

---

# Syntax-Directed Variational Autoencoder for Molecule Generation

---

Hanjun Dai<sup>1\*</sup>, Yingtao Tian<sup>2\*</sup>, Bo Dai<sup>1</sup>, Steven Skiena<sup>2</sup>, Le Song<sup>1</sup>

<sup>1</sup> College of Computing, Georgia Institute of Technology

<sup>2</sup> Department of Computer Science, Stony Brook University

<sup>1</sup> {hanjundai, bodai}@gatech.edu, lsong@cc.gatech.edu

<sup>2</sup> yittian, skiena@cs.stonybrook.edu

## Abstract

Deep generative models have been enjoying success in modeling continuous data. However it remains challenging to capture the representations for discrete structures with formal grammars and semantics. How to generate both syntactically and semantically correct data still remains largely an open problem. Inspired by the theory of compiler where syntax and semantics check is done via syntax-directed translation (SDT), we propose a novel syntax-directed variational autoencoder (SD-VAE) by introducing *stochastic lazy attributes*. This approach converts the offline SDT check into on-the-fly generated guidance for constraining the decoder. Comparing to the state-of-the-art methods, our approach enforces constraints on the output space so that the output will be not only *syntactically* valid, but also *semantically* reasonable. We evaluate the proposed model with applications including reconstruction and molecule optimization. The results demonstrate the effectiveness in incorporating syntactic and semantic constraints in discrete generative models, which is significantly better than current state-of-the-art approaches.

## 1 Introduction

Recent advances in deep representation learning have resulted in powerful probabilistic generative models which have demonstrated their ability on modeling continuous data, *e.g.*, time series signals [Oord et al., 2016, Dai et al., 2017] and images [Radford et al., 2015]. Despite the success in these domains, it is still challenging to correctly generate discrete structured data, such as graphs and molecules. Since many of the structures have syntax and semantic formalisms, the generative models without explicit constraints often produces invalid ones.

Conceptually an approach in generative model for structured data can be divided in two parts, one being the formalization of the structure generation and the other one being a (usually deep) generative model producing parameters for stochastic process in that formalization. Often the hope is that with the help of training samples and capacity of deep models, the loss function will prefer the valid patterns and encourage the mass of the distribution of the generative model towards the desired region automatically.

Arguably the simplest structured data are sequences, whose generation with deep model has been well studied under the seq2seq [Sutskever et al., 2014] framework that models the generation of sequence as a series of token choices parameterized by recurrent neural networks (RNNs). Its widespread success has encourage several pioneer works that consider the serialization of more complex structure data into sequences and apply sequence models to the represented sequences. Gómez-Bombarelli et al. [2016] (CVAE) is a representative work of such paradigm for the chemical molecule generation, using the SMILES line notation [Weininger, 1988] for representing molecules. However, because of the lack of formalization of syntax and semantics serving as the restriction of the *particular* structured data, underfitted *general-purpose* string generative models will often lead

---

\*Both authors contributed equally to the paper.

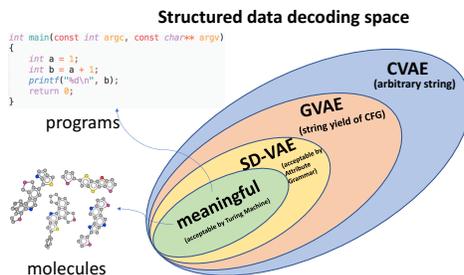


Figure 1: Illustrations of structure data decoding space with respect to different algorithms. The seq2seq model is the most flexible one with least constraints. Model with CFG constraints will output syntax correct data. Our proposed model with attribute grammar reshapes the output space tighter to the target meaningful space.

to invalid outputs. Therefore, to obtain a reasonable model via such training procedure, we need to prepare large amount of valid combinations of the structures, which in general is not practical.

To tackle such a challenge, one approach is to incorporate the structure restrictions explicitly into the generative model. For the considerations of computational cost and model generality, context-free grammars (CFG) have been taken into account in the decoder parametrization. For instance, in molecule generation tasks, Kusner et al. [2017] proposes a grammar variational autoencoder (GVAE) in which the CFG of SMILES notation is embedded into the decoder. The model generates the parse trees directly in a top-down direction, by repeatedly expanding any nonterminal with its production rules. Although the CFG provides a mechanism for generating *syntactic valid* objects, it is still incapable to regularize the model for generating *semantic valid* objects [Kusner et al., 2017]. For example, in molecule generation, the semantic of the SMILES languages requires that the rings generated must be closed. This example requires cross-serial dependencies which are not enforceable by CFG, implying that more constraints beyond CFG are needed to achieve semantic valid production.

In the theory of compiler, attribute grammars, or syntax-directed definition has been proposed for attaching semantics to a *tree yield* of context-free grammar. That is, semantics are attached to an already materialized sequence and its CFG generating tree, thus one straightforward but not practical application of attribute grammars is, after generating a syntactic valid molecule candidate, to conduct offline semantic checking. This process needs to be repeated until a semantically valid one is discovered, which is computationally inefficient and in the worst case infeasible, due to extremely low rate of passing checking. As a remedy, we propose the *syntax-direct variational autoencoder* (SD-VAE), in which a semantic restriction component is advanced to the stage of syntax tree generator. This allows the generator with both syntactic and semantic validation. The proposed syntax-direct generative mechanism in the decoder further constrains the output space to ensure the semantic correctness in the tree generation process. The relationships between our proposed model and previous models can be characterized in Figure 1.

