

Variational Flow Graphical Model

Shaogang Ren, Belhal Karimi, Dingcheng Li, Ping Li

Cognitive Computing Lab
Baidu Research
10900 NE 8th St. Bellevue, WA 98004, USA

{renshaogang, belhal.karimi, dingchengl, pingli98}@gmail.com

Abstract

¹This paper introduces a novel approach to embed flow-based models with hierarchical structures. The proposed framework is named Variational Flow Graphical (VFG) Model. VFGs learn the representation of high dimensional data via a message-passing scheme by integrating flow-based functions through variational inference. By leveraging the expressive power of neural networks, VFGs produce a representation of the data using a lower dimension, thus overcoming the drawbacks of many flow-based models, usually requiring a high dimensional latent space involving many trivial variables. Aggregation nodes are introduced in the VFG models to integrate forward-backward hierarchical information via a message passing scheme. Maximizing the evidence lower bound (ELBO) of data likelihood aligns the forward and backward messages in each aggregation node achieving a consistency node state. Algorithms have been developed to learn model parameters through gradient updating regarding the ELBO objective.

The consistency of aggregation nodes enable VFGs to be applicable in tractable inference on graphical structures. Besides representation learning and numerical inference, VFGs provide a new approach for distribution modeling on datasets with graphical latent structures. Additionally, theoretical study shows that VFGs are universal approximators by leveraging the implicitly invertible flow-based structures. With flexible graphical structures and superior expressive power, VFGs could potentially be used to improve probabilistic inference.

In the experiments, VFGs achieves improved evidence lower bound (ELBO) and likelihood values on multiple datasets. We also highlight the benefits of our VFG model on missing entry imputation for datasets with graph structures. Multiple experiments on synthetic and real-world datasets confirm the benefits of the proposed method and potentially broad applications.

¹This work was initially submitted in 2020.

1 Introduction

Learning tractable distribution or density functions from datasets has broad applications. Probabilistic graphical models (PGMs) provide a unifying framework for capturing complex dependencies among random variables (Bishop and Nasrabadi, 2006; Wainwright and Jordan, 2008; Koller and Friedman, 2009). There are two general approaches for probabilistic inference with PGMs and other models: exact inference and approximate inference. In most cases, exact inference is either computationally involved or simply intractable. Variational inference (VI), stemmed from statistical physics, is computationally efficient and is applied to tackle large-scale inference problems (Anderson and Peterson, 1987; Hinton and van Camp, 1993; Jordan et al., 1999; Ghahramani and Beal, 1999; Hoffman et al., 2013; Blei et al., 2017; Fang and Li, 2021). In variational inference, mean-field approximation (Anderson and Peterson, 1987; Hinton and van Camp, 1993; Xing et al., 2003) and variational message passing (Bishop et al., 2003; Winn and Bishop, 2005) are two common approaches. These methods are limited by the choice of distributions that are inherently unable to recover the true posterior, often leading to a loose approximation.

To tackle the probabilistic inference problem, alternative models have been developed under the name of *tractable probabilistic models (TPMs)*. They include probabilistic decision graphs (Jaeger et al., 2006), arithmetic circuits (Darwiche, 2003), and-or search spaces (Marinescu and Dechter, 2005), multi-valued decision diagrams (Dechter and Mateescu, 2007), sum-product nets (Sánchez-Cauce et al., 2021), probabilistic sentential decision diagrams (Kisa et al., 2014), and probabilistic circuits (PCs) (Choi et al., 2020). PCs leverage the recursive mixture models and distributional factorization to establish tractable probabilistic inference. PCs also aim to attain a TPM with improved expressive power. The recent GFlowNets (Bengio et al., 2021) also target tractable probabilistic inference on different structures.

Apart from probabilistic inference, generative models have been developed to model high dimensional datasets and to learn meaningful hidden data representations by leveraging the approximation power of neural networks. These models also provide a possible approach to generate new samples from underlining distributions. Variational Auto-Encoders (VAEs) (Kingma and Welling, 2014) and Generative Adversarial Networks (GAN) (Goodfellow et al., 2014; Arjovsky and Bottou, 2017; Karras et al., 2019; Zhu et al., 2017; Yin et al., 2020; Ren et al., 2020) are widely applied to different categories of datasets. Flow-based models (Dinh et al., 2017, 2015; Rezende and Mohamed, 2015; van den Berg et al., 2018; Ren et al., 2021) leverage invertible neural networks and can estimate the density values of data samples as well. Energy-based models (EBMs) (Zhu et al., 1998; LeCun et al., 2006; Hinton, 2012; Xie et al., 2016; Nijkamp et al., 2019; Zhao et al., 2021; Zheng et al., 2021) define an unnormalized probability density function of data, which is the exponential of the negative energy function. Unlike TPMs, it is usually difficult to directly use generative models to perform probabilistic inference on datasets.

In this paper, we introduce VARIATIONAL FLOW GRAPHICAL (VFG) models. By leveraging the expressive power of neural networks, VFGs can learn latent representations from data. VFGs also follow the stream of tractable neural networks that allow to perform inference on graphical structures. Sum-product networks (Sánchez-Cauce et al., 2021) and probabilistic circuits (Choi et al., 2020) are falling into this type of models as well. Sum-product networks and probabilistic circuits depend on mixture models and probabilistic factorization in graphical structure for inference. Whereas, VFGs rely on the consistency of aggregation nodes in graphical structures to achieve tractable inference. Our contributions are summarized as follows.

