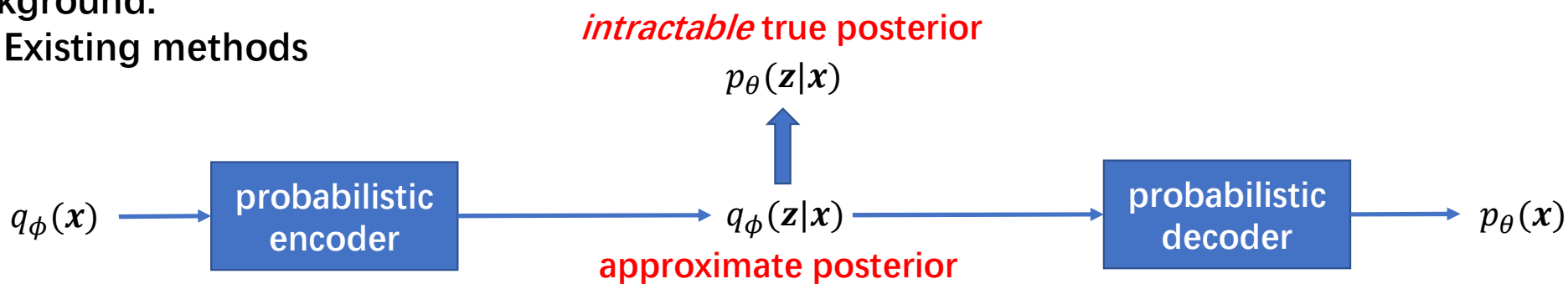
→ there's still little work on this, and we found that:

1. the vanilla VAEs (with $p_\theta(z) = N(\mathbf{0}, \mathbf{I})$) for text generation has no hole problem;
2. existing methods can solve posterior collapse effectively **at the cost of** introducing hole problem;

Improving Variational Autoencoders with Density Gap-based Regularization

1. Background:

c. Existing methods



Improving Variational Autoencoders with Density Gap-based Regularization

2. Methodology:

a. Regularization on the aggregated posterior distribution

rethink of $\mathcal{L}_{ELBO}(\theta, \phi, \mathbf{x})$:

$$\mathcal{L}_{ELBO}(\theta, \phi, \mathbf{x}) = E_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \underline{D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}))}$$

Q1: Since $q_{\phi}(\mathbf{z}|\mathbf{x})$ should not be too close to $p_{\theta}(\mathbf{z})$ (otherwise it will lead to posterior collapse), what should be close to $p_{\theta}(\mathbf{z}) = E_{p_{\theta}(\mathbf{x})}[p_{\theta}(\mathbf{z}|\mathbf{x})]$?

A1: The aggregated posterior distribution $q_{\phi}(\mathbf{z}) = E_{q_{\phi}(\mathbf{x})}[q_{\phi}(\mathbf{z}|\mathbf{x})]$.

Q2: So, how about regularizing $q_{\phi}(\mathbf{z})$ towards $p_{\theta}(\mathbf{z})$ instead in VAEs?

A2: It turns out to maximize $E_{q_{\phi}(\mathbf{x})}\mathcal{L}_{ELBO}(\theta, \phi, \mathbf{x}) + \mathbb{I}_{q_{\phi}(\mathbf{n}, \mathbf{z})}[\mathbf{n}, \mathbf{z}]$ (Hoffman et al. 2016):

$$E_{q_{\phi}(\mathbf{x})}\mathcal{L}_{ELBO}(\theta, \phi, \mathbf{x}) + \mathbb{I}_{q_{\phi}(\mathbf{n}, \mathbf{z})}[\mathbf{n}, \mathbf{z}] = E_{q_{\phi}(\mathbf{x})}E_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \underline{D_{KL}(q_{\phi}(\mathbf{z})||p_{\theta}(\mathbf{z}))}$$

$$\mathbb{I}_{q_{\phi}(\mathbf{n}, \mathbf{z})}[\mathbf{n}, \mathbf{z}] = E_{q_{\phi}(\mathbf{n}, \mathbf{z})}[\log \frac{q_{\phi}(\mathbf{n}, \mathbf{z})}{q_{\phi}(\mathbf{n})q_{\phi}(\mathbf{z})}]$$

where \mathbf{n} is the identity of datapoints in \mathbf{x} , i.e., $q_{\phi}(\mathbf{n} = n) = \frac{1}{N}$, ($n = 1, 2, \dots, N$)

effect: 1. weaken the regularization on $q_{\phi}(\mathbf{z}|\mathbf{x})$; 2. ensure $q_{\phi}(\mathbf{z}) \approx p_{\theta}(\mathbf{z})$.

Improving Variational Autoencoders with Density Gap-based Regularization

2. Methodology:

a. Regularization on the aggregated posterior distribution

Q3: Has anyone tried “regularizing $q_\phi(\mathbf{z})$ towards $p_\theta(\mathbf{z})$ instead in VAEs”?

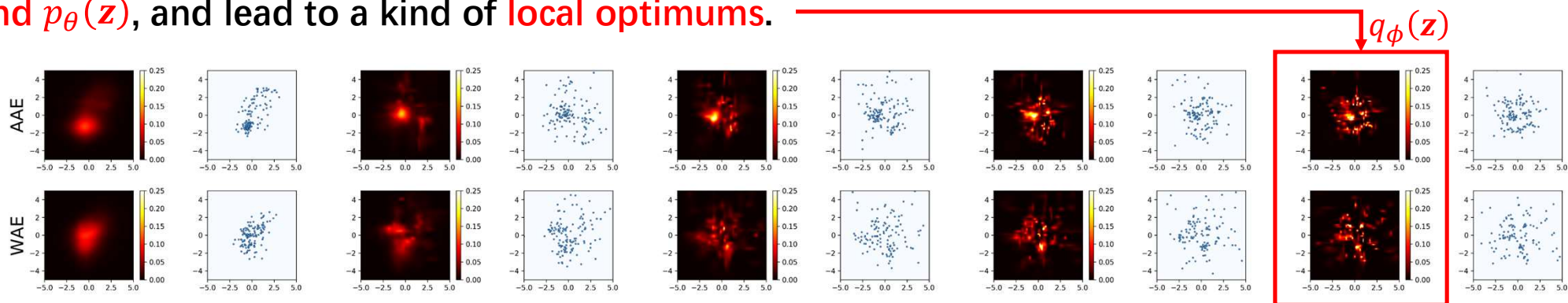
A3: Yes, as below:

AAE (Adversarial Auto-Encoder): minimize their JS divergence in the framework of GAN

WAE (Wasserstein Auto-Encoder): minimize the Maximum Mean Discrepancy between them

iVAE_{MI} (implicit VAE + MI regularization): minimize a dual form of KL divergence between them

But all their implementations of regularization are **based on merely sampling sets from $q_\phi(\mathbf{z})$ and $p_\theta(\mathbf{z})$** , and lead to a kind of **local optimum**s.



