
Flexible Text Modeling with Semi-Implicit Latent Representations

Hao Fu
Duke University
hao.fu@duke.edu

Chunyuan Li
Microsoft Research
chunyl@microsoft.com

Ke Bai
Duke University
ke.bai@duke.edu

Jianfeng Gao
Microsoft Research
jfgao@microsoft.com

Lawrence Carin
Duke University
lcarin@duke.edu

Abstract

Variational autoencoders (VAEs) have been used recently to learn probabilistic representations in many natural language processing (NLP) applications. Training deficiency has been witnessed when an auto-regressive decoder is used: the learned latent codes becomes almost identical to the prior distribution (termed “posterior collapse,” resulting in a vanishing of the Kullback-Leibler term in the variational expression). We hypothesize that the source of deficiency is due partially to the approximate Gaussian posterior often used in variational inference. We use *semi-implicit* (SI) representations for the latent distributions of the natural languages. It extends the commonly used Gaussian distribution family, by mixing the variational parameter with a flexible implicit distribution. The increasing representational power of SI is demonstrated on two NLP tasks, showing that it provides more informative latent codes in language modeling, and generates more diverse dialog responses.

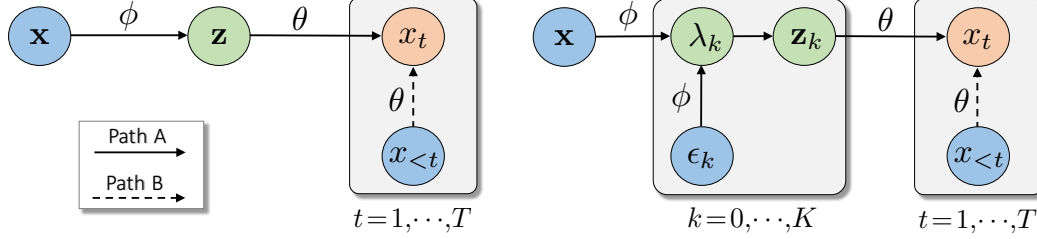
1 Introduction

Deep latent-variable models such as variational autoencoders (VAEs) [13, 21] are becoming increasingly popular in natural language processing (NLP). They have contributed to fundamental advances in many NLP tasks, such as language modeling [1, 17] and dialog response generation [27, 25]. For a text sequence of length T , $\mathbf{x} = [x_1, \dots, x_T]$, neural language models [18] generate the t -th token x_t conditioned on the previously generated tokens: $p(\mathbf{x}) = \prod_{t=1}^T p(x_t | x_{<t})$, where $x_{<t}$ indicates all tokens before t . However, language models lack an efficient inference mechanism, preventing them from reasoning about data at an abstract level. For instance, language models don’t allow the sort of neural sentence manipulations showcased in [1].

The VAE was introduced to fill the gap, simultaneously learning generative models with higher-quality sentence samples while learning an efficient inference network [13, 21, 1]. The *decoder* draws a continuous latent vector \mathbf{z} from prior $p(\mathbf{z})$, and generates the text sequence \mathbf{x} from a conditional distribution $p_{\theta}(\mathbf{x} | \mathbf{z})$; $p(\mathbf{z})$ is typically assumed a multivariate Gaussian, and θ represents the network parameters. The following auto-regressive decoding process is usually used:

$$p_{\theta}(\mathbf{x} | \mathbf{z}) = \prod_{t=1}^T p_{\theta}(x_t | x_{<t}, \mathbf{z}). \quad (1)$$

The true posterior $p_{\theta}(\mathbf{z} | \mathbf{x}) \propto p_{\theta}(\mathbf{x} | \mathbf{z})p(\mathbf{z})$ is approximated via the variational distribution $q_{\phi}(\mathbf{z} | \mathbf{x})$. This distribution is often produced with an *encoder*, implemented via a ϕ -parameterized neural



(a) Vanilla VAE with an auto-regressive decoder (b) SI-VAE with an auto-regressive decoder

Figure 1: Illustration of learning bottleneck in the two different methods. For the VAE model with an auto-regressive decoder, two paths are considered from x to its reconstruction: Path A via $\{\phi, \theta\}$ and Path B via θ . The circles with blue, green and red indicate observed, latent, reconstructed variables, respectively. (a) Vanilla VAE with Gaussian posteriors represents the observation x as Gaussian latent code z , which encodes all information flow in Path A to decode its reconstruction. (b) Semi-implicit VAE represents the latent code of observation x as the mixture of Gaussian $\{z_k\}_{k=0}^K$, whose variational parameter λ_k is generated using x and auxiliary noise samples ϵ_k .

network. It yields the *evidence lower bound* (ELBO) as an objective:

$$\log p_{\theta}(x) \geq \mathcal{L} = \mathcal{E} - \mathcal{R}, \text{ with} \quad (2)$$

$$\mathcal{E} = \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] \quad (3)$$

$$\mathcal{R} = \text{KL}(q_{\phi}(z|x) || p(z)) \quad (4)$$

where \mathcal{E} is the log-likelihood term, and \mathcal{R} is a Kullback-Leibler (KL) regularization term.