Our method brings theory of formal language into stochastic generative model. The contribution of our paper can be summarized as follows:

- *Syntax and semantics enforcement*: We propose a new formalization of semantics that systematically converts the offline semantic check into online guidance for stochastic generation using the proposed *stochastic lazy attribute*. This allows us effectively address both syntax and semantic constraints.
- *Efficient learning and inference*: Our approach has computational cost  $O(n)$  where  $n$  is the length of structured data. This is the same as existing methods like CVAE and GVAE which do not enforce semantics in generation.
- *Strong empirical performance*: We demonstrate the effectiveness of the SD-VAE through applications in molecules. Our approach consistently and significantly improves the results in evaluations including generation, reconstruction and optimization.

## 2 Background

Before introducing our model and the learning algorithm, we first provide some background knowledge which is important for understanding the proposed method.

## 2.1 Variational Autoencoder

The variational autoencoder [Kingma and Welling, 2013, Rezende et al., 2014] provides a framework for learning the probabilistic generative model as well as its posterior, respectively known as decoder and encoder. We denote the observation as  $x$ , which is the structured data in our case, and the latent variable as  $z$ . The decoder is modeling the probabilistic generative processes of  $x$  given the continuous representation  $z$  through the likelihood  $p_\theta(x|z)$  and the prior over the latent variables  $p(z)$ , where  $\theta$  denotes the parameters. The encoder approximates the posterior  $p_\theta(z|x) \propto p_\theta(x|z)p(z)$  with a model  $q_\psi(z|x)$  parametrized by  $\psi$ . The decoder and encoder are learned simultaneously by maximizing the evidence lower bound (ELBO) of the marginal likelihood, *i.e.*,

$$\mathcal{L}(X; \theta, \psi) := \sum_{x \in X} \mathbb{E}_{q(z|x)} [\log p_\theta(x|z)p(z) - \log q_\psi(z|x)] \leq \sum_{x \in X} \log \int p_\theta(x|z)p(z) dz, \quad (1)$$

where  $X$  denotes the training datasets containing the observations.

## 2.2 Context Free Grammar and Attribute Grammar

**Context free grammar** A context free grammar (CFG) is defined as  $G = \langle \mathcal{V}, \Sigma, \mathcal{R}, s \rangle$ , where symbols are divided into  $\mathcal{V}$ , the set of non-terminal symbols,  $\Sigma$ , the set of terminal symbols and  $s \in \mathcal{V}$ , the start symbol. Here  $\mathcal{R}$  is the set of production rules. Each production rule  $r \in \mathcal{R}$  is denoted as  $r = \alpha \rightarrow \beta$ , where  $\alpha \in \mathcal{V}$  is a nonterminal, and  $\beta = u_1 u_2 \dots u_{|\beta|} \in (\mathcal{V} \cup \Sigma)^*$  is a sequence of terminals and/or nonterminals.

**Attribute grammar** To enrich the CFG with “semantic meaning”, Knuth [1968] formalizes attribute grammar that introduces attributes and rules to CFG. The attribute is an attachment to the corresponding nonterminal symbol in CFG, written in the format  $\langle v \rangle . a$  where  $\langle v \rangle \in \mathcal{V}$ . There can be two types of attributes assigned to non-terminals in  $G$ : the *inherited* attributes and the *synthesized* attributes. An inherited attribute depends on the attributes from its parent and siblings, while a synthesized attribute is computed based on the attributes of its children. Formally, for a production  $u_0 \rightarrow u_1 u_2 \dots u_{|\beta|}$ , we denote  $I(u_i)$  and  $S(u_i)$  be the (disjoint) sets of *inherited* and *synthesized* attributes of  $u_i$ ,  $i \in \{0, \dots, |\beta|\}$ .

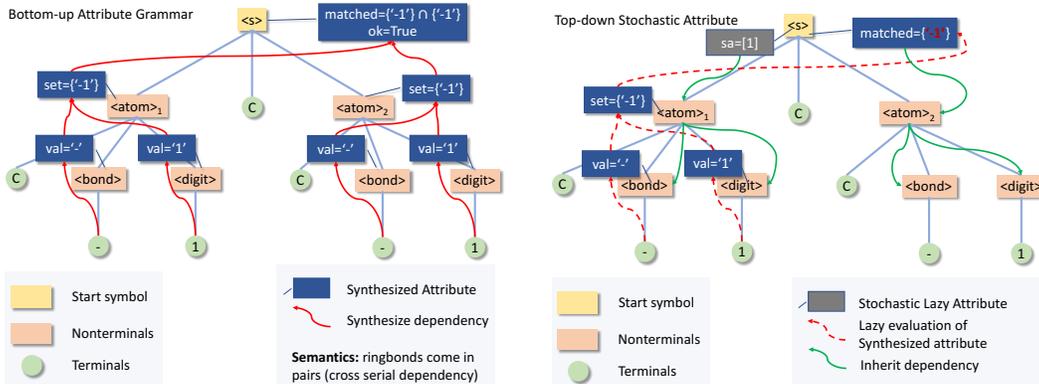
### 2.2.1 A motivational example

We here exemplify how the above defined attribute grammar enriches CFG with non-context-free semantics. We use the following toy grammar, a subset of SMILES that generates either a chain or a cycle with three carbons.

