

Re-balancing Variational Autoencoder Loss for Molecule Sequence Generation

Chaochao Yan, Sheng Wang, Jinyu Yang, Tingyang Xu, Junzhou Huang *

University of Texas at Arlington
Tencent AI Lab

Abstract

Molecule generation is to design new molecules with specific chemical properties and further to optimize the desired chemical properties. Following previous work, we encode molecules into continuous vectors in the latent space and then decode the vectors into molecules under the variational autoencoder (VAE) framework. We investigate the posterior collapse problem of current RNN-based VAEs for molecule sequence generation. For the first time, we find that underestimated reconstruction loss leads to posterior collapse, and provide both theoretical and experimental evidence. We propose an effective and efficient solution to fix the problem and avoid posterior collapse. Without bells and whistles, our method achieves SOTA reconstruction accuracy and competitive validity on ZINC 250K dataset. When generating 10,000 unique valid SMILES from random prior sampling, it costs JT-VAE 1450s while our method only needs 9s. Our implementation will be made public. Our implementation is at <https://github.com/chaoyan1037/Re-balanced-VAE>.

Discovering new molecules that have desired target properties is the key challenge of drug and material design. This can be considered as an optimization problem, and the goal is to search for molecules with the best desired property score (Gómez-Bombarelli et al. 2018). However, exhaustive exploration in the molecule space is infeasible, as the number of estimated drug-like molecules is in the order of 10^{60} (Polishchuk, Madzhidov, and Varnek 2013). Additionally, molecule synthesis and validation are time-consuming and expensive in practice.

The majority of molecule generation methods heavily rely on the variational autoencoder (VAE) which is a combination of a deep latent variable model and an accompanying variational learning technique (Kingma and Welling 2013) (Rezende, Mohamed, and Wierstra 2014). As shown in Figure 1, drug molecules can be first embedded by the encoder into the continuous latent space which can be further utilized for property prediction and optimization. After that, the decoder maps a continuous latent vector to reconstruct the input molecule. Thanks to the clustering ability of VAE, the latent representations of semantically similar molecules (with similar chemical structures and properties) are grouped to-

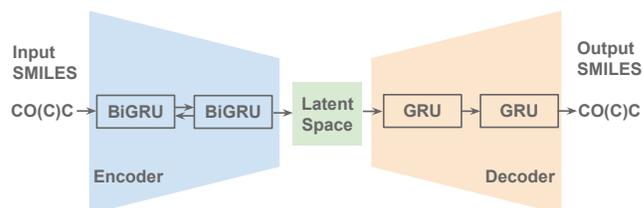


Figure 1: Overview of our VAE model. The encoder and decoder are built upon the bi-directional GRU and uni-directional GRU, respectively. Both input and output are SMILES sequences.

gether in latent space. In consequence, it allows semantically meaningful sampling and smooth interpolation in the latent space. Therefore, new molecules can be generated by randomly sampling from the prior and can be further optimized by exploring the latent space. The key idea behind the optimization above is to search for molecules that maximize an property score objective, given molecules' latent representation as input (Gómez-Bombarelli et al. 2018).

However, existing VAE models mainly suffer from the posterior collapse issue, where the decoder tends to ignore the latent vectors (Bowman et al. 2016) (Gómez-Bombarelli et al. 2018). This problem is more frequently observed in those models with RNN-based backbone (He et al. 2019). As a consequence, the generated molecules tend to be in low diversity and are weakly relevant to the latent vectors (Gómez-Bombarelli et al. 2018) (Kusner, Paige, and Hernández-Lobato 2017). This phenomenon has also been observed in Natural Language Processing (NLP) tasks, such as text generation (Bowman et al. 2016). To alleviate this problem, the major focus of the previous studies is to adopt various training strategies, such as KL cost annealing (Bowman et al. 2016) or aggressively optimizing the decoder before each encoder update (He et al. 2019). However, such methods can not be simply extended to molecule generation, mainly stemming from the fact that the molecule sequences are strictly structured and any mutations can result in invalid sequences. Motivated by the success of parse trees in the NLP field and attribute grammars in compiler design, recent

*Corresponding author: J. Huang, jzhuang@uta.edu

work (Kusner, Paige, and Hernández-Lobato 2017) (Dai et al. 2018) incorporate grammar or syntax rules to guarantee that syntactically valid SMILES sequences can be generated. As an alternative, a molecule can also be represented by a graph in order to mitigate the posterior collapse (Li et al. 2018) (Jin, Barzilay, and Jaakkola 2018).

Inspired by the essential pitfalls of the contemporary RNN-based VAE models in the molecule generation, here, we propose a new method to alleviate the posterior collapse issue. To achieve this goal, we first analyze the posterior collapse of vanilla VAE model for SMILES sequence generation. For the first time, we find that the posterior collapse is largely triggered by the underestimated reconstruction loss. We, therefore, propose to use a novel loss function to leverage the trade-off between the reconstruction loss and the KL loss in the VAE training. Without making any changes on the VAE network structures or introducing additional computational complexity, our method is extremely simple yet effective in preventing posterior collapse. We also provide the theoretical analysis of our method, and empirically demonstrate its state-of-the-art (SOTA) reconstruction accuracy and competitive validity score on the ZINC 250K dataset. Our primary contributions can be summarized as:

- We diagnose the main reason causing the posterior collapse within the RNN-based VAE model for molecule generation, with both theoretical and intuitive analyses been provided.
- We propose an effective and efficient method to eliminate the posterior collapse in VAE by leveraging the associations between the reconstruction loss and the KL loss.
- Extensive empirical studies demonstrate our method’s superiority over SOTA molecule generation approaches on the ZINC 250K dataset.