Summary of contributions. Dealing with high dimensional data using graph structures exacerbates the systemic inability for effective distribution modeling and efficient inference. To overcome these limitations, we propose the VFG model to achieve the following goals:

- **Hierarchical and flow-based:** VFG is a novel graphical architecture uniting the hierarchical latent structures and flow-based models. Our model outputs a tractable posterior distribution used as an approximation of the true posterior of the hidden node states in the considered graph structure.
- **Distribution modeling:** Our theoretical analysis shows that VFGs are universal approximators. In the experiments, VFGs can achieve improved evidence lower bound (ELBO) and likelihood values by leveraging the implicitly invertible flow-based model structure.
- **Numerical inference:** Aggregation nodes are introduced in the model to integrate hierarchical information through a variational forward-backward message passing scheme. We highlight the benefits of our VFG model on applications: the missing entry imputation problem and the numerical inference on graphical data.

Moreover, experiments show that our model achieves to disentangle the factors of variation underlying high dimensional input data.

Roadmap: Section 2 presents important concepts used in the paper. Section 3 introduces the Variational Flow Graphical (VFG) model. The approximation property of VFGs is discussed in Section 4. Section 5 provides the algorithms used to train VFG models. Section 6 discusses how to perform inference with a VFG model. Section 7 showcases the advantages of VFG on various tasks. Section 8 and Section 9 provide a discussion and conclusion of the paper.

2 Preliminaries

We introduce the general principles and notations of variational inference and flow-based models in this section.

Notation: We use $[L]$ to denote the set $\{1, \dots, L\}$, for all $L > 1$. $\mathbf{KL}(p||q) := \int_{\mathcal{Z}} p(z) \log(p(z)/q(z))dz$ is the Kullback-Leibler divergence from q to p , two probability density functions defined on the set $\mathcal{Z} \subset \mathbb{R}^m$ for any dimension $m > 0$.

Variational Inference: Following the setting discussed above, the functional mapping $\mathbf{f} : \mathcal{Z} \rightarrow \mathcal{X}$ can be viewed as a decoding process and the mapping $\mathbf{f}^{-1} : \mathcal{X} \rightarrow \mathcal{Z}$ as an encoding one between random variables $\mathbf{z} \in \mathcal{Z}$ and $\mathbf{x} \in \mathcal{X}$ with densities $\mathbf{z} \sim p(\mathbf{z})$, $\mathbf{x} \sim p_{\theta}(\mathbf{x}|\mathbf{z})$. To learn the parameters θ , VI employs a parameterized family of so-called variational distributions $q_{\phi}(\mathbf{z}|\mathbf{x})$ to approximate the true posterior $p(\mathbf{z}|\mathbf{x}) \propto p(\mathbf{z})p_{\theta}(\mathbf{x}|\mathbf{z})$. The optimization problem of VI can be shown to be equivalent to maximizing the following evidence lower bound (ELBO) objective, noted $\mathcal{L}(\mathbf{x}; \theta, \phi)$:

$$\log p(\mathbf{x}) \geq \mathcal{L}(\mathbf{x}; \theta, \phi) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \mathbf{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})). \quad (1)$$

In Variational Auto-Encoders (VAEs, (Kingma and Welling, 2014; Rezende et al., 2014)), the calculation of the reconstruction term requires sampling from the posterior distribution along with using the reparameterization trick, i.e.,

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] \simeq \frac{1}{U} \sum_{u=1}^U \log p(\mathbf{x}|\mathbf{z}_u). \quad (2)$$

Here U is the number of latent variable samples drawn from the posterior $q_\phi(\mathbf{z}|\mathbf{x})$ regarding data \mathbf{x} .

Flow-based Models: Flow-based models (Dinh et al., 2017, 2015; Rezende and Mohamed, 2015; van den Berg et al., 2018) correspond to a probability distribution transformation using a sequence of invertible and differentiable mappings, noted $\mathbf{f} : \mathcal{Z} \rightarrow \mathcal{X}$. By defining the aforementioned invertible maps $\{\mathbf{f}_\ell\}_{\ell=1}^L$, and by the chain rule and inverse function theorem, the variable $\mathbf{x} = \mathbf{f}(\mathbf{z})$ has a tractable probability density function (pdf) given as:

$$\log p_\theta(\mathbf{x}) = \log p(\mathbf{z}) + \sum_{i=1}^L \log \left| \det \left(\frac{\partial \mathbf{h}^i}{\partial \mathbf{h}^{i-1}} \right) \right|, \quad (3)$$

where we have $\mathbf{h}^0 = \mathbf{x}$ and $\mathbf{h}^L = \mathbf{z}$ for conciseness. The scalar value $\log |\det(\partial \mathbf{h}^i / \partial \mathbf{h}^{i-1})|$ is the logarithm of the absolute value of the determinant of the Jacobian matrix $\partial \mathbf{h}^i / \partial \mathbf{h}^{i-1}$, also called the log-determinant. Eq. (3) yields a simple mechanism to build families of distributions that, from an initial density and a succession of invertible transformations, returns tractable density functions that one can sample from. Rezende and Mohamed (2015) propose an approach to construct flexible posteriors by transforming a simple base posterior with a sequence of flows. Firstly a stochastic latent variable is draw from base posterior $\mathcal{N}(\mathbf{z}_0|\mu(\mathbf{x}), \sigma(\mathbf{x}))$. With K flows, latent variable \mathbf{z}_0 is transformed to \mathbf{z}_k . The reformed EBLO is given by

$$\begin{aligned} \mathcal{L}(\mathbf{x}; \theta, \phi) &= \mathbb{E}_{q_\phi} [\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_{q_0} [\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_0(\mathbf{z}_0|\mathbf{x})] + \mathbb{E}_{q_0} \left[\sum_{k=1}^K \log \left| \det \left(\frac{\partial \mathbf{f}_k(\mathbf{z}_k; \psi_k)}{\partial \mathbf{z}_k} \right) \right| \right]. \end{aligned}$$