Production	Semantic Rule
$\langle s \rangle \rightarrow \langle atom \rangle_1 \text{'C'} \langle atom \rangle_2$	$\langle s \rangle . \text{matched} \leftarrow \langle atom \rangle_1 . \text{set} \cap \langle atom \rangle_2 . \text{set}$ $\langle s \rangle . \text{ok} \leftarrow \langle atom \rangle_1 . \text{set} = \langle s \rangle . \text{matched} = \langle atom \rangle_2 . \text{set}$
$\langle atom \rangle \rightarrow \text{'C'} \mid \text{'C'} \langle bond \rangle \langle digit \rangle$	$\langle atom \rangle . \text{set} \leftarrow \emptyset \mid \text{concat}(\langle bond \rangle . \text{val}, \langle digit \rangle . \text{val})$
$\langle bond \rangle \rightarrow \text{'-'} \mid \text{'='} \mid \text{'\#}'$	$\langle bond \rangle . \text{val} \leftarrow \text{'-'} \mid \text{'='} \mid \text{'\#'}$
$\langle digit \rangle \rightarrow \text{'1'} \mid \text{'2'} \mid \dots \mid \text{'9'}$	$\langle digit \rangle . \text{val} \leftarrow \text{'1'} \mid \text{'2'} \mid \dots \mid \text{'9'}$

where we show the production rules in CFG with  $\rightarrow$  on the left, and the calculation of attributes in attribute grammar with  $\leftarrow$  on the left. Here we leverage the attribute grammar to check (with attribute *matched*) whether the ringbonds come in pairs: a ringbond generated at  $\langle atom \rangle_1$  should match the bond type and bond index that generated at  $\langle atom \rangle_2$ , also the semantic constraint expressed by  $\langle s \rangle . \text{ok}$  requires that there is no difference between the *set* attribute of  $\langle atom \rangle_1$  and  $\langle atom \rangle_2$ . Actually such constraint in SMILES is known as *cross-serial dependencies* (CSD) [Bresnan et al., 1982] which is non-context-free [Shieber, 1985]. Another example of CSD is a sequence of multiple different types of parentheses where each separately balanced disregarding the others. Figure 2a further illustrates the example. Here all the attributes are *synthetic*, *i.e.*, calculated in a bottom-up direction.

In the semantic correctness checking procedure, one need to perform (possibly multiple) bottom-up and top-down procedures for calculating the attributes *after* the parse tree is generated, however, in the structure generating process, the parse tree is not ready for semantic checking, since the synthesized attributes coming from children are not generated yet. Due to such dilemma, it is nontrivial to use the attribute grammar to guide the top-down generation of the tree-structured data.



(a) Bottom-up semantics check.

(b) Top-down tree generation.

Figure 2: Illustrations of (a) the bottom-up semantic checking following attribute grammar, and (b) the top-down tree generation with semantic validation with stochastic attribute grammar.

---

### Algorithm 1 Decoding with Stochastic Syntax-Directed Decoder

---

- 1: **Global variables:** CFG:  $G = (\mathcal{V}, \Sigma, \mathcal{R}, s)$ , decoder network parameters  $\theta$
  - 2: **procedure** GENTREE( $node, \mathcal{T}$ )
  - 3:   Sample stochastic lazy attribute  $node.sa \sim \mathcal{B}_\theta(sa|node, \mathcal{T})$    ▷ when introduced on  $node$
  - 4:   Sample production rule  $r = (\alpha \rightarrow \beta) \in \mathcal{R} \sim p_\theta(r|ctx, node, \mathcal{T})$ .   ▷ The conditioned variables encodes the semantic constraints in tree generation.
  - 5:    $ctx \leftarrow \text{RNN}(ctx, r)$    ▷ update context vector
  - 6:   **for**  $i = 1, \dots, |\beta|$  **do**
  - 7:      $v_i \leftarrow \text{Node}(u_i, node, \{v_j\}_{j=1}^{i-1})$    ▷ node creation with parent and siblings' attributes
  - 8:     GenTree( $v_i, \mathcal{T}$ )   ▷ recursive generation of children nodes
  - 9:     Update synthetic and stochastic attributes of  $node$  with  $v_i$    ▷ Lazy linking
  - 10:   **end for**
  - 11: **end procedure**
- 

One straightforward way is using acceptance-rejection sampling scheme, *i.e.*, using the CFG decoder in grammar VAE Kusner et al. [2017] as a proposal and the semantic checking as the threshold. It is obvious that since the decoder does not include semantic guidance, the proposal distribution may raise semantically invalid candidate frequently, therefore, wasting the computational cost in vain.

## 3 Syntax-Directed Variational Autoencoder

As described in Section 2.2.1, directly using attribute grammar to address both syntax and semantics constraints is not efficient. In this section we describe how to bring forward the attribute grammar online and incorporate it into variational autoencoders such that our VAE generates both *syntactic* and *semantic* valid outputs by definition. We name our proposed method Syntax-Directed Variational Autoencoder (SD-VAE).