KL Vanishing Issues The ELBO can be viewed as a regularized version of the autoencoder (AE) [7]. It is thus natural to extend the negative of \mathcal{L} in (2) by introducing a hyper-parameter β to control the strength of regularization [10]: $\mathcal{L}_{\beta} = \mathcal{E} - \beta\mathcal{R}$. When $\beta = 1$ (*constant*) during the training procedure, \mathcal{R} quickly becomes vanishingly small during training [1]. This is known as the *KL vanishing* issue, which causes undesirable outcomes: the encoder produces posteriors almost identical to the Gaussian prior, for all observations; the decoder ignores the z , and the VAE model reduces to a simpler language model. To address this issue, a *monotonic annealing schedule* was proposed to increase β from 0 to 1 during the early stage of VAE learning [1]. The *cyclical annealing schedule* further proposed to repeat this annealing process multiple time [5] during training.

2 Semi-Implicit Text Representations

2.1 Diagnostics of KL Vanishing

It has been shown recently that the auto-regressive decoder in the VAE has two paths working together to generate text sequences, and the KL vanishing issue happens due to the destructive competition between the two paths for the information flows [5], illustrated in Figure 1(a). Path A consists of $\{\phi, \theta\}$ passing through z . It first extracts the global representation z , then generates x directly. Path B employs the partial ground-truth information of x at every time step of the sequential decoding, and it generates x_t conditioned on $x_{<t}$. Hence, both paths compete to generate the observed text; the KL vanishing issue appears when Path B dominates Path A.

However, it is unknown why Path A tends to be weak in this competition. In this paper, we hypothesize that there is a *bottleneck* on Path A that may weaken the information flow from x to x_t . Typically, $q_{\phi}(z|x)$ is modeled as a fully factorized Gaussian distribution:

$$q_{\phi}(z|x) = \mathcal{N}(\mu, \text{diag}(\sigma^2)), \mu, \sigma^2 = f_{\phi}(x) \quad (5)$$

where f_{ϕ} is a ϕ -parameterized inference network, which outputs the mean μ and variance σ^2 of the Gaussian distribution.

The inference network aims to find the variational parameters for all sentences in the space of a Gaussian distribution. The performance is determined by two factors [3]: (i) the capacity of the Gaussian form to match the true posterior; and (ii) the ability of the inference network to produce good variational parameters for each sentence. In practice, the inference networks may always have a

capacity limit, which can result in imperfect variational parameters. In certain extreme cases, the limited capacity may produce all variational Gaussian parameters that share a similar configuration such as the prior, *i.e.*, KL vanishing. It is thus promising to relax the Gaussian approximation to more expressive forms.

2.2 From GMM to Semi-Implicit Posteriors

We seek to extend the Gaussian posterior in (5) to more flexible distribution via Gaussian Mixture Models (GMM) [8]:

$$q_\phi(\mathbf{z}|\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\boldsymbol{\mu}_k, \text{diag}(\boldsymbol{\sigma}_k^2)),$$

$$\text{with } \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\sigma}_k^2\}_{k=1}^K = f_\phi(\mathbf{x}), \quad (6)$$

where K is the number of Gaussian components, $\boldsymbol{\mu}_k$ and $\boldsymbol{\sigma}_k^2$ is the mean and variance of the k -th component, and $\boldsymbol{\pi} = [\pi_1, \dots, \pi_K]$ is the mixture coefficient. The traditional GMM is restrictive in that (i) K needs to be pre-defined, and (ii) the size of network ϕ for amortized inference will increase as K increases.

To bypass these issues, we propose the following *semi-implicit* (SI) form of GMM [26]:

$$q_\phi(\mathbf{z}|\mathbf{x}) = \int_{\boldsymbol{\lambda}} q(\mathbf{z}|\boldsymbol{\lambda}) q_\phi(\boldsymbol{\lambda}|\mathbf{x}) \quad (7)$$

where $\boldsymbol{\lambda} = [\boldsymbol{\mu}, \boldsymbol{\sigma}^2]$ is the mean and variance for Gaussian distribution $q(\mathbf{z}|\boldsymbol{\lambda}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2)$, whose distribution parameters are drawn from an *implicit* distribution $\boldsymbol{\lambda} \sim q_\phi(\boldsymbol{\lambda}|\mathbf{x})$. Specifically, the sampling process is:

$$\boldsymbol{\lambda} = f_\phi(\mathbf{x}, \boldsymbol{\epsilon}), \text{ with } \boldsymbol{\epsilon} \sim q_0(\boldsymbol{\epsilon}) \quad (8)$$

where $q_0(\boldsymbol{\epsilon})$ is an easy-to-sample distribution, f_ϕ is a ϕ -parameterized network with input $\{\mathbf{x}, \boldsymbol{\epsilon}\}$.