1. a sampling set from such a $q_\phi(\mathbf{z})$ can already stimulate that from $p_\theta(\mathbf{z})$ to some degree;

2. but such a $q_\phi(\mathbf{z})$ still have evident difference from $p_\theta(\mathbf{z})$

Intuitively, a sampling set from $q_\phi(\mathbf{z})$ can hardly be the same as that from $p_\theta(\mathbf{z})$, even when $q_\phi(\mathbf{z}) = p_\theta(\mathbf{z})$

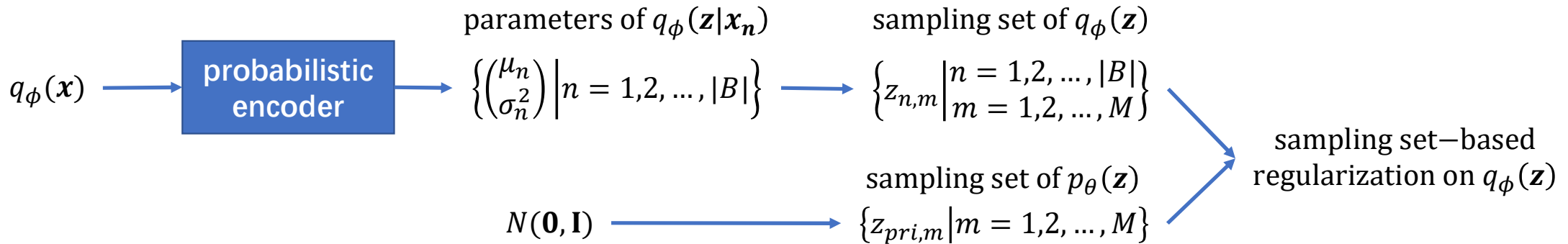
Improving Variational Autoencoders with Density Gap-based Regularization

2. Methodology:

b. Density Gap-based regularization

For example,

$$q_\phi(\mathbf{z}|\mathbf{x}_n) = N(\boldsymbol{\mu}_n, \boldsymbol{\sigma}_n^2), p_\theta(\mathbf{z}) = N(\mathbf{0}, \mathbf{I})$$

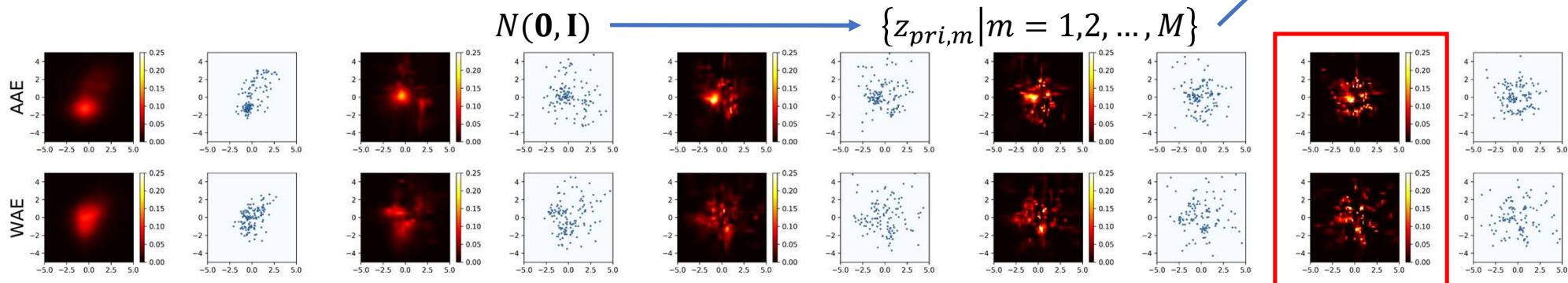
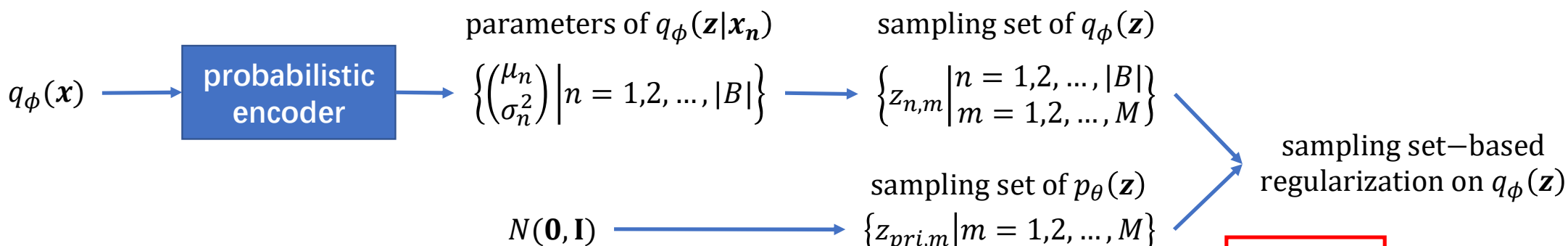


2. Methodology:

b. Density Gap-based regularization

For example,

$$q_\phi(\mathbf{z}|\mathbf{x}_n) = N(\boldsymbol{\mu}_n, \boldsymbol{\sigma}_n^2), p_\theta(\mathbf{z}) = N(\mathbf{0}, \mathbf{I})$$



1. a sampling set from such a $q_\phi(\mathbf{z})$ can already stimulate that from $p_\theta(\mathbf{z})$ to some degree;

2. but such a $q_\phi(\mathbf{z})$ still have evident difference from $p_\theta(\mathbf{z})$

Intuitively, a sampling set from $q_\phi(\mathbf{z})$ can hardly be the same as that from $p_\theta(\mathbf{z})$, even when $q_\phi(\mathbf{z}) = p_\theta(\mathbf{z})$

Improving Variational Autoencoders with Density Gap-based Regularization

2. Methodology:

b. Density Gap-based regularization

Intuitively, a sampling set from $q_\phi(\mathbf{z})$ can hardly be the same as that from $p_\theta(\mathbf{z})$, even when $q_\phi(\mathbf{z}) = p_\theta(\mathbf{z})$