Here \mathbf{f}_k is the k -th flow with parameter ψ_k , i.e., $\mathbf{z}_K = \mathbf{f}_K \circ \dots \circ \mathbf{f}_2 \circ \mathbf{f}_1(\mathbf{z}_0)$. The flows are considered as functions of data sample \mathbf{x} , and they determine the final distribution in amortized inference. Several recent models have been proposed by leveraging the invertible flow-based models. Graphical normalizing flow (Wehenkel and Louppe, 2021) learns a DAG structure from the input data under sparse penalty and maximum likelihood estimation. The bivariate causal discovery method proposed in Khemakhem et al. (2021) relies on autoregressive structure of flow-based models and the asymmetry of log-likelihood ratio for cause-effect pairs. In this paper, we propose a framework that generalizes flow-based models (Dinh et al., 2017, 2015; Rezende and Mohamed, 2015; van den Berg et al., 2018) to graphical variable inference.

3 Variational Flow Graphical Model

Assume \mathbf{k} sections in the data samples, i.e., $\mathbf{x} = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)}]$, and a relationship among these sections and the corresponding latent variable. Then, it is possible to define a graphical model using normalizing flows, as introduced Section 2, leading to exact latent variable inference and log-likelihood evaluation of data samples.

A VFG model $\mathbb{G} = \{\mathcal{V}, \mathbf{f}\}$ consists of a node set (\mathcal{V}) and an edge set (\mathbf{f}). An edge can be either a flow function or an identity function. There are two types of nodes in a VFG: *aggregation* nodes and *non-aggregation* nodes. A non-aggregation node connects with another node with a flow function or an identity function. An aggregation node has multiple children, and it connects each of them with an identity function. Figure 1-Left gives an illustration of an aggregation node and Figure 1-Right shows a tree VFG model. Unlike classical graphical models, a node in a VFG model may represent a single variable or multiple variables. Moreover, each latent variable belongs to only one node in a VFG. In the following sections, identity function is considered as a special case of flow functions.

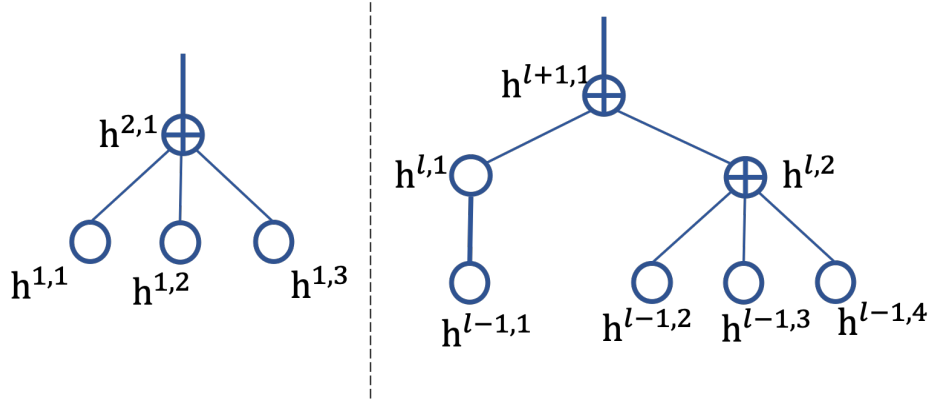


Figure 1: (Left) Node $h^{2,1}$ connects its children with invertible functions. Messages from the children are aggregated at the parent node, $h^{2,1}$. (Right) An illustration of the latent structure from layer $l-1$ to $l+1$. Thin lines are identity functions, and thick lines are flow functions. \oplus is an aggregation node, and circles stand for non-aggregation nodes.

3.1 Evidence Lower Bound of VFGs

We apply variational inference to learn model parameters θ from data samples. Different from VAEs, the recognition model (encoder) and the generative model (decoder) in a VFG share the same neural net structure and parameters. Moreover, the latent variables in a VFG lie in a hierarchy structure and are generated with deterministic flow functions.

We start with a tree VFG (Figure 2) to introduce the ELBO of the model. The hierarchical tree structure comprises L layers, \mathbf{h}^l denotes the latent state in layer l of the tree. We use $\mathbf{h}^{(j)}$ to represent node j 's latent state without specification of the layer number, and j is the node index in a tree or graph. The joint distribution for the hierarchical model is then

$$p_{\theta}(\mathbf{x}, \mathbf{h}) = p(\mathbf{h}^L)p(\mathbf{h}^{L-1}|\mathbf{h}^L) \cdots p(\mathbf{h}^1|\mathbf{h}^2)p(\mathbf{x}|\mathbf{h}^1).$$

where $\mathbf{h} = \{\mathbf{h}^1, \dots, \mathbf{h}^L\}$ denotes the set of latent states of the model. The hierarchical generative model is given by factorization $p(\mathbf{x}|\mathbf{h}^L) = p(\mathbf{x}|\mathbf{h}^1)\prod_{l=1}^{L-1}p(\mathbf{h}^l|\mathbf{h}^{l+1})$, and the prior distribution is $p(\mathbf{h}^L)$. Note that only the *root nodes* have *prior distributions*. The probabilistic density function $p(\mathbf{h}^{l-1}|\mathbf{h}^l)$ in the generative model is parameterized with one or multiple invertible flow functions. By leveraging the invertible flow functions, we use variational inference to approximate the posterior