Background Information

The Variational Autoencoder

The VAE (Kingma and Welling 2013) (Rezende, Mohamed, and Wierstra 2014) is a specially regularized version of the standard autoencoder (AE). It is appealing because it can learn complex distribution in an unsupervised manner and later act as a generative model defined by a prior $p(z)$ and a conditional distribution $p_\theta(x|z)$. Since the true data likelihood is usually intractable, so the VAE instead optimizes an evidence lower bound (ELBO) which is a valid lower bound of the true data log likelihood:

$$\begin{aligned} \mathcal{L}(x; \theta, \phi) &= \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)] - D_{\text{KL}}(q_\phi(z|x)||p(z)) \\ &\leq \log p(x). \end{aligned} \quad (1)$$

where the encoder $q_\phi(z|x)$ is parameterized with ϕ and learns to map the input x to a variational distribution, and the decoder $p_\theta(x|z)$ parameterized with θ tries to reconstruct the input x given the latent vector z from the learned distribution. Usually, $q_\phi(z|x)$ is modeled as a Gaussian distribution and optimized to approximate the true posterior $p_\theta(z|x)$.

The VAE is optimized to maximize ELBO (1), where (i) negative reconstruction loss $\mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)]$ enforces the encoder to generate meaningful latent vector z , so that the decoder can reconstruct the input x from the z , and (ii) the KL regularization loss $D_{\text{KL}}(q_\phi(z|x)||p(z))$ minimizes the KL divergence between the approximate posterior $q_\phi(z|x)$ and the prior $p(z) \sim \mathcal{N}(0, \mathbf{I})$.

Related Work

Text Generation with VAE

Motivated by the ubiquitous posterior collapse problems observed in VAE-based models for text generation, various methods have been proposed and investigated recently (Bowman et al. 2016) (Yang et al. 2017) (Higgins et al. 2017) (Kim et al. 2018). Bowman et al. (Bowman et al. 2016) propose to anneal the KL loss weight to enable the model to learn meaningful encoding before applying KL loss to cluster the encoding. They also weaken the decoder with word dropout and historyless decoding to force the decoder rely on the latent vectors. (He et al. 2019) concludes that the posterior collapse mainly attributes to the lagging encoder network’s inability in approximating the true posterior. To overcome this limitation, they propose a novel training strategy which aggressively optimize the encoder network. Inspired by (Bowman et al. 2016), (Hao Fu 2019) proposes a cyclical annealing strategy which repetitively starts training from a pretrained model resulted in the previous cycle. They claim this procedure can make the model learn more meaningful latent representations progressively.

Molecule Generation

Thanks to the development of NLP text generation, the VAE model is applied for molecule generation for the first time in CVAE (Gómez-Bombarelli et al. 2018). They build a VAE encoder and decoder with GRU layers, representing molecules in the SMILES sequences. However, their model suffers from generating invalid SMILES sequences which makes their model impracticable. To improve the prior validity, context-free grammars for SMILES are introduced in GVAE (Kusner, Paige, and Hernández-Lobato 2017) to represent a molecule in the sparse tree. However, the validity score is still unsatisfactory. Inspired by this method, Syntax-directed VAE (SD-VAE) (Dai et al. 2018) incorporates extra semantic rules to ensure generated SMILES valid, and it achieves the best performance among all SMILES-based methods. However, these models did not solve model posterior collapse problem and there is a large space to improve.

Except for SMILES representations, molecules can also be represented in graph. (Li et al. 2018) employs graph-structured representations for molecules and models the probabilistic dependencies among a graphs nodes and edge with graph neural networks. Molecules are generated node by node by their model. Chemical sub-graphs instead of atoms are used as the basic building blocks in JT-VAE (Jin, Barzilay, and Jaakkola 2018), their methods can incrementally generate molecules to ensure chemical validity at each step. JT-VAE is the SOTA model for molecule generation.

Problem and Solution

Posterior Collapse

Prior work (Bowman et al. 2016) (Yang et al. 2017) (Higgins et al. 2017) (Kim et al. 2018) on NLP text generation has observed the posterior collapse phenomenon, in which the decoder tends to ignore z when training the VAE model. When posterior collapse happens, the model training falls into the local optimum of the ELBO objective (1), in which the variational posterior $q_\phi(z|x)$ naively mimics the model prior $p(z)$. Note that the KL loss in ELBO can be further decomposed (Hoffman and Johnson 2016) as:

$$\mathbb{E}_{p_d(x)}[D_{\text{KL}}(q_\phi(z|x)||p(z))] = I_q + D_{\text{KL}}(q_\phi(z)||p(z)), \quad (2)$$

where I_q is the mutual information between x and z given $q_\phi(z|x)$, and $p_d(x)$ is empirical data distribution. When posterior collapse occurs, the KL loss decreases nearly to zero so that I_q is also close to zero (both items on the right-hand side in (2) are non-negative) during the VAE model training process. It is especially evident when modelling discrete data with a strong auto-regressive network such as LSTM (Hochreiter and Schmidhuber 1997) and GRU (Chung et al. 2014), which is exactly our case. This is undesirable since the VAE model fails to learn meaningful latent representations for input sequences.