### 3.1 Stochastic Syntax-Directed Decoder

By scrutinizing the tree generation, the major difficulty in incorporating the attributes grammar into the processes is the appearance of the synthesized attributes. For instance, when expanding the start symbol  $\langle s \rangle$ , the corresponding synthesized attribute  $\langle s \rangle$ .`matched` is not ready yet. Since none of its children is generated, their synthesized attributes are also absent at this time, making the  $\langle s \rangle$ .`matched` unable to be computed. To enable the on-the-fly computation of the synthesized attributes for semantic validation during tree generation, besides the two types of attributes, we introduce the *stochastic lazy attributes* to enlarge the existing attribute grammar, so that the synthesized attributes will be transformed to inherited constraints in generating procedure and instantiated once all the dependent attributes are ready (also named as *lazy linking* in the following content).

We demonstrate how the decoder with *stochastic lazy attributes* will generate semantic valid output through a pedagogical example with the subset of SMILES grammar in figure 2(b). Following the terminology in compiler theory, we named it as stochastic syntax-directed decoder.

The tree generation procedure is indeed sampling from the decoder  $p_\theta(x|z)$ , which can be decomposed into several steps that elaborated below:

**i) stochastic predetermination:** in figure 2(b), we start from the node  $\langle s \rangle$  with the synthesized attributes  $\langle s \rangle.\text{matched}$  determining the index and bond type of the ringbond that will be matched at node  $\langle s \rangle$ . Since we know nothing about the children nodes right now, the only thing we can do is to ‘guess’ a value. That is to say, we associate a stochastic attribute  $\langle s \rangle.\text{sa} \in \{0, 1\}^{C_a} \sim \prod_{i=1}^{C_a} \mathcal{B}(sa_i|z; p)$  as a predetermination for the sake of the absence of synthesized attribute  $\langle s \rangle.\text{matched}$ .  $\mathcal{B}(\cdot)$  is the bernoulli distribution. Here  $C_a$  is the maximum cardinality possible<sup>1</sup> for the corresponding attribute  $a$ . In above example, the 0 indicates no ringbond and 1 indicates one ringbond at both  $\langle atom \rangle_1$  and  $\langle atom \rangle_2$ , respectively.

**ii) constraints as inherited attributes:** we pass the  $\langle s \rangle.\text{sa}$  as inherited constraints to the children of node  $\langle s \rangle$ , *i.e.*,  $\langle atom \rangle_1$  and  $\langle atom \rangle_2$  to ensure the semantic validation in the tree generation.

**iii) sampling under constraints:** assume in the valid order,  $\langle atom \rangle_1$  is selected before  $\langle atom \rangle_2$ , we then sample the rules from  $p_\theta(r|\langle atom \rangle_1, \langle s \rangle, z)$  for expanding  $\langle atom \rangle_1$ , and so on and so forth to generate the subtree recursively. Since we carefully designed sampling distribution that is conditioning on the stochastic property, the inherited constraints will be eventually satisfied. In the example, due to the  $\langle s \rangle.\text{sa} = \text{‘1’}$ , when expanding  $\langle atom \rangle_1$ , the sampling distribution  $p_\theta(r|\langle atom \rangle_1, \langle s \rangle, z)$  only has positive mass on rule  $\langle atom \rangle \rightarrow \text{‘C’} \langle bond \rangle \langle digit \rangle$ .

**iv) lazy linking:** once we complete the generation of the subtree rooted at  $\langle atom \rangle_1$ , the synthesized attribute  $\langle atom \rangle_1.\text{set}$  is now available. According to the semantic rule for  $\langle s \rangle.\text{matched}$ , we can instantiate  $\langle s \rangle.\text{matched} = \langle atom \rangle_1.\text{set} = \{\text{‘-1’}\}$ . When expanding  $\langle atom \rangle_2$ , the  $\langle s \rangle.\text{matched}$  will be passed down as inherited attribute to regulate the generation of  $\langle atom \rangle_2$ .

In summary, the general syntax tree  $\mathcal{T} \in L(G)$  can be constructed step by step, within the languages  $L(G)$  covered by grammar  $G$ . In the beginning,  $\mathcal{T}^{(0)} = \text{root}$ , where  $\text{root}.\text{symbol} = s$  which contains only the start symbol  $s$ . At step  $t$ , we will choose an nonterminal node in the *frontier*<sup>2</sup> of partially generated tree  $\mathcal{T}^{(t)}$  to expand. The generative process in each step  $t = 0, 1, \dots$  can be described as:

1. Pick node  $v^{(t)} \in Fr(\mathcal{T}^{(t)})$  where its attributes needed are either satisfied, or are stochastic attributes that should be sampled first according to bernoulli distribution  $\mathcal{B}(\cdot|v^{(t)}, \mathcal{T}^{(t)})$ ;
2. Sample rule  $r^{(t)} = \alpha^{(t)} \rightarrow \beta^{(t)} \in \mathcal{R}$  according to distribution  $p_\theta(r^{(t)}|v^{(t)}, \mathcal{T}^{(t)})$ , where  $v^{(t)}.\text{symbol} = \alpha^{(t)}$ , and  $\beta^{(t)} = u_1^{(t)}u_2^{(t)} \dots u_{|\beta^{(t)}|}^{(t)}$ , *i.e.*, expand the nonterminal with production rules defined in CFG.
3.  $\mathcal{T}^{(t+1)} = \mathcal{T}^{(t)} \cup \{(v^{(t)}, u_i^{(t)})\}_{i=1}^{|\beta^{(t)}|}$ , *i.e.*, grow the tree by attaching  $\beta^{(t)}$  to  $v^{(t)}$ . Now the node  $v^{(t)}$  will have children represented by symbols in  $\beta^{(t)}$ .