→ The probability density of $q_\phi(\mathbf{z})$ and $p_\theta(\mathbf{z})$ are the same everywhere when $q_\phi(\mathbf{z}) = p_\theta(\mathbf{z})$

→ Density Gap-based regularization:

$$KL(q_\phi(\mathbf{z}) || p_\theta(\mathbf{z})) = E_{q_\phi(\mathbf{z})}[\log \frac{q_\phi(\mathbf{z})}{p_\theta(\mathbf{z})}]$$

$$\approx \frac{1}{S} \sum_{s=1}^S [\log q_\phi(z_s) - \log p_\theta(z_s)]$$

we refer to this as

the density gap
between $q_\phi(\mathbf{z})$ and $p_\theta(\mathbf{z})$
at position z_s

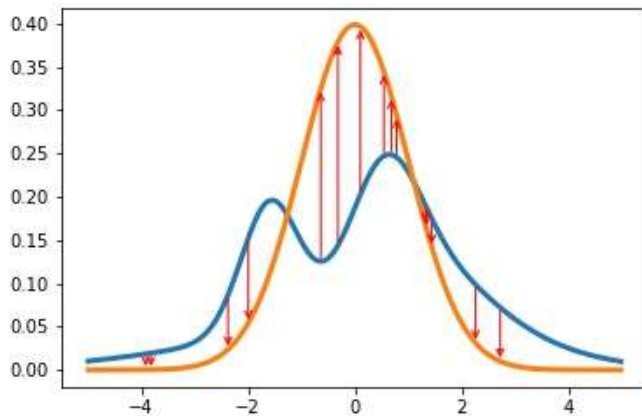
where,

z_s is the s^{th} sample from $q_\phi(\mathbf{z})$

$q_\phi(z_s)$ and $p_\theta(z_s)$ are the values of
corresponding PDFs

stratified sampling &
reparameterization trick

parametric differentiable PDFs



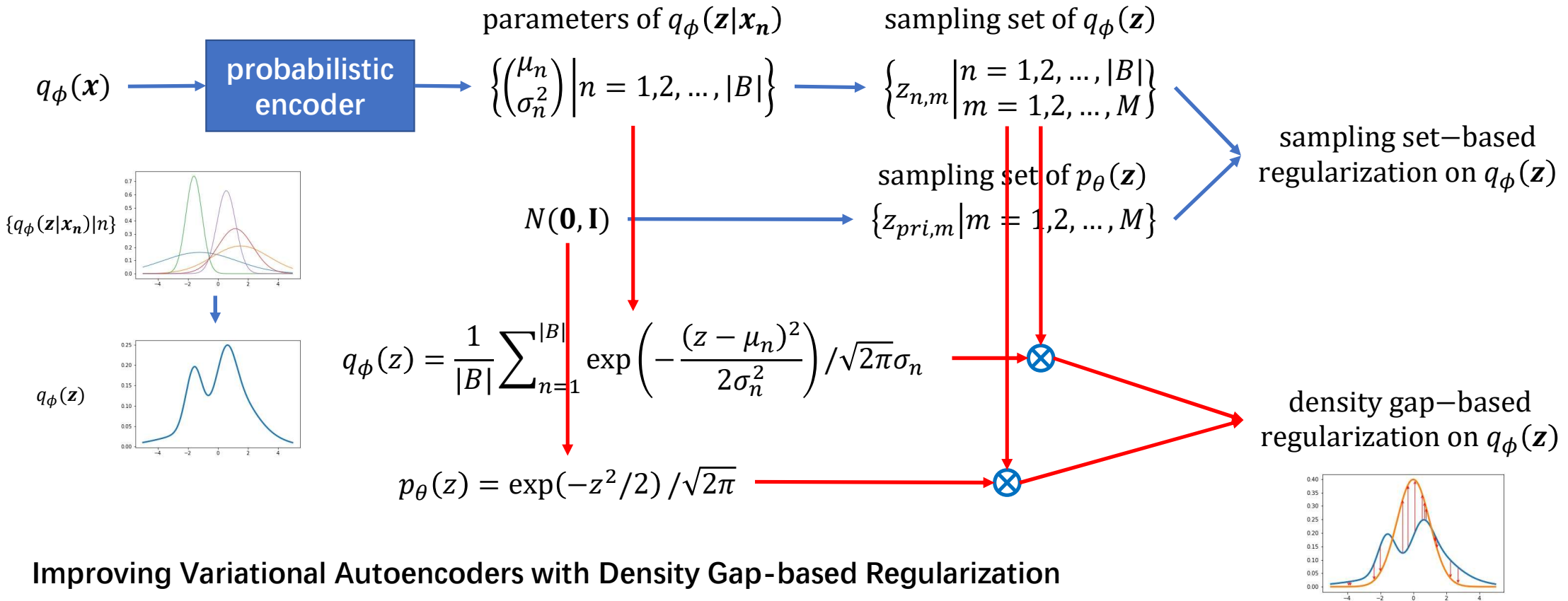
Improving Variational Autoencoders with Density Gap-based Regularization