Note that (7) can be viewed as a GMM with an infinite number of components: once ϕ is well trained, one can draw an arbitrary number of samples via (8), each of which corresponds to a new Gaussian component. In contrast to (6), this sampling process only requires adding a small constant number of network parameters in ϕ .

We show the probabilistic graphical model of SI in Figure 1(b). Path A is now composed of $\{\mathbf{x}, \boldsymbol{\epsilon}\} \rightarrow \boldsymbol{\lambda} \rightarrow \mathbf{z} \rightarrow x_t$, where $q_\phi(\mathbf{z}|\mathbf{x})$ has a much more flexible distribution family to choose from than a Gaussian form. It allows \mathbf{z} to encode \mathbf{x} more easily. The enriched \mathbf{z} can strength path A, and help it stay informative in the competition.

2.3 Learning and Inference with SI

Though flexible, the SI representation in (7) introduces additional difficulty in learning: the evaluation of the KL term in (2) becomes inefficient because its closed form disappears. Fortunately, there is a upper bound for the KL term:

$$\mathcal{R} = \text{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) \quad (9)$$

$$= \text{KL}(\mathbb{E}_{q_\phi(\boldsymbol{\lambda}|\mathbf{x})} q(\mathbf{z}|\boldsymbol{\lambda}) || p(\mathbf{z}))$$

$$\leq \mathbb{E}_{q_\phi(\boldsymbol{\lambda}|\mathbf{x})} [\text{KL}(q(\mathbf{z}|\boldsymbol{\lambda}) || p(\mathbf{z}))] \triangleq \underline{\mathcal{R}}. \quad (10)$$

Note that (10) is efficient to evaluate due to the Gaussian form of $q(\mathbf{z}|\boldsymbol{\lambda})$. Therefore, we can maximize the lower bound of the ELBO:

$$\mathcal{L} = \mathcal{E} - \mathcal{R} \geq \mathcal{E} - \underline{\mathcal{R}} \triangleq \underline{\mathcal{L}} \quad (11)$$

However, directly optimizing (11) causes a *degeneracy* issue, characterized by $q_\phi(\boldsymbol{\lambda}|\mathbf{x})$ converging to a point mass density, making $q_\phi(\mathbf{z}|\mathbf{x})$ reduce to the vanilla Gaussian form. To prevent degeneracy, a repulsive term is added to regularize $\underline{\mathcal{L}}$ [26]:

$$\mathcal{B}_K = \mathbb{E}_{\boldsymbol{\lambda}, \boldsymbol{\lambda}^{(1)}, \dots, \boldsymbol{\lambda}^{(K)} \sim q_\phi(\boldsymbol{\lambda}|\mathbf{x})} \text{KL}(q(\mathbf{z}|\boldsymbol{\lambda}) || \tilde{q}_K(\mathbf{z}|\boldsymbol{\lambda}))$$

$$\text{where } \tilde{q}_K(\mathbf{z}|\boldsymbol{\lambda}) \triangleq \frac{q(\mathbf{z}|\boldsymbol{\lambda}) + \sum_{k=1}^K q(\mathbf{z}|\boldsymbol{\lambda}^{(k)})}{K+1} \quad (12)$$

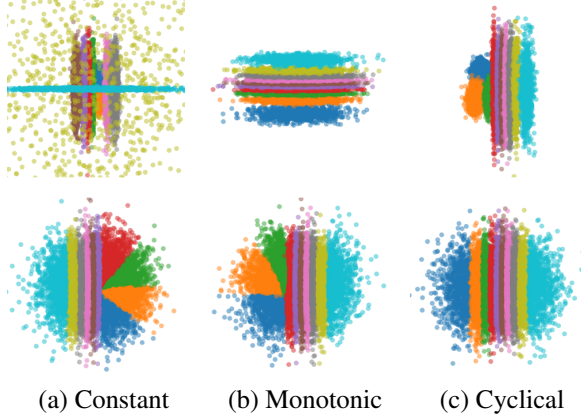


Figure 2: Comparison of VAE (top row) and SI (bottom row) in the learned latent spaces for the three schedules.