The above process continues until all the nodes in the frontier of  $\mathcal{T}^{(T)}$  are all terminals after  $T$  steps. Then, we obtain the algorithm 1 for sampling both syntactic and semantic valid structures.

In fact, in the model training phase, we need to compute the likelihood  $p_\theta(x|z)$  given  $x$  and  $z$ . The probability computation procedure is similar to the sampling procedure in the sense that both of them requires tree generation. The only difference is that in the likelihood computation procedure, the tree structure, *i.e.*, the computing path, is fixed since  $x$  is given, while in sampling procedure, it is sampled following the learned model. Specifically, the generative likelihood can be written as:

$$p_\theta(x|z) = \prod_{t=0}^T p_\theta(r_t|ctx^{(t)}, \text{node}^{(t)}, \mathcal{T}^{(t)}) \mathcal{B}_\theta(sa_t|\text{node}^{(t)}, \mathcal{T}^{(t)}) \quad (2)$$

where  $ctx^{(0)} = z$  and  $ctx^{(t)} = \text{RNN}(r_t, ctx^{(t-1)})$ . Here RNN can be commonly used LSTM, *etc.*.

### 3.2 Structure-Based Encoder

As we introduced in section 2, the encoder,  $q_\psi(z|x)$  approximates the posterior of the latent variable through the model with some parametrized function with parameters  $\psi$ . Since the structure in the observation  $x$  plays an important role, the encoder parametrization should take care of such information. The recently developed deep learning models [Duvenaud et al., 2015, Dai et al., 2016,

<sup>1</sup>Note that setting threshold for  $C_a$  assumes a *mildly context sensitive grammar* (e.g., limited CSD).

<sup>2</sup>Here frontier is the set of all nonterminal leaves in current tree.

Lei et al., 2017] provide powerful candidates as encoder. However, to demonstrate the benefits of the proposed syntax-directed decoder in incorporating the attribute grammar for semantic restrictions, we will exploit the same encoder in Kusner et al. [2017] for a fair comparison later.

We provide a brief introduction to the particular encoder model used in Kusner et al. [2017] for a self-contained purpose. Given a SMILES sequence, we obtain the corresponding parse tree using CFG and decompose it into a sequence of productions through a pre-order traversal on the tree. Then, we convert these productions into one-hot indicator vectors, in which each dimension corresponds to one production in the grammar. We will use a deep convolution neural networks which maps this sequence of one-hot vectors to a continuous vector as the encoder.

### 3.3 Model Learning

Our learning goal is to maximize the evidence lower bound in Eq 1. During training, each instance  $x$  is first parsed into syntax tree with CFG parser. Given the encoder, we can then map the structure input into latent space  $z$ . The variational posterior  $q(z|x)$  is parameterized with Gaussian distribution, where the mean and variance are the output of corresponding neural networks. The prior of latent variable  $p(z) = \mathcal{N}(0, I)$ . Since both the prior and posterior are Gaussian, we use the closed form of KL-divergence that proposed in Kingma and Welling [2013]. In the decoding stage, our goal is to maximize  $p_\theta(x|z)$ . Using the Algorithm 1, we can compute the corresponding conditional likelihood.

For efficient learning, during the training time we divide the calculation in our stochastic decoder into two phases: the first phase generates tree and the second phase only consists of a sequence of updates to the context vector. This decoupling is possible since in training time we know the all decisions in the decoder since the tree’s calculation is deterministic. In doing so, the first phase can be accelerated using multiple CPU cores in parallel and the second one can effectively computed using any mini-batch optimization on GPU. In practice, we observe no significant time penalty measured in wall clock time compared to previous works.

## 4 Related work

Generative models with discrete structured data have raised increasing interests among researchers in different domains. The classical sequence to sequence model [Sutskever et al., 2014] and its variations have also been applied to molecules [Gómez-Bombarelli et al., 2016]. Since the model is quite flexible, it is hard to generate valid structures with limited data. Techniques including data augmentation [Bjerrum, 2017], active learning [Janz et al., 2017] and reinforcement learning [Guimaraes et al., 2017] have also been proposed to tackle this issue. However, according to the empirical evaluations from Benhenda [2017], the validity is still not satisfactory. Even when the validity is enforced, the models tend to overfit to simple structures while neglect the diversity.