2. Methodology:

b. Density Gap-based regularization

For example,

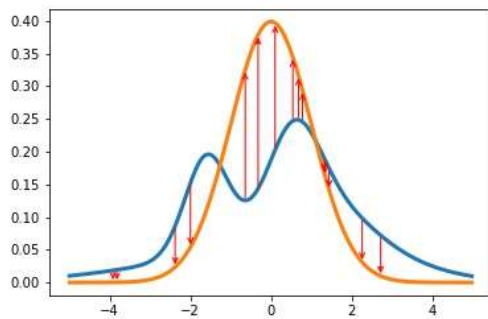
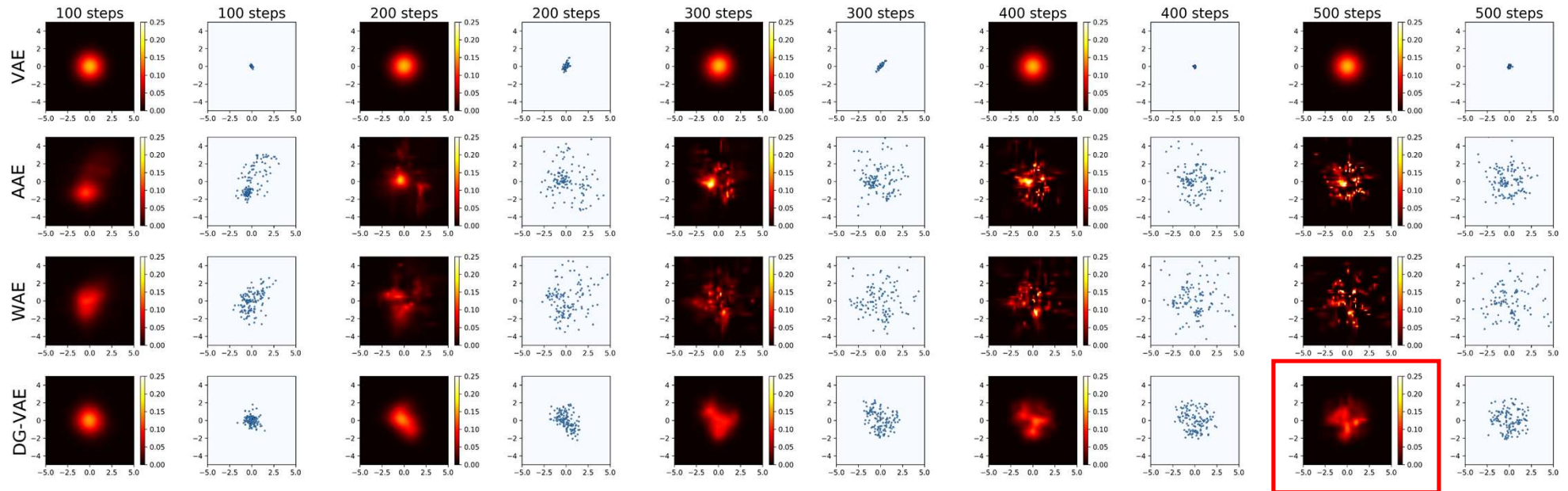
$$q_\phi(\mathbf{z}|\mathbf{x}_n) = N(\boldsymbol{\mu}_n, \boldsymbol{\sigma}_n^2), p_\theta(\mathbf{z}) = N(\mathbf{0}, \mathbf{I})$$



Improving Variational Autoencoders with Density Gap-based Regularization

2. Methodology:

b. Density Gap-based regularization



regularize $q_\phi(z)$ towards $p_\theta(z)$ in the perspective of their mismatch in PDFs

Improving Variational Autoencoders with Density Gap-based Regularization

2. Methodology:

c. Marginal regularization for more Mutual Information

We can apply the proposed regularization in training with mini-batch gradient descent:

$$E_{q_\phi(\mathbf{x})} \mathcal{L}_{ELBO}(\theta, \phi, \mathbf{x}) + \mathbb{I}_{q_\phi(\mathbf{n}, \mathbf{z})}[\mathbf{n}, \mathbf{z}] = E_{q_\phi(\mathbf{x})} E_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - D_{KL}(q_\phi(\mathbf{z}) || p_\theta(\mathbf{z}))$$

where the data distribution $q_\phi(\mathbf{x})$ is described by the current mini-batch B :

$$B = \{x_1, x_2, \dots, x_{|B|}\}$$
$$q_\phi(\mathbf{x} = x_n) = q_\phi(\mathbf{n}) = \frac{1}{|B|}$$

→ the mutual information term to maximize has a limited upper bound:

$$\mathbb{I}_{q_\phi(\mathbf{n}, \mathbf{z})}[\mathbf{n}, \mathbf{z}] = H_{q_\phi(\mathbf{n})}(\mathbf{n}) - H_{q_\phi(\mathbf{n}, \mathbf{z})}(\mathbf{n}|\mathbf{z}) \leq H_{q_\phi(\mathbf{n})}(\mathbf{n}) = \log|B| < \log N$$

→ for a high dimensional prior distribution, it still have limited effect on solving posterior collapse (it is already enough for $\mathbb{I}_{q_\phi(\mathbf{n}, \mathbf{z})}[\mathbf{n}, \mathbf{z}]$ to reach $\log|B|$ with limited dimensions of \mathbf{z} being activated)

→ in order to activate all dimensions of \mathbf{z} , we propose marginal regularization:

$$E_{q_\phi(\mathbf{x})} \mathcal{L}_{ELBO}(\theta, \phi, \mathbf{x}) + \sum_{i=1}^{Dim} \mathbb{I}_{q_\phi(\mathbf{n}, \mathbf{z}_i)}[\mathbf{n}, \mathbf{z}_i] = E_{q_\phi(\mathbf{x})} E_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - \sum_{i=1}^{Dim} D_{KL}(q_\phi(\mathbf{z}_i) || p_\theta(\mathbf{z}_i))$$

where $i = 1, 2, \dots, Dim$ denotes the index of dimension, \mathbf{z}_i denotes the i^{th} component of \mathbf{z} , $q_\phi(\mathbf{z}_i)$ and $p_\theta(\mathbf{z}_i)$ denote the marginal distribution of $q_\phi(\mathbf{z})$ and $p_\theta(\mathbf{z})$ on the i^{th} dimension respectively.

Improving Variational Autoencoders with Density Gap-based Regularization

2. Methodology:

c. Marginal regularization for more Mutual Information

→ in order to activate all dimensions of \mathbf{z} , we propose marginal regularization:

$$E_{q_\phi(x)} \mathcal{L}_{ELBO}(\theta, \phi, \mathbf{x}) + \sum_{i=1}^{Dim} \mathbb{I}_{q_\phi(\mathbf{n}, \mathbf{z}_i)}[\mathbf{n}, \mathbf{z}_i] = E_{q_\phi(x)} E_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - \sum_{i=1}^{Dim} D_{KL}(q_\phi(\mathbf{z}_i) || p_\theta(\mathbf{z}_i))$$

where $i = 1, 2, \dots, Dim$ denotes the index of dimension, \mathbf{z}_i denotes the i^{th} component of \mathbf{z} , $q_\phi(\mathbf{z}_i)$ and $p_\theta(\mathbf{z}_i)$ denote the marginal distribution of $q_\phi(\mathbf{z})$ and $p_\theta(\mathbf{z})$ on the i^{th} dimension respectively.