For NLP text generation task, the posterior collapse problem has been mainly attributed to the low quality of latent representations z at the early stage of model training (Bowman et al. 2016) (He et al. 2019) (Hao Fu 2019). To be more specific, the decoder $p_\theta(x|z)$ falls behind the encoder $q_\phi(z|x)$ at the initial training procedure, and $q_\phi(z|x)$ generates low-quality latent representations so that it is very hard for $p_\theta(x|z)$ to recover the input sequences. In consequence, the model is forced to ignore z . Many solutions have been proposed to solve the problem and they have demonstrated satisfactory improvement on NLP datasets.

However, the molecule generation is a quite different scenario though it appears to be same as the NLP text generation. First of all, its token size is far more less than the NLP text generation. The token size for NLP text is usually tens of thousands or even more, while it is less than 100 for chemical molecule data. The smaller token size makes the molecule reconstruction task much easier. Second, the molecule sequence is composed strictly following the SMILES grammar or syntax rules, and the reconstructed sequence must be exactly the same as the input to be matched. Any token mutations can result in a completely different sequence. However, there are no rigid grammar rules applied to the NLP text and exact match is not required.

We have found existing solutions (He et al. 2019) (Hao Fu 2019) to posterior collapse in NLP text generation does not work well for chemical molecule generation. This motivates us to propose such a solution for molecule generation.

The Problem in Previous Solutions

To avoid posterior collapse, which will cause a VAE losing reconstruction ability, previous SMILES-based meth-

ods CVAE (Gómez-Bombarelli et al. 2018), GVAE (Kusner, Paige, and Hernández-Lobato 2017), and SD-VAE (Dai et al. 2018) reduce the standard deviation σ of prior Gaussian distribution to a small value 0.01 (can be found in their public implementation CVAE^{1,2}, GVAE³, SD-VAE⁴), which makes their models more like AEs instead of VAEs. That is why CVAE and GVAE have a decent reconstruction accuracy but extremely low validity scores as shown in Table 1. If we set the $\sigma=1$, all these three models will suffer from model posterior collapse and lose the reconstruct ability (similar to the vanilla VAE in Figure 2(e)). In following our analysis and experiments, we strictly keep the $\sigma=1$.

Underestimated Reconstruction Loss

To investigate the cause of posterior collapse within the VAE for molecule generation, we conduct convincing analysis and investigation into posterior collapse. We hypothesize it is the underestimated reconstruction loss that causes posterior collapse of a VAE during training process. Both theoretical analysis and experimental support are provided.

From the perspective of theory, reconstruction loss term $\mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)]$ measures the reconstruction ability of the decoder given latent vector z . The decoder should only receive information from z and tries to reconstruct the full sequence accurately from the given starting z . However, in practice RNN model is usually optimized with teacher forcing (Williams and Zipser 1989), in which the current input is the ground truth instead of prediction from a prior time step.

We can rewrite the reconstruction loss term in (1) as:

$$\mathbb{E}_{q_\phi(z|x)}\left[\sum_{t=1}^T \log p_\theta(x_t|z, \tilde{x}_{<t})\right], \quad (3)$$

where the T is the maximum time step, $\tilde{x}_{<t}$ is the prediction prefix before time t and the current input is the previous time step output \tilde{x}_{t-1} , and \tilde{x}_0 is the start symbol.

With teacher forcing, the actual reconstruction loss is:

$$\mathbb{E}_{q_\phi(z|x)}\left[\sum_{t=1}^T \log p_\theta(x_t|z, \tilde{x}_{<t}, x_{<t})\right], \quad (4)$$

where $x_{<t}$ is the ground-truth prefix before time t and the ground-truth token of previous time step x_{t-1} is the RNN input, and x_0 is also the start symbol.

Since the ground-truth information is incorporated additionally in (4) when training the VAE, which can make the prediction easier since the ground-truth prefix is given, we can expect that the reconstruction ability of decoder is largely overestimated compared with (3). Therefore, we can assume the reconstruction loss term is underestimated, which will potentially breaks the balance between reconstruction loss and KL loss in (1). We will verify the assumption and also demonstrate quantitatively how much the reconstruction loss is underestimated in the experiment section. Let us agree on the claim for now.

¹https://github.com/aspuru-guzik-group/chemical_vae

²<https://github.com/HIPS/molecule-autoencoder>

³<https://github.com/mkusner/grammarVAE>

⁴<https://github.com/Hanjun-Dai/sdvae>

Re-balanced VAE Loss

Since reconstruction loss is underestimated, and it breaks the balance with KL loss, which leads to the posterior collapse finally. We can recover the balance by applying a reconstruction loss weight α to the ELBO (1):

$$\mathcal{L}(x; \theta, \phi) = \alpha \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)] - D_{\text{KL}}(q_\phi(z|x)||p(z)), \alpha > 1, \quad (5)$$

where α can be estimated using Monte Carlo in every training iteration. Specifically, we can sample a batch of data as input and run a VAE with/without teacher forcing, respectively. Since the reconstruction loss without teacher forcing can be regarded as the "true" reconstruction loss, we approximate α as the ratio of reconstruction loss without teacher forcing to that with teacher forcing. However, estimating α in every training iteration is too expensive. We can set α as a hyper-parameter for simplicity.