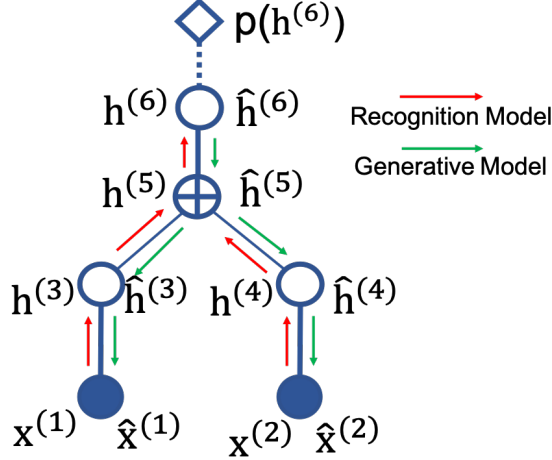


Figure 2: Forward message from data to approximate posterior distributions; generative model is realized by backward message from the root and generates the samples or reconstructions at each layer.

distribution of latent states. The hierarchical posterior (recognition model) is factorized as

$$q_{\theta}(\mathbf{h}|\mathbf{x}) = q(\mathbf{h}^1|\mathbf{x})q(\mathbf{h}^2|\mathbf{h}^1) \cdots q(\mathbf{h}^L|\mathbf{h}^{L-1}). \quad (4)$$

Evaluation of the posterior (recognition model) (4) involves forward information flows from the bottom of the tree to the top, and similarly, sampling the generative model takes the reverse direction.

By leveraging the hierarchical conditional independence in both generative model and posterior, the ELBO regarding the model is

$$\log p_{\theta}(\mathbf{x}) \geq \mathcal{L}(\mathbf{x}; \theta) = \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{h}^{1:L})] - \sum_{l=1}^L \mathbf{KL}^l. \quad (5)$$

Here \mathbf{KL}^l is the Kullback-Leibler divergence between the posterior and generative model in layer l . The first term in (5) evaluates data reconstruction. When $1 \leq l \leq L$,

$$\mathbf{KL}^l = \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log q(\mathbf{h}^l|\mathbf{h}^{l-1}) - \log p(\mathbf{h}^l|\mathbf{h}^{l+1})]. \quad (6)$$

When $l = L$, $\mathbf{KL}^L = \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log q(\mathbf{h}^L|\mathbf{h}^{L-1}) - \log p(\mathbf{h}^L)]$. It is easy to extend the computation of the ELBO (5) to DAGs with topology ordering of the nodes (and thus of the layers). Let $ch(i)$ and $pa(i)$ denote node i 's child set and parent set, respectively. Then, the ELBO for a DAG structure reads:

$$\mathcal{L}(\mathbf{x}; \theta) = \mathbb{E}_{q(\mathbf{h}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{h})] - \sum_{i \in \mathcal{V} \setminus \mathcal{R}_{\mathbb{G}}} \mathbf{KL}^{(i)} - \sum_{i \in \mathcal{R}_{\mathbb{G}}} \mathbf{KL}(q(\mathbf{h}^{(i)}|\mathbf{h}^{ch(i)}) || p(\mathbf{h}^{(i)})). \quad (7)$$

Here $\mathbf{KL}^{(i)} = \mathbb{E}_{q(\mathbf{h}|\mathbf{x})} [\log q(\mathbf{h}^{(i)}|\mathbf{h}^{ch(i)}) - \log p(\mathbf{h}^{(i)}|\mathbf{h}^{pa(i)})]$. $\mathcal{R}_{\mathbb{G}}$ is the set of root nodes of DAG $\mathbb{G} = \{\mathcal{V}, \mathbf{f}\}$. Assuming there are k leaf nodes on a tree or a DAG model, corresponding to k sections of the input sample $\mathbf{x} = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)}]$.

Maximizing the ELBO (5) or (7) equals to optimizing the parameters of the flows, θ . Similar to VAEs, we apply forward message passing (encoding) to approximate the posterior distribution of each layer's latent variables, and backward message passing (decoding) to generate the reconstructions as shown in Figure 2. For the following sections, we use \mathbf{h}^i to represent node i 's state in the forward message, and $\hat{\mathbf{h}}^i$ for node i 's state in the backward message. For all nodes, both \mathbf{h}^i and $\hat{\mathbf{h}}^i$ are sampled from the posterior. At the root nodes, we have $\hat{\mathbf{h}}^{\mathcal{R}} = \mathbf{h}^{\mathcal{R}}$.

3.2 Aggregation Nodes

There are two approaches to aggregate signals from different nodes: average-based and concatenation-based. We rather focus on average-based aggregation in this paper, and Figure 3 gives an example denoted by the operator \oplus . Let $\mathbf{f}_{(i,j)}$ be the direct edge (function) from node i to node j , and $\mathbf{f}_{(i,j)}^{-1}$ or $\mathbf{f}_{(j,i)}$ defined as its inverse function. Then, the aggregation operation at node i reads

$$\mathbf{h}^{(i)} = \frac{1}{|ch(i)|} \sum_{j \in ch(i)} \mathbf{f}_{(j,i)}(\mathbf{h}^{(j)}), \quad \hat{\mathbf{h}}^{(i)} = \frac{1}{|pa(i)|} \sum_{j \in pa(i)} \mathbf{f}_{(j,i)}(\hat{\mathbf{h}}^{(j)}). \quad (8)$$

Note that the above two equations hold even when node i has only one child or parent.

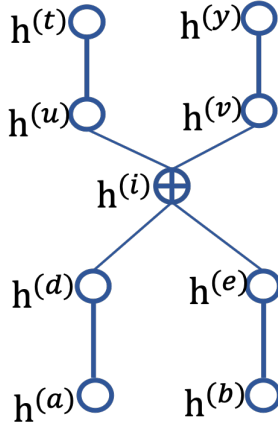


Figure 3: Aggregation node on a DAG VFG.