Methods	BLEU			BOW Embedding			Intra Distinct		Inter Distinct		L
	R	P	F1	A	E	G	dist-1	dist-2	dist-1	dist-2	
HRED	0.262	0.262	0.262	0.820	0.537	0.832	0.813	0.452	0.081	0.045	12.1
CVAE	0.295	0.258	0.275	0.836	0.572	0.846	0.803	0.415	0.112	0.102	12.4
CVAE+BOW	0.298	0.272	0.284	0.828	0.555	0.840	0.819	0.493	0.107	0.099	12.5
CVAE+CO	0.299	0.269	0.283	0.839	0.557	0.855	0.863	0.581	0.111	0.110	10.3
DialogWAE	0.394	0.254	0.309	0.897	0.627	0.887	0.713	0.651	0.245	0.413	15.5
DialogWAE+GMP	0.420	0.258	0.319	0.925	0.661	0.894	0.713	0.671	0.333	0.555	15.2
SI + M	0.376	0.248	0.299	0.888	0.640	0.890	0.832	0.461	0.209	0.298	11.5
CVAE + M	0.360	0.245	0.292	0.884	0.634	0.884	0.826	0.459	0.204	0.282	11.2
SI + C	0.441	0.246	0.315	0.937	0.683	0.905	0.888	0.782	0.499	0.77	12.1
CVAE + C	0.440	0.244	0.314	0.937	0.688	0.906	0.876	0.741	0.461	0.72	12.0

Table 1: Performance comparison on the SwitchBoard dataset (P: precision, R: recall, L: average length). **Blue** numbers indicate SI performs than VAE. **Bold** means the highest compared with all other methods.

This leads to $\mathcal{L}_K \triangleq \mathcal{L} + \mathcal{B}_K$. Note that maximizing \mathcal{L}_K with $K \geq 1$ would encourage positive \mathcal{B}_K and drive $q_\phi(\lambda|\mathbf{x})$ away from degeneracy. In other words, \mathcal{B}_K guarantees the SI in (7) to spread out as an infinite number of Gaussian mixtures, rather than collapsing into one single Gaussian.

In the Appendix A, we describe how SI mitigates the KL issue, with various training techniques, including β schedules (**Monotonic** and **Cyclical**), as well as **Aggressive** encoder training [9]. We also discuss related work in the Appendix B.

3 Experiments

We implement SI in Pytorch. We provide a detailed description of the settings, datasets, evaluation metrics in Appendix.

3.1 Visualization of Latent Spaces

To illustrate the learned \mathbf{z} of SI and VAE, in Figure 2 we visualize the latent space of a synthetic dataset of 10 sequences, where each color corresponds to $\mathbf{z} \sim q(\mathbf{z}|n)$, for $n = 1, \dots, 10$. The Gaussian VAE results are in the top row. The constant schedule produces heavily mixed latent codes \mathbf{z} for different sequences. The monotonic schedule divides the space into a mixture of 10 thin and long cluttered Gaussians. The cyclical schedule behaves similarly but with more circled distributions. This shows that the Gaussian VAE results heavily depend on the β -scheduling schemes. This is perhaps because the learning process is conducted in the restrictive space of Gaussian variational forms, and hence carefully-crafted searching schedules play an important role. However, the proposed SI method consistently produces well-divided latent representations for all three schedules. Importantly, its posteriors $q(\mathbf{z}|n)$ have non-Gaussian forms such as slice-shaped distributions.

Methods		AU \uparrow	MI \uparrow	KL \uparrow	PPL \downarrow	IWP \downarrow
M	VAE	1	0.56	0.68	94.6	93.1
	SA	-	-	1.05	109.2	-
	SI	7	2.73	4.83	105.6	96.73
C	VAE	4	0.95	1.26	95.74	96.23
	SA	-	-	2.26	108.7	-
	SI	5	4.08	2.63	106.3	96.87
M A	VAE	14	3.33	9.02	101.9	92.26
	SI	32	3.73	9.20	100.3	88.51

Table 2: Results on PTB. The SA-VAE results are from [5]. **Blue** numbers indicate SI performs than VAE. **Bold** means the highest compared with all other methods.

3.2 Language Modeling

The results on Penn Tree Bank (PTB) dataset [15] are reported in Table 2. We first compare SI with the VAE trained with a monotonic and cyclical β schedule. We see that SI achieves higher AU, MI, and KL than VAE for both schedules. This means SI can provide more informative latent codes than VAE. We compare with the semi-amortized (SA) training [11]. SI outperforms SA in terms of KL while maintaining the same PPL. This implies that learning with more flexible distribution forms yields more informative latent codes. When training VAE and SI with the aggressive encoder training, both methods are improved. SI still achieves substantially better results than VAE in terms of AU. Interestingly, it activates all 32 latent units. We argue that the aggressive encoder training can be particularly important for the proposed SI, as it allows SI getting fully optimized to reach its representational power.