Inspired by the β -VAE (Higgins et al. 2017) formulation, we can instead reduce KL loss weight β , which is equivalent to increasing reconstruction loss weight α . It is more natural and convenient to search for the optimal value of hyper-parameter β since increasing β from 0 is a gradual transition from AE to VAE. So we can have a similarly modified VAE formulation:

$$\mathcal{L}(x; \theta, \phi) = \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)] - \beta D_{\text{KL}}(q_\phi(z|x)||p(z)), 0 \leq \beta < 1. \quad (6)$$

Note that in our case $\beta < 1$, while β -VAE requires the KL weight $\beta > 1$. β -VAE is proposed in (Higgins et al. 2017) to learn disentangled representation of generative factors by enforcing a larger penalty on KL loss, since they postulate that $\beta > 1$ could place a stronger constraint on the latent representation to drive the VAE to learn more efficient latent representation of input x . While we have a completely different motivation and goal of fixing imbalanced VAE loss by reducing KL weight since we find reconstruction loss is underestimated in ELBO (1).

Except for theoretical analysis, our method can also be explained from an intuitive perspective. In previous methods CVAE, GVAE, and SD-VAE, when sampling latent vectors z they have to reduce the standard deviation σ to a small value 0.01 otherwise the model will collapse and lose the reconstruct ability. Instead of reducing sampling σ , we can anneal the KL loss weight β to make the model transform from AE to VAE gradually (Bowman et al. 2016). Different from (Bowman et al. 2016), we restrict β to be smaller than 1. By searching for the optimal β , we can arrive a trade-off between the reconstruction accuracy and validity score.

We acknowledge that previous methods have empirically tried to reduce the KL loss weight to avoid the posterior collapse (Dai et al. 2018) (He et al. 2019) (Hao Fu 2019). β -VAE ($\beta = 0.4$) alleviates the problem and achieves competitive performance on density estimation for NLP text datasets (He et al. 2019), which proves that reducing β is viable for NLP text task. It is also indicated that setting $\beta = 1/\text{LatentDimension}$ could lead to better results (Kusner, Paige, and Hernández-Lobato 2017) (Dai et al. 2018). But none of these methods provided any analysis or

explanation, they are completely empirical. We are the first to recognize the underestimated reconstruction loss leads to posterior collapse problem, and further we officially propose to reduce KL loss weight to overcome the posterior collapse with solid support, both theoretically and intuitively.

Experiments

Our proposed solution to the VAE model posterior collapse is simple but extremely effective and efficient. We do not need modify the network architecture and only adjust the training loss slightly, without introducing much extra computation. In this section, we will first train a vanilla VAE model and track the occur of model collapse, as well as experimentally verify that the reconstruction loss is underestimated. Then we will conduct extensive experiments to demonstrate the effectiveness of our proposed method.

VAE Architecture, Dataset, and Evaluation Metrics

We build our VAE model based on GRU. The VAE encoder is composed of two layers of bi-directional GRU which is better at capturing the sequence representation (Schuster and Paliwal 1997), and the hidden size of each layer is 512. The decoder is made up of four layers of uni-directional GRU with the same hidden size 512. Following previous work (Gómez-Bombarelli et al. 2018) (Jin, Barzilay, and Jaakkola 2018), we use unit Gaussian prior and set the latent vector dimension to be 56. The ELBO objective is optimized with Adam (Kingma and Ba 2014) and learning rate is 0.0001. The model is trained with teacher forcing and KL loss annealing following previous work. Since the model has a really good convergence, we train the model for 150 epochs and report the performance of the final model. We implement our model using PyTorch (Paszke et al. 2017). Experiments are conducted on a machine with a Intel Core i7-5930K@3.50GHz CPU and a GTX 1080 Ti GPU.

We conduct all our experiments on ZINC 250K dataset (Kusner, Paige, and Hernández-Lobato 2017) which is a subset of the ZINC (Sterling and Irwin 2015). Molecule sequences are tokenized with the regular expression from (Schwaller et al. 2018). We use the same training and testing split as previous work (Kusner, Paige, and Hernández-Lobato 2017) (Jin, Barzilay, and Jaakkola 2018), and have 10K hold-out data out of the training as the validation data. From now on, we will use the same experimental setting in all our experiments unless explicitly stated.

As for the model evaluation metrics, we report the reconstruction accuracy and validity score like previous work. Following (Jin, Barzilay, and Jaakkola 2018), we encode each molecule from test dataset 10 times, and then decode 10 times for each latent vector. The reconstruction accuracy is defined to be the ratio of successfully reconstructed molecule sequences to the total tried reconstruction. The reconstructed SMILES must be exactly the same as the input to be counted as successful. To calculate validity, 1000 latent vectors are randomly sampled from the prior distribution, and each is decoded 100 times. The validity is the portion of chemically valid reconstruction SMILES to the total decoded sequences. We use RDKit (Landrum and others 2006) to check if a SMILES is valid.