→ in such way, the mutual information term to maximize has an upper bound linear with Dim :

$$\sum_{i=1}^{Dim} \mathbb{I}_{q_\phi(\mathbf{n}, \mathbf{z}_i)}[\mathbf{n}, \mathbf{z}_i] \leq \sum_{i=1}^{Dim} H_{q_\phi(\mathbf{n})}(\mathbf{n}) = Dim * \log|B|$$

→ we implement this for VAEs with $p_\theta(\mathbf{z}) = N(\mathbf{0}, \mathbf{I})$, as **its marginal distributions are independent**:

$$p_\theta(\mathbf{z}) = \prod_{i=1}^{Dim} p_\theta(\mathbf{z}_i)$$

→ it should be noted that, this independency-based decomposition of $p_\theta(\mathbf{z})$ **is not established for von Mises-Fisher distributions**, e.g., $p_\theta(\mathbf{z}) = vMF(\boldsymbol{\mu}, \kappa)$, so we only implement the joint regularization for von Mises-Fisher distribution-based VAEs.

2. Methodology:

d. Aggregation size for ablation

→ to further investigate the effect of maximizing mutual information, we split the mini-batch B into non-overlapping subsets:

$$B = \bigcup_{i=1}^C b_i, \text{ s.t. } b_i \cap b_j = \emptyset \text{ iff } i \neq j$$

those subsets have the same size $|b| = |b_i| = |b_j| = \frac{|B|}{C}$ which **we refer to as the aggregation size**, as we only calculate the aggregated posterior distributions inside each subsets, and regularize them to the prior distribution respectively:

$$q_{\phi,j}(\mathbf{z}) = E_{x \sim b_j} [q_{\phi}(\mathbf{z}|\mathbf{x})]$$
$$\sum_{j=1}^C \sum_{i=1}^{Dim} D_{KL}(q_{\phi,j}(\mathbf{z}) || p_{\theta}(\mathbf{z}_i))$$

→ in such way, the maximized mutual information term has an upper bound linear with $\log |b|$:

$$\sum_{j=1}^C \sum_{i=1}^{Dim} \mathbb{I}_{q_{\phi,j}(\mathbf{n}, \mathbf{z}_i)}[\mathbf{n}, \mathbf{z}_i] \leq \sum_{j=1}^C \sum_{i=1}^{Dim} H_{q_{\phi,j}(\mathbf{n})}(\mathbf{n}) = C * Dim * \log |b|$$

when $|b| = 1$, the proposed method is equivalent to the vanilla VAE.

3. Experiment

a. Language modeling

Table 2: Results of Language Modeling on Yahoo dataset. We bold up $MI(\phi) \geq 9.0$, $AU(\phi) \geq 30$, $CU(\phi) \geq 30$, the highest $priorLL(\theta)$ and $postLL(\theta, \phi)$ for the same methods.

Models	$priorLL(\theta)$	$postLL(\theta, \phi)$	$KL(\phi)$	$MI(\phi)$	$AU(\phi)$	$CU(\phi)$
VAE (default)	-330.7	-330.7	0.0	0.0	0	32
cyclic-VAE	-329.8	-328.9	1.1	1.0	2	31
bow-VAE	-330.5	-330.5	0.0	0.0	0	32
skip-VAE	-330.1	-325.2	5.0	4.3	8	31
δ -VAE(0.15)	-330.5	-330.6	4.8	0.0	0	0
BN-VAE(0.6)	-327.6	-321.1	6.6	5.9	32	32
BN-VAE(1.2)	-330.9	-310.1	26.2	9.2	32	0
BN-VAE(1.8)	-343.5	-308.6	51.3	9.2	32	0
FB-VAE(4)	-329.8	-328.4	3.9	1.8	32	32
FB-VAE(16)	-325.7	-320.8	16.1	8.5	32	8
FB-VAE(49)	-344.6	-296.1	50.0	9.2	32	0
β -VAE(0.4)	-330.8	-324.8	7.0	6.7	3	31
β -VAE(0.2)	-338.6	-310.3	30.1	9.2	22	25
β -VAE(0.1)	-369.9	-289.6	83.7	9.2	32	0
DG-VAE ($ b = 1$)	-330.7	-330.7	0.0	0.0	0	32
DG-VAE ($ b = 4$)	-330.4	-318.3	14.3	9.1	11	32
DG-VAE ($ b = 32$)	-355.4	-294.1	65.2	9.1	32	32
DG-VAE (default)	-358.0	-290.8	70.8	9.1	32	32

Table 1: Statistics of sentences in the datasets

Dataset	Train	Valid	Test	Vocab size	Length (avg \pm std)
Yelp	100,000	10,000	10,000	19997	98.01 \pm 48.86
Yahoo	100,000	10,000	10,000	20001	80.76 \pm 46.21
Short-Yelp	100,000	10,000	10,000	8411	10.96 \pm 3.60
SNLI	100,000	10,000	10,000	9990	11.73 \pm 4.33

$$priorLL(\theta) = E_x \log E_{p_{\theta}(z)} [p_{\theta}(x|z)]$$

$$postLL(\theta, \phi) = E_x \log E_{q_{\phi}(z|x)} [p_{\theta}(x|z)]$$

$$KL(\phi) = E_x KL(q_{\phi}(z|x) || p_{\theta}(z))$$

$$MI(\phi) = H(q_{\phi}(z)) - E_x H(q_{\phi}(z|x))$$

$$AU(\phi) = |\{i | Var_x E_{q_{\phi}(z|x)} [z_i] > 0.01\}|$$

Small values indicate posterior collapse

$$CU(\phi) = |\{i | KL(q_{\phi}(z_i) || p_{\theta}(z_i)) < 0.03\}|$$

Small values indicate the hole problem

Improving Variational Autoencoders with Density Gap-based Regularization

3. Experiment

a. Language modeling

Table 1: Statistics of sentences in the datasets

Dataset	Train	Valid	Test	Vocab size	Length (avg \pm std)
Yelp	100,000	10,000	10,000	19997	98.01 \pm 48.86
Yahoo	100,000	10,000	10,000	20001	80.76 \pm 46.21
Short-Yelp	100,000	10,000	10,000	8411	10.96 \pm 3.60
SNLI	100,000	10,000	10,000	9990	11.73 \pm 4.33