3.3 Dialogue Response Generation

Table 1 summarizes the results for the various methods. All baseline results are from [8]. We compare SI with Gaussian CVAE for both monotonic and cyclical β schedules. SI improves CVAE in terms of all evaluation metrics for the monotonic schedule. When the cyclical schedule is used, SI provides higher BLEU scores than CVAE, indicating that SI is able to generate more relevant responses. This also implies that SI is able to learn more informative z , and suffers less from the KL vanishing issue. When comparing with the state-of-the-art methods, SI provides the highest BLEU Recall and BOW Embedding (A). Meanwhile, SI achieves the best Intra/Inter Distinct values, showing evidence that SI can produce the most diverse response generation.

Table 4 in Appendix shows examples of generated responses from the CVAE and SI. Given a context, five samples are drawn for each method. SI-CVAE can generate more coherent responses that cover multiple plausible aspects.

4 Conclusion

We have introduced a semi-implicit approximation in the posterior learning of VAEs to infer text representations. The flexibility of SI helps reduce the representational bottleneck in the latent space and thus alleviates the KL vanishing issue. The effectiveness of SI is validated with the clear latent space division in the synthetic dataset and improved performance on two NLP tasks: providing more informative latent codes in VAE language modeling, and promoting the diversity in conditioned dialogue response generation.

References

- [1] Samuel R Bowman, Luke Vilnis, Oriol Vinyals, Andrew M Dai, Rafal Jozefowicz, and Samy Bengio. Generating sentences from a continuous space. *ACL*, 2016.
- [2] Yuri Burda, Roger Grosse, and Ruslan Salakhutdinov. Importance weighted autoencoders. *arXiv preprint arXiv:1509.00519*, 2015.
- [3] Chris Cremer, Xuechen Li, and David Duvenaud. Inference suboptimality in variational autoencoders. *ICML*, 2018.
- [4] Yihao Feng, Dilin Wang, and Qiang Liu. Learning to draw samples with amortized stein variational gradient descent. *UAI*, 2017.

- [5] Hao Fu, Chunyuan Li, Xiaodong Liu, Jianfeng Gao, Asli Celikyilmaz, and Lawrence Carin. Cyclical annealing schedule: A simple approach to mitigating KL vanishing. *NAACL*, 2019.
- [6] John J Godfrey and Edward Holliman. Switchboard-1 release 2. *Linguistic Data Consortium, Philadelphia*, 926:927, 1997.
- [7] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*, volume 1. MIT press Cambridge, 2016.
- [8] Xiaodong Gu, Kyunghyun Cho, Jungwoo Ha, and Sunghun Kim. DialogWAE: Multimodal response generation with conditional Wasserstein auto-encoder. *ICLR*, 2019.
- [9] Junxian He, Daniel Spokoyny, Graham Neubig, and Taylor Berg-Kirkpatrick. Lagging inference networks and posterior collapse in variational autoencoders. *ICLR*, 2019.
- [10] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. Beta-VAE: Learning basic visual concepts with a constrained variational framework. In *ICLR*, 2017.
- [11] Yoon Kim, Sam Wiseman, Andrew C Miller, David Sontag, and Alexander M Rush. Semi-amortized variational autoencoders. *ICML*, 2018.
- [12] Yoon Kim, Sam Wiseman, and Alexander M Rush. A tutorial on deep latent variable models of natural language. *arXiv preprint arXiv:1812.06834*, 2018.
- [13] Diederik P Kingma and Max Welling. Auto-encoding variational Bayes. *ICLR*, 2013.
- [14] Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. Adversarial autoencoders. *ICLR workshop*, 2016.
- [15] Mitchell P Marcus, Mary Ann Marcinkiewicz, and Beatrice Santorini. Building a large annotated corpus of english: The penn treebank. *Computational linguistics*, 1993.
- [16] Lars Mescheder, Sebastian Nowozin, and Andreas Geiger. Adversarial variational bayes: Unifying variational autoencoders and generative adversarial networks. In *ICML*, 2017.
- [17] Yishu Miao, Lei Yu, and Phil Blunsom. Neural variational inference for text processing. In *ICML*, 2016.
- [18] Tomáš Mikolov, Martin Karafiát, Lukáš Burget, Jan Černocký, and Sanjeev Khudanpur. Recurrent neural network based language model. In *Eleventh Annual Conference of the International Speech Communication Association*, 2010.
- [19] Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representation. In *EMNLP*, 2014.
- [20] Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows. *arXiv preprint arXiv:1505.05770*, 2015.
- [21] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. *ICML*, 2014.
- [22] Iulian Vlad Serban, Alessandro Sordoni, Yoshua Bengio, Aaron C Courville, and Joelle Pineau. Building end-to-end dialogue systems using generative hierarchical neural network models. In *AAAI*, 2016.
- [23] Iulian Vlad Serban, Alessandro Sordoni, Ryan Lowe, Laurent Charlin, Joelle Pineau, Aaron Courville, and Yoshua Bengio. A hierarchical latent variable encoder-decoder model for generating dialogues. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- [24] Xiaoyu Shen, Hui Su, Shuzi Niu, and Vera Demberg. Improving variational encoder-decoders in dialogue generation. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [25] Tsung-Hsien Wen, Yishu Miao, Phil Blunsom, and Steve Young. Latent intention dialogue models. *ICML*, 2017.