With the identity function between the parent and its children, there are *node consistency rules* regarding an average aggregation node: (a) a parent node's backward state equals the mean of its children's forward states, i.e., $\hat{\mathbf{h}}^{(i)} = \frac{1}{|ch(i)|} \sum_{j \in ch(i)} \mathbf{h}^{(j)}$; (b) a child node's forward state equals to the average of its parents' backward states, i.e., $\mathbf{h}^{(i)} = \frac{1}{|pa(i)|} \sum_{j \in pa(i)} \hat{\mathbf{h}}^{(j)}$. These rules empower VFGs with implicit invertibility.

We use aggregation node i in the DAG presented in Figure 3 as an example to illustrate node consistency. Node i has two parents, u and v ; and two children, d and e . Node i connects its parents and children with identity functions. According to (8), we have $\mathbf{h}^{(i)} = (\mathbf{h}^{(d)} + \mathbf{h}^{(e)})/2$ and $\hat{\mathbf{h}}^{(i)} = (\hat{\mathbf{h}}^{(u)} + \hat{\mathbf{h}}^{(v)})/2$. Here aggregation *consistency* means, for i 's children, their forward state should be consistent with i 's backward state, i.e.,

$$\mathbf{h}^{(d)} = \mathbf{h}^{(e)} = \hat{\mathbf{h}}^{(i)}. \quad (9)$$

For i 's parents, their backward state should be consistent with i 's forward state, i.e.,

$$\hat{\mathbf{h}}^{(u)} = \hat{\mathbf{h}}^{(v)} = \mathbf{h}^{(i)}. \quad (10)$$

We utilize the **KL** term in the ELBO (7) to ensure (9) and (10) can be satisfied during parameter updating. The **KL** term regarding node i is

$$\begin{aligned} \mathbf{KL}^{(i)} &= \mathbb{E}_{q(\mathbf{h}, \hat{\mathbf{h}}|\mathbf{x})} [\log q(\mathbf{h}^{(i)}|\mathbf{h}^{ch(i)}) - \log p(\mathbf{h}^{(i)}|\hat{\mathbf{h}}^{pa(i)})] \\ &\simeq \log q(\mathbf{h}^{(i)}|\mathbf{h}^{ch(i)}) - \log p(\mathbf{h}^{(i)}|\hat{\mathbf{h}}^{pa(i)}). \end{aligned} \quad (11)$$

As the term $\log q(\mathbf{h}^{(i)}|\mathbf{h}^{ch(i)})$ involves node states that are deterministic according to (8), it is omitted in the computation of (11). With Laplace as the latent state distribution, here

$$\begin{aligned} & \log p(\mathbf{h}^{(i)}|\hat{\mathbf{h}}^{pa(i)}) \\ &= \frac{1}{2}(\log p(\mathbf{h}^{(i)}|\hat{\mathbf{h}}^{(u)}) + p(\mathbf{h}^{(i)}|\hat{\mathbf{h}}^{(v)})) \\ &= \frac{1}{2}(-\|\mathbf{h}^{(i)} - \hat{\mathbf{h}}^{(u)}\|_1 - \|\mathbf{h}^{(i)} - \hat{\mathbf{h}}^{(v)}\|_1 - 2m \cdot \log 2). \end{aligned}$$

Hence minimizing $\mathbf{KL}^{(i)}$ is equal to minimizing $\{\|\mathbf{h}^{(i)} - \hat{\mathbf{h}}^{(u)}\|_1 + \|\mathbf{h}^{(i)} - \hat{\mathbf{h}}^{(v)}\|_1\}$ which achieves the consistent objective in (10).

Similarly, **KL**s of i 's children intend to realize consistency given in (9). We use node d as an example. The **KL** term regarding node d is

$$\begin{aligned} \mathbf{KL}^{(d)} &= \mathbb{E}_{q(\mathbf{h}, \hat{\mathbf{h}}|\mathbf{x})} [\log q(\mathbf{h}^{(d)}|\mathbf{h}^{ch(d)}) - \log p(\mathbf{h}^{(d)}|\hat{\mathbf{h}}^{pa(d)})] \\ &\simeq \log q(\mathbf{h}^{(d)}|\mathbf{h}^{ch(d)}) - \log p(\mathbf{h}^{(d)}|\hat{\mathbf{h}}^{pa(d)}). \end{aligned}$$

The first term $\log q(\mathbf{h}^{(d)}|\mathbf{h}^{ch(d)})$ is omitted in the calculation of $\mathbf{KL}^{(d)}$ due to the deterministic relation with (8). Knowing that

$$\begin{aligned} \log p(\mathbf{h}^{(d)}|\hat{\mathbf{h}}^{pa(d)}) &= \log p(\mathbf{h}^{(d)}|\hat{\mathbf{h}}^{(i)}) \\ &= -\|\mathbf{h}^{(d)} - \hat{\mathbf{h}}^{(i)}\|_1 - m \cdot \log 2, \end{aligned}$$

we notice that minimizing $\mathbf{KL}^{(d)}$ boils down to minimizing $\|\mathbf{h}^{(d)} - \hat{\mathbf{h}}^{(i)}\|_1$ that targets at (9). In summary, by maximizing the ELBO of a VFG, the aggregation consistency can be attained along with fitting the model to the data.

3.3 Implementation Details

The calculation of the data reconstruction term in (7) requires node states \mathbf{h}^i and $\hat{\mathbf{h}}^i$ ($\forall i \in \mathcal{V}$) from the posterior. They correspond to the encoding and decoding procedures in VAE model as shown in Eq. (2). At the root node, we have $\hat{\mathbf{h}}^{\mathcal{R}} = \mathbf{h}^{\mathcal{R}}$. The reconstruction terms in ELBO (7) can be computed with the backward message in the generative model $p(\mathbf{x}|\hat{\mathbf{h}}^1)$, i.e.,

$$\mathbb{E}_{q(\mathbf{h}, \hat{\mathbf{h}}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{h}, \hat{\mathbf{h}})] \simeq \frac{1}{U} \sum_{u=1}^U \log p(\mathbf{x}|\hat{\mathbf{h}}_u^{1:L}) = \frac{1}{U} \sum_{u=1}^U \log p(\mathbf{x}|\hat{\mathbf{h}}_u^{pa(x)}).$$

For a VFG model, we set $U = 1$. In the last term, $p(\mathbf{x}|\hat{\mathbf{h}}^{pa(x)})$ is either Gaussian or binary distribution parameterized with $\hat{\mathbf{x}}$ generated via the flow function with $\hat{\mathbf{h}}^{pa(x)}$ as the input.