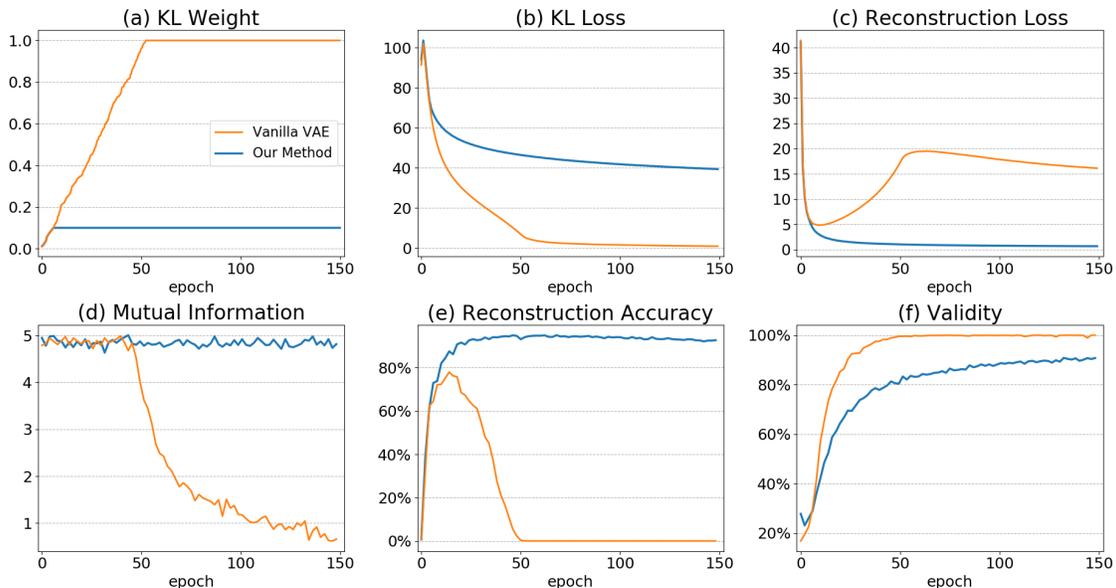


Figure 2: Training dynamic of vanilla VAE model on validation data. We track (a) KL weight β , (b) KL loss $D_{\text{KL}}(q_{\phi}(z|x)||p(z))$, (c) reconstruction loss $-\mathbb{E}_{q_{\phi}(z|x)}[\log p_{\theta}(x|z)]$, (d) mutual information I_q , (e) reconstruction accuracy, and (f) validity during the training. The orange line is the vanilla VAE with training KL loss annealing, and the maximum KL weight β is 1. Our method (Blue) reduces the maximum value of β to 0.1. Both models are trained with KL weight annealing and teacher forcing.

VAE Training Dynamic

We track the training process of a vanilla VAE model for SMILES sequences, as well as that of our proposed method. By investigating the training dynamic like KL weight, KL loss, reconstruction loss, mutual information, as well as the model performance (reconstruction accuracy and validity), we conclude that the underestimated reconstruction loss causes the posterior collapse of vanilla VAE model during the training process. Mutual information I_q can be calculated using Monte Carlo sampling as proposed in (Hoffman and Johnson 2016) (Dieng et al. 2019):

$$I_q = \mathbb{E}_{p_d(x)}[D_{\text{KL}}(q_{\phi}(z|x)||p(z))] - D_{\text{KL}}(q_{\phi}(z)||p(z)), \quad (7)$$

which is actually the same as the (2). We approximate the aggregated posterior $q_{\phi}(z) = \mathbb{E}_{p_d(x)}[q_{\phi}(z|x)]$ using Monte Carlo sampling. $D_{\text{KL}}(q_{\phi}(z)||p(z))$ can also be estimated by the Monte Carlo, and we can obtain samples from $q_{\phi}(z)$ by ancestral sampling. More details about I_q computation can be found in (Hoffman and Johnson 2016).

As a comparison, we also illustrate the training dynamic when our proposed method is applied. For our method, we set the KL weight $\beta = 0.1$ which is the optimal parameter we found. We keep all the other experimental settings the same as the vanilla VAE to make a fair comparison.

Results of two models run are plot in the Figure 2. The vanilla VAE model performances well on the validation data at the early stage of the KL weight annealing. As the KL weight increases, KL loss drops quickly as expected since more penalty is added to the KL loss term, while the small reconstruction loss starts to rise at the same time. The mutual

information I_q decreases to 0.65 at the end, which means the decoder does not absorb much information from the latent vectors when generating the output. This evidence indicates the posterior collapse has happened. When looking at the model performance on validation data, we can notice that the reconstruction accuracy is close 0% while the validity score is almost perfect. This indicates that too much pressure has been placed on the KL loss, which breaks the balance between the reconstruction loss and KL loss and results in the model posterior collapse.

Our method achieves lower reconstruction loss early and can maintain it during model training. Although the KL loss of our method is larger than the vanilla VAE, considering that we have a much smaller KL weight β now, the equivalent KL loss added to the training objective should still be in the normal range. Especially, our method maintains the mutual information to be around 4.8, which means output sequences are strongly related to latent vectors. As for the model performance, our method achieves 92.7% reconstruction accuracy and 90.7% validity score, which proves the superiority of our method.