- [26] Mingzhang Yin and Mingyuan Zhou. Semi-implicit variational inference. *ICML*, 2018.
- [27] Tiancheng Zhao, Ran Zhao, and Maxine Eskenazi. Learning discourse-level diversity for neural dialog models using conditional variational autoencoders. *ACL*, 2017.

Table 3: Comparisons of vanilla VAE and SI.

	VAE	SI
K	$K = 0$	$K \geq 1$
Speed	Fast	Larger K yields higher compute cost
Flexibility	Limited	Larger K yields higher flexibility

A Mitigating KL Vanishing with SI

The KL vanishing issue happens due to a lack of good latent codes \mathbf{z} in the initial stages of decoder training [9, 5], and the decoder passively chooses to use the previous word tokens $x_{<t}$ alone to generate the next word x_t . Based on this observation, two algorithms are proposed independently:

- **Aggressive encoder training:** He et al. [9] propose to aggressively optimize the encoder T_a times ($T_a > 0$) before performing each model update.
- **Cyclical β -annealing schedule:** Fu et al. [5] propose simply repeats the monotonic β annealing procedure M times ($M > 1$). Within each cycle, only a proportion $R \in (0, 1)$ is used to increase β from 0 to 1, and while the $1 - R$ for fixing $\beta = 1$. A linear annealing scheme is:

$$\beta_t = \begin{cases} \tau/R, & \tau \leq R \\ 1, & \tau > R \end{cases} \quad \text{with } \tau = \frac{\text{mod}(t-1, \lceil T/M \rceil)}{T/M}. \quad (13)$$

Both works share a similar motivation: improving the quality of \mathbf{z} before it is used for decoder training. However, they optimize the latent codes in the space of Gaussian distribution, due to the variational assumption.

In this paper, we propose to further mitigate the KL issue through improving the latent codes using SI. The update of SI is based on the tractable bound:

$$\underline{\mathcal{L}}_K = \underline{\mathcal{L}} + \mathcal{B}_K = \mathcal{E} + \mathcal{R} + \mathcal{B}_K \quad (14)$$

$$\begin{aligned} &= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z}) + p(\mathbf{z}) - \tilde{q}_K(\mathbf{z}|\boldsymbol{\lambda})]. \\ &\approx \frac{1}{J} \sum_{j=1}^J \left[\log p_\theta(\mathbf{x}|\mathbf{z}_j) + \log p(\mathbf{z}_j) - \log \frac{1}{K+1} \left[\sum_{k=1}^K q(\mathbf{z}_j|\boldsymbol{\lambda}_k) + q(\mathbf{z}_j|\boldsymbol{\lambda}_j) \right] \right], \end{aligned} \quad (15)$$

where the expectation in (14) is estimated via J samples in (15), with the j -th sample drawn as:

$$\mathbf{z}_j \sim q(\mathbf{z}|\boldsymbol{\lambda}_j), \quad \boldsymbol{\lambda}_j \sim q_\phi(\boldsymbol{\lambda}|\mathbf{x}) \quad (16)$$

On the impact of K The representational power of SI relies on K , the number of samples used to approximate (7). Each sample is used to construct the variational parameters for one Gaussian mixture, thus larger K indicates more Gaussian mixtures. It has been shown in [26] that the regularized lower bound $\underline{\mathcal{L}}_K$ is an asymptotically exact ELBO that satisfies $\underline{\mathcal{L}}_0 = \underline{\mathcal{L}}$ and $\lim_{K \rightarrow \infty} \underline{\mathcal{L}}_K = \mathcal{L}$.

We summarize the connection/comparison of vanilla VAE and SI in Table 3. Note that vanilla VAE can be considered as a special case of SI-VAE when $K = 0$. SI can balance the trade-off between the posterior flexibility and the computational cost using K . To leverage both advantages, we suggest an *annealing training* procedure via scheduling K . It performs in two stages: (i) A small K is first used. We can leverage the speed advantage to push the variational distribution towards the target distribution at a fast pace. (ii) We gradually increase K to a pre-defined K_{\max} . It refines the distribution learned in the first stage into highly flexible SI representations.