4 Universal Approximation Property

A universal approximation power of coupling-layer based flows has been highlighted in Teshima et al. (2020). Following the analysis for flows Teshima et al. (2020), we prove that coupling-layer based VFGs have universal approximation as well. We first give several additional definitions regarding universal approximation. For a measurable mapping $\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^n$ and a subset $K \subset \mathbb{R}^m$, we define the following,

$$\|\mathbf{f}\|_{p,K} = \left(\int_K \|f(x)\|^p dx \right)^{1/p}.$$

Here $\|\cdot\|$ is the Euclidean norm of \mathbb{R}^n and $\|\mathbf{f}\|_{\text{sup},K} := \sup_{x \in K} \|\mathbf{f}(x)\|$.

Definition 4.1. (L^p -/ \sup -universality) Let \mathcal{M} be a model which is a set of measurable mappings from \mathbb{R}^m to \mathbb{R}^n . Let $p \in [1, \infty)$, and let \mathcal{G} be a set of measurable mappings $\mathbf{g} : U_{\mathbf{g}} \rightarrow \mathbb{R}^n$, where $U_{\mathbf{g}}$ is a measurable subset of \mathbb{R}^m which may depend on \mathbf{g} . We say that \mathcal{M} has the L^p -universal approximation property for \mathcal{G} if for any $\mathbf{g} \in \mathcal{G}$, any $\epsilon > 0$, and any compact subset $K \in U_{\mathbf{g}}$, there exists $\mathbf{f} \in \mathcal{M}$ such that $\|\mathbf{f} - \mathbf{g}\|_{p,K} < \epsilon$. We define the \sup -universality analogously by replacing $\|\cdot\|_{p,K}$ with $\|\cdot\|_{\sup,K}$.

Definition 4.2. (Immersion and submanifold) $\mathbf{g} : \mathfrak{M} \rightarrow \mathfrak{N}$ is said to be an immersion if $\text{rank}(\mathbf{g}) = m = \dim(\mathfrak{M})$ everywhere. If \mathbf{g} is injective (one-to-one) immersion, then \mathbf{g} establish an one-to-one correspondence of \mathfrak{M} and the subset $\tilde{\mathfrak{M}} = \mathbf{g}(\mathfrak{M})$ of \mathfrak{N} . If we use this correspondence to endow $\tilde{\mathfrak{M}}$ with a topology and C^∞ structure, then $\tilde{\mathfrak{M}}$ will be called a submanifold (or immersed submanifold) and $\mathbf{g} : \mathfrak{M} \rightarrow \tilde{\mathfrak{M}}$ is a diffeomorphism.

Definition 4.3. (C^r -diffeomorphisms for submanifold: \mathcal{Q}^r). We define \mathcal{Q}^r as the set of all C^r -diffeomorphisms $\mathbf{g} : U_{\mathbf{g}} \rightarrow \mathfrak{U}$, where $U_{\mathbf{g}} \subset \mathbb{R}^m$ is an open set C^r -diffeomorphic to \mathfrak{U} , which may depend on \mathbf{g} , and \mathfrak{U} is a submanifold of \mathbb{R}^n .

We use m to represent the root node dimension of a VFG, and n to denote the dimension of data samples. VFGs learn the data manifold embedded in \mathbb{R}^n . We define $\mathcal{C}_c^\infty(\mathbb{R}^{m-1})$ as the set of all compactly-supported C^∞ mappings from \mathbb{R}^{m-1} to \mathbb{R} . For a function set \mathcal{T} , we define \mathcal{T} -ACF as the set of affine coupling flows [Teshima et al. \(2020\)](#) that are assembled with functions in \mathcal{T} , and we use $\text{VFG}_{\mathcal{T}\text{-ACF}}$ to represent the set of VFGs constructed using flows in \mathcal{T} -ACF.

Theorem 4.1. (L^p -universality) Let $p \in [0, \infty)$. Assume \mathcal{H} is a \sup -universal approximator for $\mathcal{C}_c^\infty(\mathbb{R}^{m-1})$, and that it consists of C^1 -functions. Then $\text{VFG}_{\mathcal{H}\text{-ACF}}$ is an L^p -universal approximator for \mathcal{Q}_c^0 .

Proof. We construct a VFG structure that forms a mapping from \mathbb{R}^m to \mathbb{R}^n . Let $r = n \bmod m$.

If $r = 0$, it is easy to construct a one-layer tree VFG \mathbf{f} (\mathbf{f} also represents the function/edge set) and the root as an aggregation node. The children divide the n input entries into $\tau = n/m$ even sections, and each section connects the aggregation node with a flow function.