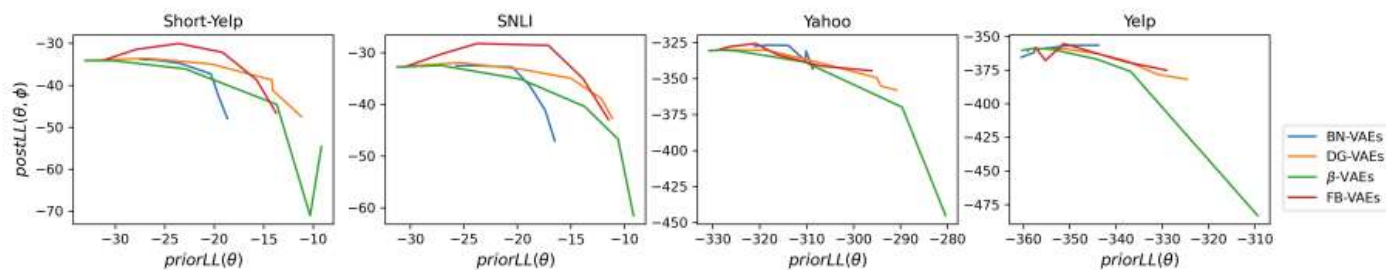


Figure 2: The curves of $priorLL(\theta)$ and $postLL(\theta, \phi)$ in Gaussian distribution-based VAEs.

$$priorLL(\theta) = E_x \log E_{p_{\theta}(z)} [p_{\theta}(x|z)]$$

$$postLL(\theta, \phi) = E_x \log E_{q_{\phi}(z|x)} [p_{\theta}(x|z)]$$

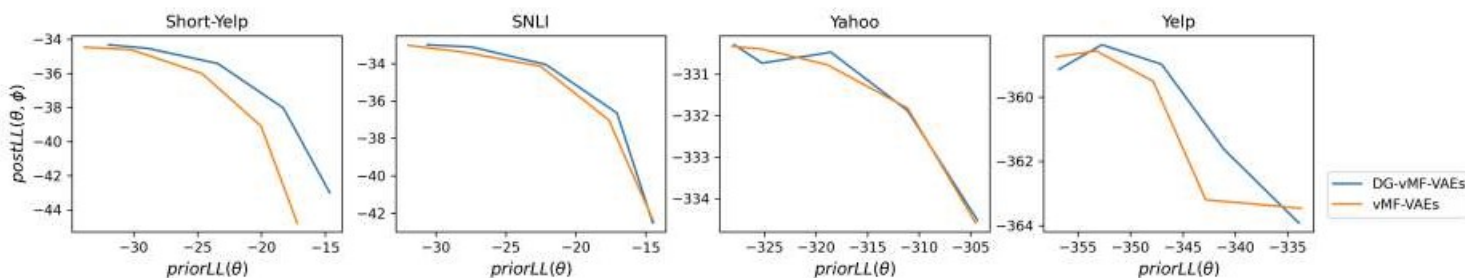


Figure 3: The curves of $priorLL(\theta)$ and $postLL(\theta, \phi)$ in vMF distribution-based VAEs.

Improving Variational Autoencoders with Density Gap-based Regularization

3. Experiment

b. Visualization of the posterior

Table 1: Statistics of sentences in the datasets

Dataset	Train	Valid	Test	Vocab size	Length (avg \pm std)
Yelp	100,000	10,000	10,000	19997	98.01 \pm 48.86
Yahoo	100,000	10,000	10,000	20001	80.76 \pm 46.21
Short-Yelp	100,000	10,000	10,000	8411	10.96 \pm 3.60
SNLI	100,000	10,000	10,000	9990	11.73 \pm 4.33

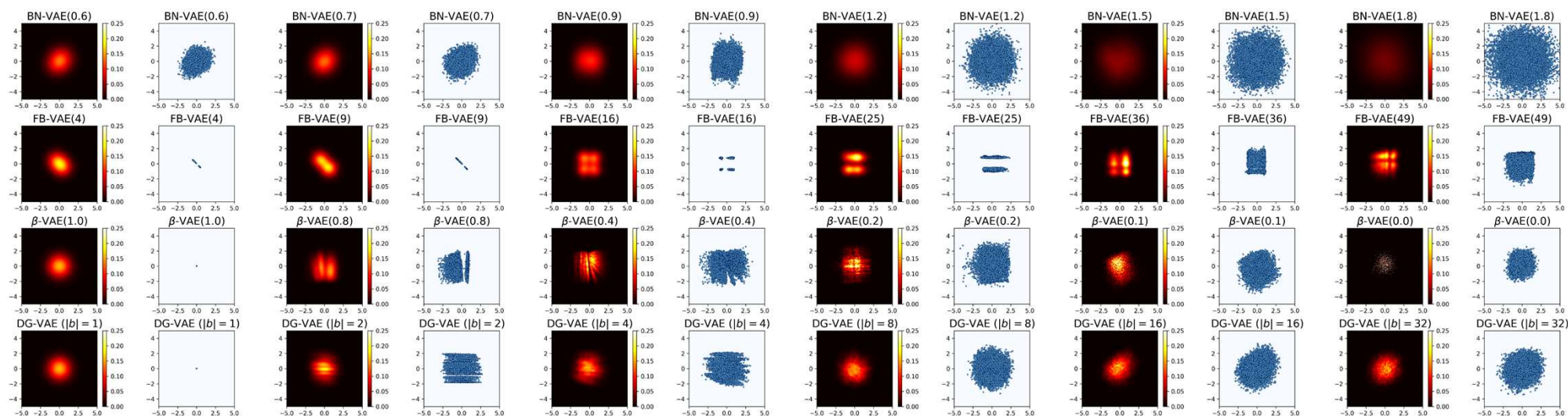


Figure 4: The visualization of the aggregated posterior distributions (red-in-black) and the posterior centers distributions (blue-in-white) for BN-VAEs, FB-VAEs, β -VAEs, and DG-VAEs on the Yahoo test-set. Illustrations for more datasets, more models, and more dimensions, are shown in Appendix G.