Since the structured data often comes with formal grammars, it is very helpful to generate its parse tree derived from CFG, instead of generating sequence of tokens directly. The Grammar VAE[Kusner et al., 2017] introduced the CFG constrained decoder for simple math expression and SMILES string generation. The rules are used to mask out invalid syntax such that the generated sequence is always from the language defined by its CFG. Parisotto et al. [2016] uses a RecursiveReverse-Recursive Neural Network (R3NN) to capture global context information while expanding with CFG production rules. Although these works follows the syntax via CFG, the context sensitive information can only be captured using variants of sequence/tree RNNs [Alvarez-Melis and Jaakkola, 2016, Dong and Lapata, 2016, Zhang et al., 2015], which may not be time and sample efficient.

In our work, we capture the semantics with proposed stochastic lazy attributes when generating structured outputs. By addressing the most common semantics to harness the deep networks, it can greatly reshape the output domain of decoder [Hu et al., 2016]. As a result, we can also get a better generative model for discrete structures.

## 5 Experiments

We compare our method with GVAE Kusner et al. [2017] and CVAE Gómez-Bombarelli et al. [2016] (a character variational autoencoder) using the protocols for experiments that are set up in Kusner et al. [2017]. The dataset contains 250,000 SMILES string, prepared by Kusner et al. [2017] using randomly extraction from the ZINC database [Gómez-Bombarelli et al., 2016]. We use 5000 SMILES strings as the holdout set for testing and the rest for training. For syntax, our formalization of SMILES follows the grammar specified in Appendix A. For our SD-VAE, we address some of the

Method	Reconstruction %	Valid Prior %	Method	LL	RMSE
SD-VAE	<b>76.2</b>	<b>43.5</b>	CVAE	-1.812 ± 0.004	1.504 ± 0.006
GVAE	53.7	7.2	GVAE	-1.739 ± 0.004	1.404 ± 0.006
CVAE	44.6	0.7	SD-VAE	<b>-1.697 ± 0.015</b>	<b>1.366 ± 0.023</b>

(a) Reconstruction Accuracy and Prior Validity using Monte Carlo estimation. Our proposed method (SD-VAE) performance significantly better than baselines.

(b) Predictive performance using encoded mean latent vector. Test LL and RMSE are reported.

most common semantics: *a)* ringbonds should satisfy cross-serial dependencies, *b)* explicit valence of atoms should not go beyond permitted.

### 5.1 Reconstruction Accuracy and Prior Validity

One metric of the soundness of VAE is to measure the ability of the model to encode data into a representation in the latent space and reconstruct the input by decoding from that point. Another metric is how often the model can decode into a valid data when the prior is randomly sampled. Since both encoding and decoding are stochastic, we follow the estimation by Monte Carlo method similar to that proposed in Kusner et al. [2017]. The CVAE and GVAE results are included directly from Kusner et al. [2017]. We show in the right part of Table 1a that our model produces a much higher rate of successful reconstruction, and a large increase in ratio of valid prior. Note that the results we reported only partially take the semantics of aromaticity into account. If we use an equivalent kekulization form of SMILES to train the model, then the valid portion of prior can go up to 97.3%.

### 5.2 Bayesian Optimization

The variational autoencoder realizes the conversion from molecule space to a continuous latent space via the encoder and vice versa via the decoder. This naturally leads to following two important applications: First, we can now train an extra model that predicts the data’s property from the representation in latent space, as suggested in Gómez-Bombarelli et al. [2016]. Second and more importantly, the continuous nature of latent space makes possible the optimization of finding new data with better properties. Following the protocol used in Kusner et al. [2017], we use Bayesian Optimization (BO) to search the molecules with desired properties in latent space.

In this section, we ask the model to optimize for octanol-water partition coefficients (a.k.a  $\log P$ ), an important measurement of drug-likeness of a given molecule. As Gómez-Bombarelli et al. [2016] suggests, for drug-likeness assessment  $\log P$  is penalized by other properties including synthetic accessibility score [Ertl and Schuffenhauer, 2009]. In Figure 3 we show the the top-3 best molecules found by each method, where our method found molecules with much better scores than previous works. Also the structures are richer than baselines (which mainly contain a chain structure).

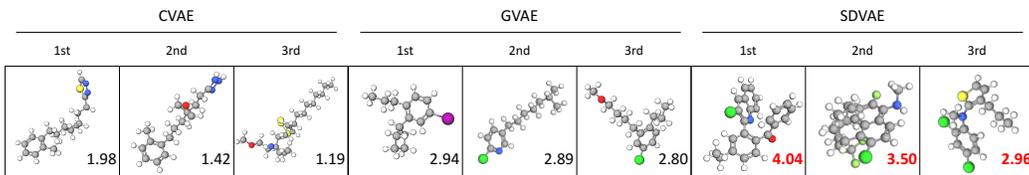


Figure 3: Best top-3 molecules and the corresponding scores found by each method using Bayesian Optimization.

### 5.3 Predictive performance of latent representation

We seek to know how well our latent space predicts the properties of molecules. We train the same sparse Gaussian Process as in Sec 5.2, with the same target value (namely the drug-likeness for molecules) for regression. We test the performance in the hold-out test dataset. In Table 1b, we report the result in Log Likelihood (LL) and Regression Mean Square Error (RMSE), which show that our SD-VAE always produces latent space that are more discriminative than both CVAE and GVAE baselines. This also shows that, with a properly designed decoder, the quality of encoder will also be improved via end2end training.