Given an injective immersion $\mathbf{g} : \mathfrak{M} \rightarrow \mathfrak{N}$, function \mathbf{g} can be represented with the concatenation of a set of functions, i.e., $\mathbf{g} = [\mathbf{g}_1, \dots, \mathbf{g}_\tau]^\top$, each invertible \mathbf{g}_i has dimension m . According to the function decomposition theory [Kuo et al. \(2010\)](#), its inverse can be represent as the summation of functions $\mathbf{g}_i^{-1}, 1 \leq i \leq \tau$, i.e., $\mathbf{g}^{-1} = \frac{1}{\tau} \sum_{i=1}^{\tau} \mathbf{g}_i^{-1}$. For each \mathbf{g}_i , and $\tilde{\mathfrak{M}}_i = \mathbf{g}_i(\mathfrak{M})$ is a submanifold in \mathfrak{N} , and it is diffeomorphic to \mathfrak{M} . According to Theorem 2 in [Teshima et al. \(2020\)](#), \mathcal{H} -ACF is an universal approximator for each $\mathbf{g}_i, 1 \leq i \leq \tau$. Therefore, VFG \mathbf{f} has universal approximation for immersion $\mathbf{g} : \mathfrak{M} \rightarrow \mathfrak{N}$.

If $r \neq 0$, let $\tau = \lfloor n/m \rfloor$. We divide the τ -th section and the remaining r entries into two equal small sections that are denoted with τ and $\tau + 1$. Sections τ and $\tau + 1$ have r overlapped entries. Similarly, we can construct an one-layer VFG \mathbf{f} with $\tau + 1$ children, and each child takes a section as the input.

The input coordinate index of \mathbf{g}_τ in \mathbb{R}^m is $I_\tau = [1, 2, \dots, \lceil (m+r)/2 \rceil]$, and the output index of \mathbf{g}_τ in \mathbb{R}^n is $I_\tau + \gamma = [\gamma + 1, \gamma + 2, \dots, \gamma + \lceil (m+r)/2 \rceil]$, and $\gamma = (\tau - 1)m$. The input coordinate index of $\mathbf{g}_{\tau+1}$ in \mathbb{R}^m is $I_{\tau+1} = [m - \lceil (m+r)/2 \rceil + 1, \dots, m - 1, m]$, and the output index of $\mathbf{g}_{\tau+1}$ in \mathbb{R}^n is $I_{\tau+1} + \gamma$. We can see that the m dimensions are divided into two sets, the overlapped set $O = [m - \lceil (m+r)/2 \rceil + 1, \lceil (m+r)/2 \rceil]$, and the remaining set R containing the rest dimensions.

The mapping $\mathbf{g} : \mathfrak{M} \rightarrow \mathfrak{N}$ can be decomposed into $\tau + 1$ functions, i.e., $\mathbf{g} = [\mathbf{g}_1, \dots, \mathbf{g}_\tau, \mathbf{g}_{\tau+1}]^\top$, and the inverse \mathbf{g}^{-1} is adjusted here: $\mathbf{g}_j^{-1} = \frac{1}{\omega} \sum_{i=1}^{\omega} \mathbf{g}_{i(j)}^{-1}$. When $j \in O$, $\omega = \tau + 1$, and all \mathbf{g}_i^{-1} s will

be involved; when $j \in R$, $\omega = \tau$, and either \mathbf{g}_τ^{-1} or $\mathbf{g}_{\tau+1}^{-1}$ is omitted due to the missing of entry j in the function output. The mapping \mathbf{g}_τ is a diffeomorphism from manifold \mathfrak{M}_τ ($\mathfrak{M}_\tau \subset \mathfrak{M}$) to sub-manifold $\tilde{\mathfrak{M}}_\tau$ in \mathfrak{N} . Similarly $\mathbf{g}_{\tau+1}$ is a diffeomorphism from $\mathfrak{M}_{\tau+1}$ to manifold $\tilde{\mathfrak{M}}_{\tau+1}$. For each \mathbf{g}_i , $1 \leq i \leq \tau + 1$, it can be universally approximated with a function in $\mathcal{H} - ACF$ Teshima et al. (2020). Hence, we construct a VFG with universal approximation for any \mathbf{g} in \mathcal{Q}_c^0 . \square

With the conditions in Theorem 4.1, $\text{VFG}_{\mathcal{H}-ACF}$ is a distributional universal approximator as well (Teshima et al., 2020).

5 The Proposed Algorithms

In this section, we develop the training algorithm (Algorithm 1) to maximize the ELBO objective function (7). In Algorithm 1, the inference of the latent states is performed via forwarding message passing, cf. Line 6, and their reconstructions are computed in backward message passing, cf. Line 11. A VFG is a deterministic network passing latent variable values between nodes. Ignoring explicit neural network parameterized variances for all latent nodes enables us to use flow-based models as both the encoders and decoders. Hence, we obtain a deterministic ELBO objective (5)-(7) that can efficiently be optimized with standard stochastic optimizers.

Algorithm 1 Inference model parameters with forward and backward message propagation

```

1: Input: Data distribution  $\mathcal{D}$ ,  $\mathbb{G} = \{\mathcal{V}, \mathbf{f}\}$ 
2: for  $s = 0, 1, \dots$  do
3:   Sample minibatch  $b$  samples  $\{\mathbf{x}_1, \dots, \mathbf{x}_b\}$  from  $\mathcal{D}$ ;
4:   for  $i \in \mathcal{V}$  do
5:     // forward message passing
6:      $\mathbf{h}^{(i)} = \frac{1}{|ch(i)|} \sum_{j \in ch(i)} \mathbf{f}_{(j,i)}(\mathbf{h}^{(j)})$ ;
7:   end for
8:    $\hat{\mathbf{h}}^{(i)} = \mathbf{h}^{(i)}$  if  $i \in \mathcal{R}_{\mathbb{G}}$  or  $i \in \text{layer L}$ ;
9:   for  $i \in \mathcal{V}$  do
10:    // backward message passing
11:     $\hat{\mathbf{h}}^{(i)} = \frac{1}{|pa(i)|} \sum_{j \in pa(i)} \mathbf{f}_{(i,j)}^{-1}(\hat{\mathbf{h}}^{(j)})$ ;
12:   end for
13:    $\mathbf{h} = \{\mathbf{h}^{(t)} | t \in \mathcal{V}\}$ ,  $\hat{\mathbf{h}} = \{\hat{\mathbf{h}}^{(t)} | t \in \mathcal{V}\}$ ;
14:   Approximate the KL terms in ELBO for each layer with  $b$  samples;
15:   Updating VFG model  $\mathbb{G}$  with gradient ascending:  $\theta_{\mathbf{f}}^{(s+1)} = \theta_{\mathbf{f}}^{(s)} + \nabla_{\theta_{\mathbf{f}}} \frac{1}{b} \sum_{i=1}^b \mathcal{L}(\mathbf{x}_b; \theta_{\mathbf{f}}^{(s)})$ .
16: end for
```