3. Experiment

c. Interpolation study

$$z_a, z_b \sim q_\phi(\mathbf{z}|x_a), q_\phi(\mathbf{z}|x_b)$$

$$z_\lambda = \lambda * z_a + (1 - \lambda) * z_b$$

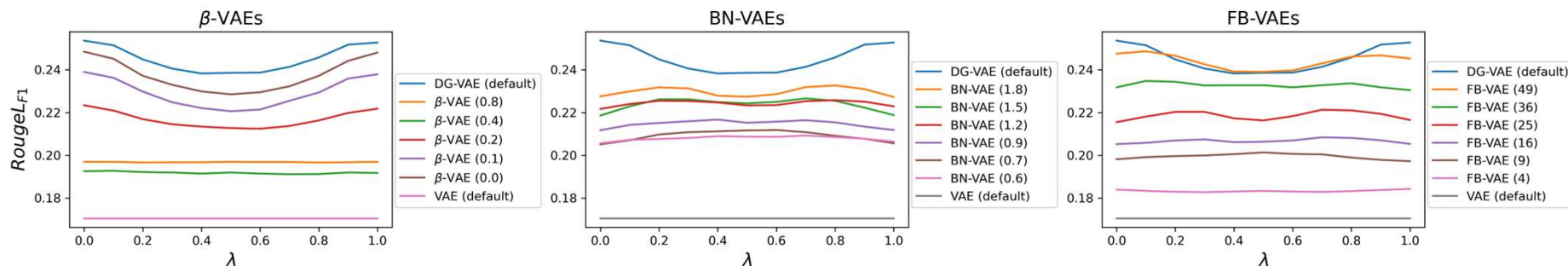
$$x_\lambda \sim p_\theta(\mathbf{x}|z_\lambda)$$

$$RougeL_{F1} = \frac{1}{2} (F_{lcs}(x_a, x_\lambda) + F_{lcs}(x_b, x_\lambda))$$

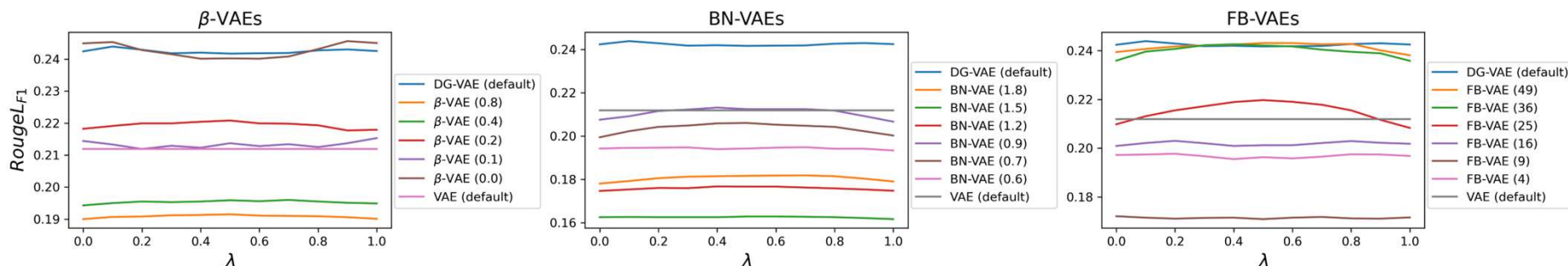
Table 1: Statistics of sentences in the datasets

Dataset	Train	Valid	Test	Vocab size	Length (avg \pm std)
Yelp	100,000	10,000	10,000	19997	98.01 \pm 48.86
Yahoo	100,000	10,000	10,000	20001	80.76 \pm 46.21
Short-Yelp	100,000	10,000	10,000	8411	10.96 \pm 3.60
SNLI	100,000	10,000	10,000	9990	11.73 \pm 4.33

Yahoo



Yelp



Improving Variational Autoencoders with Density Gap-based Regularization

3. Experiment

c. Interpolation study

$$z_a, z_b \sim q_\phi(\mathbf{z}|x_a), q_\phi(\mathbf{z}|x_b)$$

$$z_\lambda = \lambda * z_a + (1 - \lambda) * z_b$$

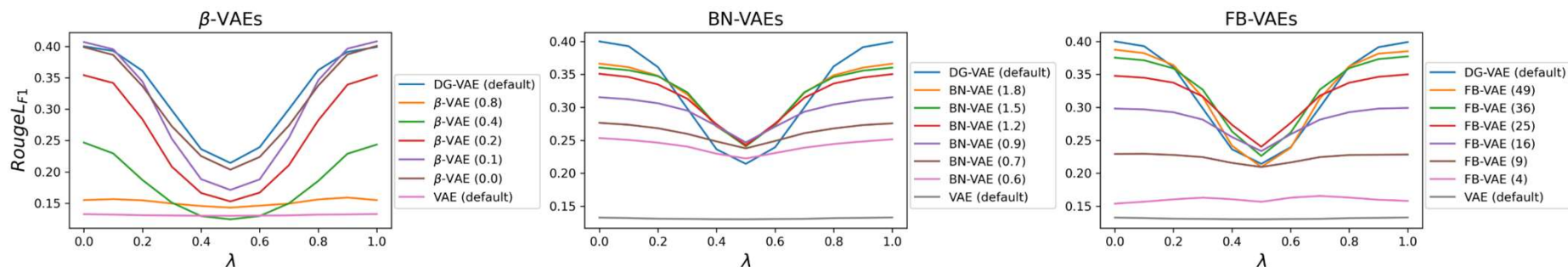
$$x_\lambda \sim p_\theta(\mathbf{x}|z_\lambda)$$

$$RougeL_{F1} = \frac{1}{2} (F_{lcs}(x_a, x_\lambda) + F_{lcs}(x_b, x_\lambda))$$

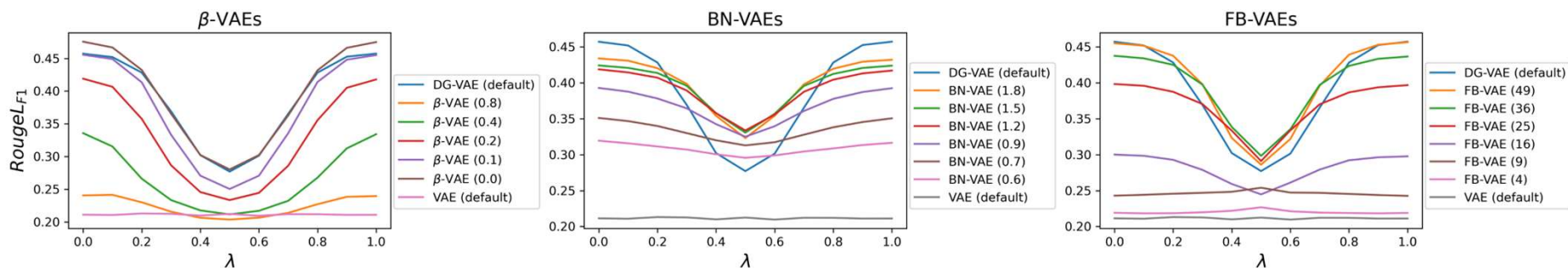
Table 1: Statistics of sentences in the datasets