Similarity Metric	MorganFp	MACCS	PairFp	TopologicalFp
GVAE	<b>0.92 ± 0.10</b>	<b>0.83 ± 0.15</b>	0.94 ± 0.10	0.71 ± 0.14
SD-VAE	<b>0.92 ± 0.09</b>	<b>0.83 ± 0.13</b>	<b>0.95 ± 0.08</b>	<b>0.75 ± 0.14</b>

Table 2: Diversity as statistics from pair-wise distances measured as  $1 - s$ , where  $s$  is one of the similarity metrics. So higher values indicate better diversity. We show mean  $\pm$  stddev of  $\binom{100}{2}$  pairs among 100 molecules. Note that we report results from GVAE and our SD-VAE, because CVAE has very low valid priors, thus completely only failing this evaluation protocol.

#### 5.4 Diversity of generated molecules

Inspired by Benhenda [2017], here we seek to measure the diversity of generated molecules as an assessment of our methods. The intuition is that a good generative model should be able to generate diverse data and avoid model collapse in the learned space. In detail, we conduct this experiment in SMILES dataset, where we sample 100 points from the prior distribution, and for each point, we associate it as a molecule, which is the most frequent occurring valid SMILES decoded from 200 decoding attempts since the decoding is stochastic. We then, with one of several molecular similarity, compute the pair-wise similarity and report the mean and standard deviation in Table 2. We see both our method and] baseline do not have the model collapse problem, and the diversities with respect to different measurements are comparable. It indicates that although our method has more restricted decoding space than baselines, the diversity is not sacrificed. This is because we never rule-out the valid molecules.

#### 5.5 Visualizing the Latent Space

We seek to visualize the latent space as an assessment of how well our generative model is able to produce a coherent and smooth space of molecules. We visualize the latent space in 2 dimensions. We first embed a random molecule into latent space. Then we randomly generate 2 orthogonal unit vectors  $A$ . To get the latent representation of neighborhood, we interpolate the 2-D grid and project back to latent space with pseudo inverse of  $A$ . Finally we show decoded molecules. In Figure 4, we present two of such grid visualizations. Subjectively compared with figures in Kusner et al. [2017], our visualization is characterized by having smooth differences between neighboring molecules, and more complicated decoded structures.



Figure 4: Latent Space visualization. We start from the center molecule and decode the neighborhood latent vectors (neighborhood in projected 2D space).

#### References

David Alvarez-Melis and Tommi S Jaakkola. Tree-structured decoding with doubly-recurrent neural networks. 2016.

Mostapha Benhenda. Chemgan challenge for drug discovery: can ai reproduce natural chemical diversity? *arXiv preprint arXiv:1708.08227*, 2017.

- Esben Jannik Bjerrum. Smiles enumeration as data augmentation for neural network modeling of molecules. *arXiv preprint arXiv:1703.07076*, 2017.
- Joan Bresnan, Ronald M Kaplan, Stanley Peters, and Annie Zaenen. Cross-serial dependencies in dutch. 1982.
- Hanjun Dai, Bo Dai, and Le Song. Discriminative embeddings of latent variable models for structured data. In *ICML*, 2016.
- Hanjun Dai, Bo Dai, Yan-Ming Zhang, Shuang Li, and Le Song. Recurrent hidden semi-markov model. 2017.
- Li Dong and Mirella Lapata. Language to logical form with neural attention. *arXiv preprint arXiv:1601.01280*, 2016.
- David K Duvenaud, Dougal Maclaurin, Jorge Iparraguirre, Rafael Bombarell, Timothy Hirzel, Alán Aspuru-Guzik, and Ryan P Adams. Convolutional networks on graphs for learning molecular fingerprints. In *Advances in Neural Information Processing Systems*, pages 2215–2223, 2015.
- Peter Ertl and Ansgar Schuffenhauer. Estimation of synthetic accessibility score of drug-like molecules based on molecular complexity and fragment contributions. *Journal of cheminformatics*, 1(1):8, 2009.
- Rafael Gómez-Bombarelli, David Duvenaud, José Miguel Hernández-Lobato, Jorge Aguilera-Iparraguirre, Timothy D Hirzel, Ryan P Adams, and Alán Aspuru-Guzik. Automatic chemical design using a data-driven continuous representation of molecules. *arXiv preprint arXiv:1610.02415*, 2016.
- Gabriel Lima Guimaraes, Benjamin Sanchez-Lengeling, Pedro Luis Cunha Farias, and Alán Aspuru-Guzik. Objective-reinforced generative adversarial networks (organ) for sequence generation models. *arXiv preprint arXiv:1705.10843*, 2017.
- Zhiting Hu, Xueze Ma, Zhengzhong Liu, Eduard Hovy, and Eric Xing. Harnessing deep neural networks with logic rules. *arXiv preprint arXiv:1603.06318*, 2016.
- David Janz, Jos van der Westhuizen, and José Miguel Hernández-Lobato. Actively learning what makes a discrete sequence valid. *arXiv preprint arXiv:1708.04465*, 2017.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Donald E Knuth. Semantics of context-free languages. *Theory of Computing Systems*, 2(2):127–145, 1968.
- Matt J Kusner, Brooks Paige, and José Miguel Hernández-Lobato. Grammar variational autoencoder. *arXiv preprint arXiv:1703.01925*, 2017.
- Tao Lei, Wengong Jin, Regina Barzilay, and Tommi Jaakkola. Deriving neural architectures from sequence and graph kernels. *arXiv preprint arXiv:1705.09037*, 2017.
- Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016.
- Emilio Parisotto, Abdel-rahman Mohamed, Rishabh Singh, Lihong Li, Dengyong Zhou, and Pushmeet Kohli. Neuro-symbolic program synthesis. *arXiv preprint arXiv:1611.01855*, 2016.
- Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- Danilo J Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pages 1278–1286, 2014.
- Stuart M Shieber. Evidence against the context-freeness of natural language. 1985.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112, 2014.
- David Weininger. Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules. *Journal of chemical information and computer sciences*, 28(1):31–36, 1988.
- Xingxing Zhang, Liang Lu, and Mirella Lapata. Top-down tree long short-term memory networks. *arXiv preprint arXiv:1511.00060*, 2015.