Proof of Underestimated Reconstruction Loss

We hypothesize that introducing ground-truth information into the decoder will result in underestimated reconstruction loss, and have provided our detailed analysis previously. In this section, we will experimentally verify that the reconstruction loss is indeed underestimated during the training. We can estimate how much the reconstruction loss has been underestimated using Monte Carlo Sampling. Specifically, we can sample a batch of data, then run the model with

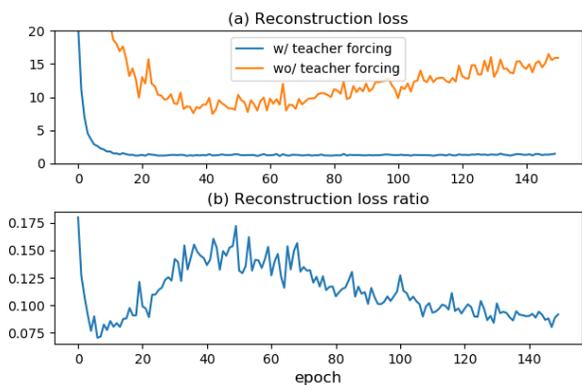


Figure 3: (a) Reconstruction loss on validation dataset. At each time step, models parameters are the same when calculating training and testing loss. (b) Reconstruction loss underestimated ratio.

and without the teacher forcing, respectively. The underestimated ratio can be approximated by the ratio of reconstruction loss with teacher forcing to that without teacher forcing.

We track the reconstruction loss on the validation dataset when the teacher forcing is applied and removed, respectively. Results are shown in the Figure 3(a). When teacher forcing is applied, the reconstruction loss drops close to 1 quickly, while the loss is much larger (at least 7.5) without teacher forcing. This is expected since without teacher forcing, any wrong prediction token as input may result in following prediction totally different from ground-truth sequences.

To figure out how much the reconstruction loss has been underestimated, we can compute the ratio as reconstruction loss w/ teacher forcing to that wo/ teacher forcing at each time step. Results are shown in Figure 3(b). It confirms our claim that the reconstruction loss is underestimated. To recover a balanced VAE loss, we can set KL loss weight exactly as the underestimated ratio in each epoch. To be simplified, we set $\beta = 0.1$ and we find it works well in practice.

Table 1: Reconstruction accuracy and validity results. Baseline results are reported in (Kusner, Paige, and Hernández-Lobato 2017) (Dai et al. 2018) (Simonovsky and Komodakis 2018) (Jin, Barzilay, and Jaakkola 2018).

Model	Reconstruction	Validity
SMILES-based		
CVAE	44.6%	0.7%
GVAE	53.7%	7.2%
SD-VAE	76.2%	43.5%
Our Method	92.7%	90.7%
Graph-based		
GraphVAE	-	13.5%
JT-VAE	76.7%	100.0%

Molecule Reconstruction Accuracy and Validity

We summarize the molecule reconstruction accuracy and validity on test dataset in the Table 1. Our method outperforms all previous models in reconstruction accuracy by a large margin (16% larger than the second best model). In the meanwhile, our method achieves 90.7 % validity, which is the second best among all the models.

Compared with other SMILES-based methods, our model is much more superior in both the reconstruction accuracy and prior validity, even if complex grammar or syntax rules are incorporated (Kusner, Paige, and Hernández-Lobato 2017) (Dai et al. 2018). Note that JT-VAE model assembles molecules by adding sub-graphs step-by-step to make sure the generated molecule graphs are always valid. However, the sub-graphs are extracted from the training dataset, which limits the JT-VAE can not generate molecules with unseen sub-graphs. Our method achieves competitive validity performance without any constraints, and is able to generate novel molecules that are not from the same distribution as the training data. That is one important reason why our method achieves the STOA reconstruction accuracy, while JT-VAE suffers from reconstructing testing molecules (Mohammadi et al. 2019). Besides, our method is much more efficient than JT-VAE. When generating 10,000 unique valid SMILES from prior random sampling, JT-VAE⁵ (faster version) takes about 1450s while our method only needs 9s.

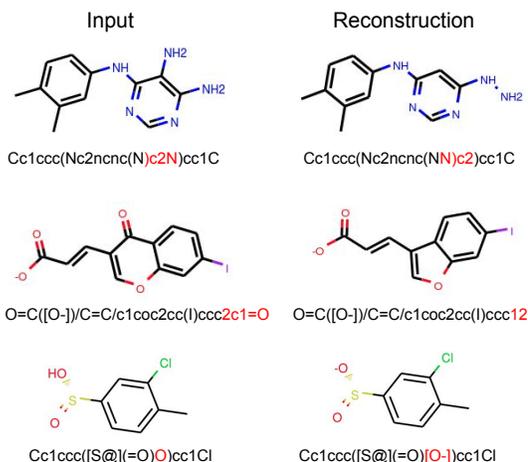


Figure 4: Reconstruction error examples. Unmatched tokens between the input and reconstruction SMILES are shown in red (“[O-]” is one token).