Dataset	Train	Valid	Test	Vocab size	Length (avg \pm std)
Yelp	100,000	10,000	10,000	19997	98.01 \pm 48.86
Yahoo	100,000	10,000	10,000	20001	80.76 \pm 46.21
Short-Yelp	100,000	10,000	10,000	8411	10.96 \pm 3.60
SNLI	100,000	10,000	10,000	9990	11.73 \pm 4.33

Short-Yelp



SNLI



Improving Variational Autoencoders with Density Gap-based Regularization

3. Experiment

c. Interpolation study

$$z_a, z_b \sim q_\phi(\mathbf{z}|x_a), q_\phi(\mathbf{z}|x_b)$$

$$z_\lambda = \lambda * z_a + (1 - \lambda) * z_b$$

$$x_\lambda \sim p_\theta(x|z_\lambda)$$

$$RougeL_{F1} = \frac{1}{2} (F_{lcs}(x_a, x_\lambda) + F_{lcs}(x_b, x_\lambda))$$

x_a : two girls walking in a park .		x_b : the two kids are playing in water	
β -VAE(0.1)		DG-VAE	
λ	$DG(z_\lambda)$	x_λ	$DG(z_\lambda)$
0.0	32.9	two girls walking in a park .	38.9
0.1	26.1	two girls walking in a park .	35.5
0.2	11.4	two women walking in a park .	29.4
0.3	-11.3	two women walking in a pool .	20.6
0.4	-42.2	two cats playing in a pool .	9.2
0.5	-55.4	an african man walks in the pool .	-4.8
0.6	-48.0	an elderly man walks in water .	1.0
0.7	-15.1	the two children are playing in water .	18.7
0.8	12.2	the two children are playing in water	31.7
0.9	30.0	the two kids are playing in water	40.1
1.0	38.2	the two kids are playing in water	43.9

Table 1: Statistics of sentences in the datasets

Dataset	Train	Valid	Test	Vocab size	Length (avg \pm std)
Yelp	100,000	10,000	10,000	19997	98.01 \pm 48.86
Yahoo	100,000	10,000	10,000	20001	80.76 \pm 46.21
Short-Yelp	100,000	10,000	10,000	8411	10.96 \pm 3.60
SNLI	100,000	10,000	10,000	9990	11.73 \pm 4.33

x_a : a man in a white shirt and black pants poses in front of a large banner .		x_b : two people hug each other to warm up while they are locked out of the house .	
β -VAE(0.1)		DG-VAE	
λ	$DG(z_\lambda)$	x_λ	$DG(z_\lambda)$
0.0	60.2	a man in a black shirt and black pants sits in front of a large gathering .	56.5
0.1	38.6	a man in a black shirt and blue pants walking in front of a large gathering .	44.2
0.2	-21.9	a man in jeans and a white shirt walking down in an orange kayak in the forest .	12.5
0.3	-121.4	a man in shorts and a black shirt walking through snow , on the street .	-38.6
0.4	-157.4	a toddler , wearing shorts and black pants walking a green scooter while looking in the water .	-47.2
0.5	-159.9	a toddler girl wearing black shorts and sandals walking through her house while on the sunny sidewalk .	-53.7
0.6	-161.7	a toddler girl wearing pink pants and boots walks across the street in front of cars .	-63.5
0.7	-164.6	a girl , who looks over her head while she sits alone on the edge .	-76.7
0.8	-46.4	two people decide whether , as they walk up in the water while looking up .	-26.2
0.9	37.5	two people decide whether to each other , and one is out out of the window .	40.8
1.0	66.8	two people hug as they walk out and sun to get out of the sun .	64.8

3. Experiment

c. Interpolation study

$$z_a, z_b \sim q_\phi(\mathbf{z}|x_a), q_\phi(\mathbf{z}|x_b)$$

$$z_\lambda = \lambda * z_a + (1 - \lambda) * z_b$$

$$x_\lambda \sim p_\theta(\mathbf{x}|z_\lambda)$$

$$RougeL_{F1} = \frac{1}{2} (F_{lcs}(x_a, x_\lambda) + F_{lcs}(x_b, x_\lambda))$$

x_a : great place for a romantic <unk> .		x_b : the asian cucumber salad was bland .	
β -VAE(0.1)		DG-VAE	
λ	$DG(z_\lambda)$	x_λ	$DG(z_\lambda)$
0.0	46.9	great place for a romantic <unk> .	52.1
0.1	32.2	great place for a romantic <unk> .	38.3
0.2	-5.6	great place for a chilly <unk> .	3.5
0.3	-55.6	oh you 're perfect and special .	-16.7
0.4	-60.4	oh you 'll enjoy the special .	-9.2
0.5	-75.0	the guys keep it clean though .	-9.4
0.6	-86.0	the apartments make it was comfortable .	-17.2
0.7	-72.0	the specialty pie are good out .	-32.7
0.8	-6.3	the wood martinis are very cheap .	11.6
0.9	34.3	the wood martinis taste was bland .	43.6
1.0	49.7	the english muffins were good bland .	56.8

Table 1: Statistics of sentences in the datasets

Dataset	Train	Valid	Test	Vocab size	Length (avg \pm std)
Yelp	100,000	10,000	10,000	19997	98.01 \pm 48.86
Yahoo	100,000	10,000	10,000	20001	80.76 \pm 46.21
Short-Yelp	100,000	10,000	10,000	8411	10.96 \pm 3.60
SNLI	100,000	10,000	10,000	9990	11.73 \pm 4.33

x_a : our server was not even <unk> familiar with the food or food preparation .		x_b : i have had just about everything on the menu and everything is delicious .	
β -VAE(0.1)		DG-VAE	
λ	$DG(z_\lambda)$	x_λ	$DG(z_\lambda)$
0.0	82.1	our server was not even warm about the food and the quality service .	69.8
0.1	58.0	our server was not even busy on the menu and the food network .	51.9
0.2	-11.4	our waitress was not even busy on the menu and the food sucked .	1.0
0.3	-126.1	still they was not even warm by the food and taste was awesome .	-33.7
0.4	-179.0	still unfortunately the cashier just kept asking the menu and was seriously awesome .	-32.8
0.5	-177.4	even we was nothing as much of the menu and food was decent .	-32.4
0.6	-169.5	then i still had all about eating at the food and everything was good .	-33.2
0.7	-63.0	did i say before they use no food and there are very quick .	-13.2
0.8	13.6	did i say before they use from the menu and menu was decent .	29.4
0.9	60.0	i have had just twice there and their food is very yummy .	56.5
1.0	76.3	i have had just because of the menu and everything is awesome .	68.1