# Appendix

## A Grammar for Modecule

Our syntax grammar is based on OpenSMILES standard, a generative context free grammar starting with  $\langle s \rangle$ .

$\langle s \rangle$	$\rightarrow \langle atom \rangle$
$\langle smiles \rangle$	$\rightarrow \langle chain \rangle$
$\langle atom \rangle$	$\rightarrow \langle bracket atom \rangle \mid \langle aliphatic organic \rangle \mid \langle aromatic organic \rangle$
$\langle aliphatic organic \rangle$	$\rightarrow \text{'B' } \mid \text{'C' } \mid \text{'N' } \mid \text{'O' } \mid \text{'S' } \mid \text{'P' } \mid \text{'F' } \mid \text{'I' } \mid \text{'Cl' } \mid \text{'Br'}$
$\langle aromatic organic \rangle$	$\rightarrow \text{'c' } \mid \text{'n' } \mid \text{'o' } \mid \text{'s'}$
$\langle bracket atom \rangle$	$\rightarrow \text{'[ ' } \langle baracekt atom (isotope) \rangle \text{' ]'}$
$\langle baracekt atom (isotope) \rangle$	$\rightarrow \langle isotope \rangle \langle symbol \rangle \langle baracekt atom (chiral) \rangle$ $\mid \langle symbol \rangle \langle baracekt atom (chiral) \rangle$ $\mid \langle isotope \rangle \langle symbol \rangle \mid \langle symbol \rangle$
$\langle baracekt atom (chiral) \rangle$	$\rightarrow \langle chiral \rangle \langle baracekt atom (h count) \rangle$ $\mid \langle baracekt atom (h count) \rangle$ $\mid \langle chiral \rangle$
$\langle baracekt atom (h count) \rangle$	$\rightarrow \langle h count \rangle \langle baracekt atom (charge) \rangle$ $\mid \langle baracekt atom (charge) \rangle$ $\mid \langle h count \rangle$
$\langle baracekt atom (charge) \rangle$	$\rightarrow \langle charge \rangle$
$\langle symbol \rangle$	$\rightarrow \langle aliphatic organic \rangle \mid \langle aromatic organic \rangle$
$\langle isotope \rangle$	$\rightarrow \langle digit \rangle \mid \langle digit \rangle \langle digit \rangle \mid \langle digit \rangle \langle digit \rangle \langle digit \rangle$
$\langle digit \rangle$	$\rightarrow \text{'1' } \mid \text{'2' } \mid \text{'3' } \mid \text{'4' } \mid \text{'5' } \mid \text{'6' } \mid \text{'7' } \mid \text{'8'}$
$\langle chiral \rangle$	$\rightarrow \text{'@' } \mid \text{'@@'}$
$\langle h count \rangle$	$\rightarrow \text{'H' } \mid \text{'H' } \langle digit \rangle$
$\langle charge \rangle$	$\rightarrow \text{'-' } \mid \text{'-' } \langle digit \rangle \mid \text{'+' } \mid \text{'+' } \langle digit \rangle$
$\langle bond \rangle$	$\rightarrow \text{'-' } \mid \text{'=' } \mid \text{'\#'} \mid \text{'/' } \mid \text{'\}'}$
$\langle ringbond \rangle$	$\rightarrow \langle digit \rangle$
$\langle branched atom \rangle$	$\rightarrow \langle atom \rangle \mid \langle atom \rangle \langle branches \rangle \mid \langle atom \rangle \langle ringbounds \rangle$ $\mid \langle atom \rangle \langle ringbounds \rangle \langle branches \rangle$
$\langle ringbounds \rangle$	$\rightarrow \langle ringbounds \rangle \langle ringbond \rangle \mid \langle ringbond \rangle$
$\langle branches \rangle$	$\rightarrow \langle branches \rangle \langle branch \rangle \mid \langle branch \rangle$
$\langle branch \rangle$	$\rightarrow \text{'(' } \langle chain \rangle \text{' )' } \mid \text{'(' } \langle bond \rangle \langle chain \rangle \text{' )'}$
$\langle chain \rangle$	$\rightarrow \langle branched atom \rangle \mid \langle chain \rangle \langle branched atom \rangle$ $\mid \langle chain \rangle \langle bond \rangle \langle branched atom \rangle$