Error Analysis and Visualization

Our model achieves 92.7% reconstruction accuracy. We investigate the reconstruction results further and find that our model can predict 97.3% of all tokens correctly, which is measured on the level of token instead of the sequence. Besides, most of unmatched sequences (62%) are valid. These evidences indicate that our model is very well learned. We show some valid but unmatched examples in Figure 4.

⁵<https://github.com/wengong-jin/icml18-jtnn>

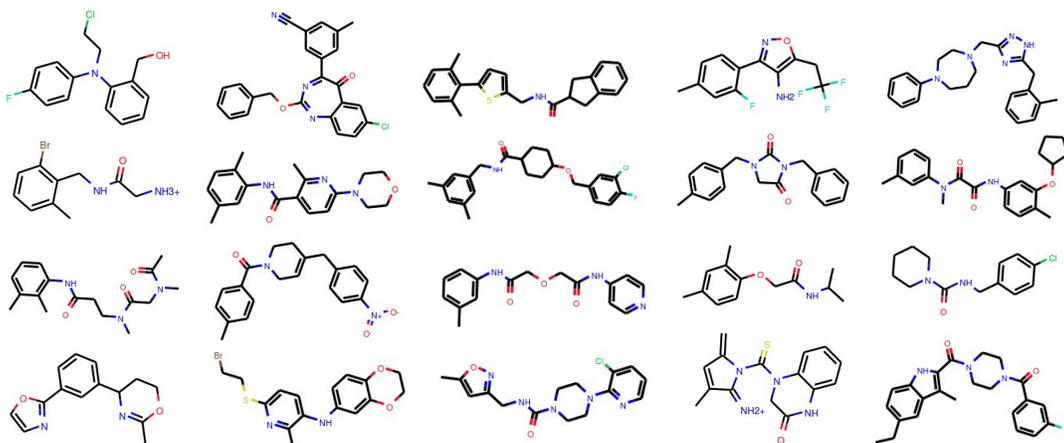


Figure 5: Generated molecules by random sampling from the prior.

As for the validity, we also investigate model outputs. We illustrate some generated molecules in the Figure 5, which demonstrates that our model can generate complicated and diverse molecules with multiple rings.

For those invalid sequences, from both the reconstruction and prior sampling, there are several typical errors: (1) unkekulized atoms, (2) valence error, (3) unclosed ring, and (4) parentheses error. We believe that advanced techniques like grammar and syntax rules (Kusner, Paige, and Hernández-Lobato 2017) (Dai et al. 2018) are necessary and can help to reduce these kind of errors, and our method is essential and complementary to these methods.

Bayesian Optimization

One of the important tasks in the drug molecule generation is to make molecules with desired chemical properties. We follow (Kusner, Paige, and Hernández-Lobato 2017) (Jin, Barzilay, and Jaakkola 2018) for all the experimental setting, and the optimization target score is:

$$y(m) = \log P(m) - SA(m) - cycle(m), \quad (8)$$

where $\log P(m)$ is the octanol-water partition coefficients of molecule m , $SA(m)$ is synthetic accessibility score, and $cycle(m)$ is number of large rings with more than six atoms.

Table 2: Top-3 molecule property scores found by the BO. Baseline results are copied from (Kusner, Paige, and Hernández-Lobato 2017) (Dai et al. 2018) (Jin, Barzilay, and Jaakkola 2018) (Jin, Barzilay, and Jaakkola 2018).

Model	1st	2nd	3rd
SMILES-based			
CVAE	1.98	1.42	1.19
GVAE	2.94	2.89	2.80
SD-VAE	4.04	3.50	2.96
Our Method	5.32	5.28	5.23
Graph-based			
JT-VAE	5.30	4.93	4.49

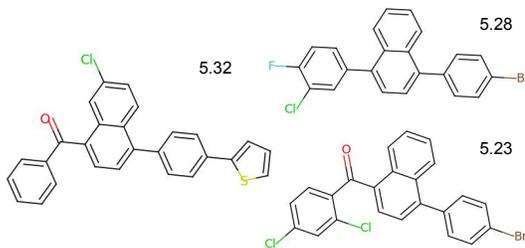


Figure 6: Top-3 molecules and associated scores found by our model with Bayesian optimization.

We first associate each molecule with a latent vector which is the mean of the learned variational encoding distribution. The latent vector for each molecule will be treated as its feature and we train a Sparse Gaussian Process (SGP) to predict target score $y(m)$ given its latent vector. After training SGP, five iterations of batched Bayesian optimization (BO) are performed with expected improvement heuristics.

We report SGP prediction performance when trained on latent representations learned by different models. We train the SGP with 10-fold cross validation considering randomness and report the top-3 molecules found by the BO.

As shown in Table 2, molecules found by our model are much better than that by previous SMILES-based methods, and our method is even superior to the graph-based method JT-VAE. Figure 6 shows top-3 molecules found by our model.

Discussion

Our method works extremely well in the molecule generation, in which SMILES sequences are highly structured and grammarly organized. Our experimental results confirm that grammar and syntax rules are necessary to generate more valid SMILES sequences. Besides, SMILES-based methods and graph-based methods may be combined together to boost the model performance further.