This annealing training procedure can be integrated into the aggressive encoder training [9] or the cyclical schedule [5] to get better performance, while maintaining low computational cost. Note the key to reduce KL vanishing is to feed high-quality \mathbf{z} to decode in the initial stage of decode learning. Here, this initial stage means the time after each aggressive training loop ends in [9], or each cyclical period ends in [5]. Hence, the SI annealing training periods can be synchronized with them. By the time the next aggressive loop or cyclical period start, we can make sure the flexible SI latent codes are fed into the decoder.

We summarize the full training in Algorithm 1. K_t is iteration-dependent; it is reserved to implement the stacked training. Note that when $M = 1$ and $R < 0.5$, we have the standard monotonical annealing scheme. When $T_a = 0$ we have the standard encoder-decoder update.

Algorithm 1: SI Training

Input : encoder ϕ , decoder θ , K_t , T_a , M , R

```
1 Initialize:  $\phi, \theta$ ;  
2 for  $t = 1, 2, \dots, T$  do  
3   %  $\beta$  annealing schedule  
4   Evaluate  $\beta_t$  using (13) ;  
5   % Aggressive encoder training  
6   for  $t_a = 1, 2, \dots, T_a$  do  
7     Sample a minibatch  $\mathbf{X}$  ;  
8     Compute gradient  $g_\phi = \nabla_{\phi} \mathcal{L}_{\beta_t, K_t}(\mathbf{X})$  ;  
9     Update  $\phi$  using gradient  $g_\phi$  ;  
10  end  
11  % SI update  
12  Sample a minibatch  $\mathbf{X}$  ;  
13  Sample  $\lambda_k$  using (8) for  $k = 1, \dots, K_t$  ;  
14  Sample  $z_j$  using (16) for  $j = 1, \dots, J$  ;  
15  Compute gradient  $g_{\phi, \theta} = \nabla_{\phi, \theta} \mathcal{L}_{\beta_t, K_t}(\mathbf{X})$  ;  
16  Update  $\theta, \phi$  using gradient  $g_{\phi, \theta}$  ;  
17 end
```

B Related Work

Good latent representations are key for a wide range of NLP tasks. We recommend [12] for a comprehensive study for deep latent variable models in natural languages. Several VAE model variants have been developed to reduce the approximate errors in Gaussian proposals via more flexible posteriors, including normalizing flows (NFs) [20], adversarial variational Bayes (AVB) [16], particle-based methods [4]. Our method is based on [26]. It bypasses the strict invertibility requirement of NFs, the instability of adversarial training in AVB, and prohibitive computation in particle-based methods. Nevertheless, all these methods can provide more accurate posterior approximate than the Gaussian variants when carefully tuned. However, to the best of our knowledge, they are only evaluated on image datasets with MLP/CNN decoders. In the text domain, the use of auto-regressive decoders in VAE brings additional difficulties in learning representations. Our paper presents the first work to investigate the flexible posterior learning for auto-regressive decoders. We gain the insights that the KL vanishing issue is less severe than previously thought when trained with more flexible posteriors such as SI.

C Visualization

Each sequence is a 10-dimensional one-hot vector with the value 1 appearing in different positions. A 2-dimensional latent space is used for the convenience of visualization. A 1-layer LSTM with hidden units is employed for the decoder in analog with the sentence decoding process. The encoder is implemented using a 2-layer MLP with 64 hidden. We use $T=4K$ total iterations, $K=J=100$. Three β -scheduling schemes are considered [5].

We refer $q(\mathbf{z}) = \sum_{n=1}^N q(\mathbf{z}|n)q(n)$ as the aggregated posterior [14]. This marginal distribution characterizes the aggregated \mathbf{z} after embedding the entire dataset into the latent space. Good latent representations should have low $\text{KL}(q(\mathbf{z})||p(\mathbf{z}))$. Ideally, it means the overall shape of \mathbf{z} samples is close to $p(\mathbf{z}) = \mathcal{N}(0, 1)$. The flexibility of SI allows the learned posteriors to adapt arbitrarily well to fit the prior, thus produces high-quality $q(\mathbf{z})$. In contrast, the limitation of Gaussian assumptions in vanilla VAE leads to some holes between $q(\mathbf{z}|n)$ and $p(\mathbf{z})$, making $q(\mathbf{z})$ violate the constraint of $p(\mathbf{z})$. The SI representations capture more clear patterns and structured information is captured in \mathbf{z} , which is beneficial in downstream applications below.

D Language Modeling

Dataset and Setup We first consider applying VAEs to the language modeling task on the Penn Tree Bank (PTB) dataset [15]. Following [11, 9], we set the word embedding dimension as 512. The encoder and decoder are LSTMs with 1024 hidden units, respectively. $J=5$, K is linearly increased to $K_{\max} = 50$.

Context mind of the jury is probably a better vehicle for sentencing than the judge himself because then there's the consensus thing rather than a single person	
Target well there is an argument in that a judge has seen a lot more cases and you know perhaps he sees one crime	
SI	CVAE
1. yeah that is true	1. yeah
2. that's a very good idea i guess is not even the jury system	2. that's right
3. i think the judges should be a little more like a man	3. um - hum
4. uh - huh	4. oh i can't believe that
5. right i should be able to get on the jury that's not necessarily the way	5. uh - huh yeah i think that's a good idea but i'm not sure about it

Table 4: Generated dialog responses from SI and CVAE methods.

Evaluation In addition to generating high-quality sentences as in the traditional language models that only, VAEs also aim to learn a good posterior distribution in the latent space. The language modeling performance is evaluated with ELBO, perplexity (PPL) and importance weighted perplexity (IWP) [9]. Note that IWP provides a tighter bound to $\log p(\mathbf{x})$ than PPL. Higher ELBO and lower PPL/IWP indicate the model fits the observed sentences better. More importantly, we are interested in the learned \mathbf{z} , which is evaluated using the following three metrics:

- **AU**: The total number of active units in \mathbf{z} , defined as $A_z = \text{Cov}_{\mathbf{x}}(\mathbb{E}_{z \sim q(\mathbf{z}|\mathbf{x})}[\mathbf{z}]) > 0.01$ [2];
- **MI**: The mutual information $I(\mathbf{x}, \mathbf{z})$;
- **KL**: The posterior-prior KL divergence, *i.e.*, \mathcal{R} defined in (2). The KL for SI is reported as its bound $\underline{\mathcal{R}}_K + \mathcal{B}_K$.

E Conditional VAE for Dialog

E.1 Model Description

Each conversation can be represented via three random variables: the dialog context \mathbf{c} composed of the dialog history, the response utterance \mathbf{x} , and a latent variable \mathbf{z} , which is used to capture the latent distribution over the valid responses ($\beta = 1$) [27]. The ELBO can be written as:

$$\begin{aligned} \log p_{\theta}(\mathbf{x}|\mathbf{c}) &\geq \mathcal{L}_{\text{ELBO}} \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{c})} [\log p_{\theta}(\mathbf{x}|\mathbf{z}, \mathbf{c})] - \beta \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{c}) || p(\mathbf{z}|\mathbf{c})) \end{aligned} \quad (17)$$

Dataset and Setup The conditional VAE (CVAE) is often used in dialogue response generation, the learned representations are crucial in generating relevant and diverse responses. We evaluate our model on Switchboard dataset [6], which contains 2,400 two-way telephone conversations under 70 specified topics. We use SI to replace the Gaussian posteriors in CVAE. Our model architecture and hyper-parameters are identical to [27] except that we concatenate the context and sentence embedding with 100-dim noise vector as the input of the SI inference network. $J=5$, K is linearly increased to $K_{\text{max}} = 30$.

Baselines We compare the performance of SI-CVAE with six recently-proposed baselines for dialogue modeling: (1) HRED: a generalized sequence-to-sequence model with hierarchical RNN encoder [22]; (2) CVAE: a conditional VAE model with KL-annealing [23] [27]; (3) CVAE-BOW: a conditional VAE model with a BOW loss [27]; (4) CVAE-CO: a collaborative conditional VAE model [24]; (5) DialogWAE: a conditional Wasserstein autoencoder model [8] and (6) DialogWAE-GMP: a conditional Wasserstein autoencoder model with Gaussian mixture prior [8].

E.2 Evaluation

We adopted the same metrics from [8] to quantify the performance for dialog response generation:

- **BLEU**: BLEU measures how much a generated response contains n -gram overlaps with the reference.
- **BOW Embedding**: The cosine similarity of bag-of-words embeddings between the hypothesis and the reference. The Glove vectors [19] are used. For each test context, we report the maximum BOW embedding score among the 10 sampled responses. We use three metrics to compute the word embedding similarity: **Average**, **Extrema** and **Greedy**.
- **Intra & Inter Distinct**: $\text{dist-}n$ is defined as the ratio of unique n -grams ($n = 1, 2$) over all n -grams in the generated responses. Note that intra-dist and inter-dist are defined as the average of distinct values within each sampled response and among all sampled responses, respectively. Higher values indicate higher diversity.