References

- Bowman, S. R.; Vilnis, L.; Vinyals, O.; Dai, A.; Jozefowicz, R.; and Bengio, S. 2016. Generating sentences from a continuous space. In *Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning*, 10–21.
- Chung, J.; Gulcehre, C.; Cho, K.; and Bengio, Y. 2014. Empirical evaluation of gated recurrent neural networks on sequence modeling. In *NIPS 2014 Workshop on Deep Learning, December 2014*.
- Dai, H.; Tian, Y.; Dai, B.; Skiena, S.; and Song, L. 2018. Syntax-directed variational autoencoder for structured data. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Dieng, A. B.; Kim, Y.; Rush, A. M.; and Blei, D. M. 2019. Avoiding latent variable collapse with generative skip models. In *The 22nd International Conference on Artificial Intelligence and Statistics*, 2397–2405.
- Gómez-Bombarelli, R.; Wei, J. N.; Duvenaud, D.; Hernández-Lobato, J. M.; Sánchez-Lengeling, B.; Sheberla, D.; Aguilera-Iparraguirre, J.; Hirzel, T. D.; Adams, R. P.; and Aspuru-Guzik, A. 2018. Automatic chemical design using a data-driven continuous representation of molecules. *ACS central science* 4(2):268–276.
- Hao Fu, Chunyuan Li, X. L. J. G. A. C. L. C. 2019. Cyclical annealing schedule: A simple approach to mitigating KL vanishing. In *NAACL*.
- He, J.; Spokoyny, D.; Neubig, G.; and Berg-Kirkpatrick, T. 2019. Lagging inference networks and posterior collapse in variational autoencoders. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Higgins, I.; Matthey, L.; Pal, A.; Burgess, C.; Glorot, X.; Botvinick, M.; Mohamed, S.; and Lerchner, A. 2017. beta-vae: Learning basic visual concepts with a constrained variational framework. In *ICLR*.
- Hochreiter, S., and Schmidhuber, J. 1997. Long short-term memory. *Neural computation* 9(8):1735–1780.
- Hoffman, M. D., and Johnson, M. J. 2016. Elbo surgery: yet another way to carve up the variational evidence lower bound. In *NIPS Workshop on Advances in Approximate Bayesian Inference*.
- Jin, W.; Barzilay, R.; and Jaakkola, T. 2018. Junction tree variational autoencoder for molecular graph generation. In *International Conference on Machine Learning*, 2328–2337.
- Kim, Y.; Wiseman, S.; Miller, A.; Sontag, D.; and Rush, A. 2018. Semi-amortized variational autoencoders. In *International Conference on Machine Learning*, 2683–2692.
- Kingma, D. P., and Ba, J. 2014. Adam: A method for stochastic optimization. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Kingma, D. P., and Welling, M. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.
- Kusner, M. J.; Paige, B.; and Hernández-Lobato, J. M. 2017. Grammar variational autoencoder. In *Proceedings of the 34th International Conference on Machine Learning—Volume 70*, 1945–1954. JMLR. org.
- Landrum, G., et al. 2006. Rdkit: Open-source cheminformatics.
- Li, Y.; Vinyals, O.; Dyer, C.; Pascanu, R.; and Battaglia, P. 2018. Learning deep generative models of graphs. *arXiv preprint arXiv:1803.03324*.
- Mohammadi, S.; O’Dowd, B.; Paulitz-Erdmann, C.; and Gorerlitz, L. 2019. Penalized variational autoencoder for molecular design. *ChemRxiv*.
- Paszke, A.; Gross, S.; Chintala, S.; Chanan, G.; Yang, E.; DeVito, Z.; Lin, Z.; Desmaison, A.; Antiga, L.; and Lerer, A. 2017. Automatic differentiation in pytorch.
- Polishchuk, P. G.; Madzhidov, T. I.; and Varnek, A. 2013. Estimation of the size of drug-like chemical space based on gdb-17 data. *Journal of computer-aided molecular design* 27(8):675–679.
- Rezende, D. J.; Mohamed, S.; and Wierstra, D. 2014. Stochastic backpropagation and approximate inference in deep generative models. In *International Conference on Machine Learning*, 1278–1286.
- Schuster, M., and Paliwal, K. K. 1997. Bidirectional recurrent neural networks. *IEEE Transactions on Signal Processing* 45(11):2673–2681.
- Schwaller, P.; Gaudin, T.; Lanyi, D.; Bekas, C.; and Laino, T. 2018. found in translation: predicting outcomes of complex organic chemistry reactions using neural sequence-to-sequence models. *Chemical science* 9(28):6091–6098.
- Simonovsky, M., and Komodakis, N. 2018. Graphvae: Towards generation of small graphs using variational autoencoders. In *International Conference on Artificial Neural Networks*, 412–422. Springer.
- Sterling, T., and Irwin, J. J. 2015. Zinc 15–ligand discovery for everyone. *Journal of chemical information and modeling* 55(11):2324–2337.
- Williams, R. J., and Zipser, D. 1989. A learning algorithm for continually running fully recurrent neural networks. *Neural computation* 1(2):270–280.
- Yang, Z.; Hu, Z.; Salakhutdinov, R.; and Berg-Kirkpatrick, T. 2017. Improved variational autoencoders for text modeling using dilated convolutions. In *Proceedings of the 34th International Conference on Machine Learning—Volume 70*, 3881–3890. JMLR. org.