In training Algorithm 1, the backward variable state $\hat{\mathbf{h}}^l$ in layer l is generated according to $p(\hat{\mathbf{h}}^l | \hat{\mathbf{h}}^{l+1})$, and at the root layer, node state $\hat{\mathbf{h}}^{\mathcal{R}}$ is set equal to $\mathbf{h}^{\mathcal{R}}$ that is from the posterior $q(\mathbf{h} | \mathbf{x})$, not from the prior $p(\mathbf{h}^{\mathcal{R}})$. So we can see all the forward and backward latent variables are sampled from the posterior $q(\mathbf{h} | \mathbf{x})$.

From a practical perspective, layer-wise training strategy can improve the accuracy of a model especially when it is constructed of more than two layers. In such a case, the parameters of only one layer are updated with backpropagation of the gradient of the loss function while keeping the other layers fixed at each optimization step. By maximizing the ELBO (7) with the above algorithm, the node consistency rules in Section 3.2 are expected to be satisfied.

5.1 Improve Training of VFG

The inference ability of VFG can be reinforced by masking out some sections of the training samples. The training objective can be changed to force the model to impute the value of the masked sections. For example in a tree model, the alternative objective function reads

$$\begin{aligned} \mathcal{L}(\mathbf{x}, O_{\mathbf{x}}; \theta) = & \sum_{t: 1 \leq t \leq k, t \notin O} \mathbb{E}_{q(\mathbf{h}, \hat{\mathbf{h}} | \mathbf{x}^{O_{\mathbf{x}}})} \left[\log p(\mathbf{x}^{(t)} | \hat{\mathbf{h}}^1) \right] \\ & - \sum_{l=1}^{L-1} \mathbb{E}_{q(\mathbf{h}, \hat{\mathbf{h}} | \mathbf{x})} \left[\log q(\mathbf{h}^l | \mathbf{h}^{l-1}) - \log p(\mathbf{h}^l | \hat{\mathbf{h}}^{l+1}) \right] \\ & - \mathbf{KL}(q(\mathbf{h}^L | \mathbf{h}^{L-1}) | p(\mathbf{h}^L)). \end{aligned} \quad (12)$$

where $O_{\mathbf{x}}$ is the index set of leaf nodes with observation, and $\mathbf{x}^{O_{\mathbf{x}}}$ is the union of observed data sections. The random-masking training procedure for objective (12) is described in Algorithm 2. In practice, we use Algorithm 2 along with Algorithm 1 to enhance the training of a VFG model. However, we only occasionally update the model parameter θ with the gradient of (12) to ensure the distribution learning running well.

Algorithm 2 Inference model parameters with random masking

```

1: Input: Data distribution  $\mathcal{D}$ ,  $\mathbb{G} = \{\mathcal{V}, \mathbf{f}\}$ 
2: for  $s = 0, 1, \dots$  do
3:   Sample minibatch  $b$  samples  $\{\mathbf{x}_1, \dots, \mathbf{x}_b\}$  from  $\mathcal{D}$ ;
4:   Optimize (5) with Line 4 to Line 15 in Algorithm 1;
5:   Sample a subset of the  $k$  data sections as data observation set  $O_{\mathbf{x}}$ ;  $O \leftarrow O_{\mathbf{x}}$ ;
6:   for  $i \in \mathcal{V}$  do
7:     // forward message passing
8:      $\mathbf{h}^{(i)} = \frac{1}{|ch(i) \cap O|} \sum_{j \in ch(i) \cap O} \mathbf{f}_{(j,i)}(\mathbf{h}^{(j)})$ ;
9:      $O \leftarrow O \cup \{i\}$  if  $ch(i) \cap O \neq \emptyset$ ;
10:  end for
11:   $\hat{\mathbf{h}}^{(i)} = \mathbf{h}^{(i)}$  if  $i \in \mathcal{R}_{\mathbb{G}}$  or  $i \in \text{layer } L$ ;
12:  for  $i \in \mathcal{V}$  do
13:    // backward message passing
14:     $\hat{\mathbf{h}}^{(i)} = \frac{1}{|pa(i)|} \sum_{j \in pa(i)} \mathbf{f}_{(i,j)}^{-1}(\hat{\mathbf{h}}^{(j)})$ ;
15:  end for
16:   $\mathbf{h} = \{\mathbf{h}^{(t)} | t \in \mathcal{V} \cap O\}$ ,  $\hat{\mathbf{h}} = \{\hat{\mathbf{h}}^{(t)} | t \in \mathcal{V}\}$ ;
17:  Approximate the  $\mathbf{KL}$  terms in ELBO for each layer with  $b$  samples;
18:  Updating VFG with gradient of (12):  $\theta_{\mathbf{f}}^{(s+1)} = \theta_{\mathbf{f}}^{(s)} + \nabla_{\theta_{\mathbf{f}}} \frac{1}{b} \sum_{i=1}^b \mathcal{L}(\mathbf{x}_b, O_{\mathbf{x}}; \theta_{\mathbf{f}}^{(s)})$ ,
19: end for
```

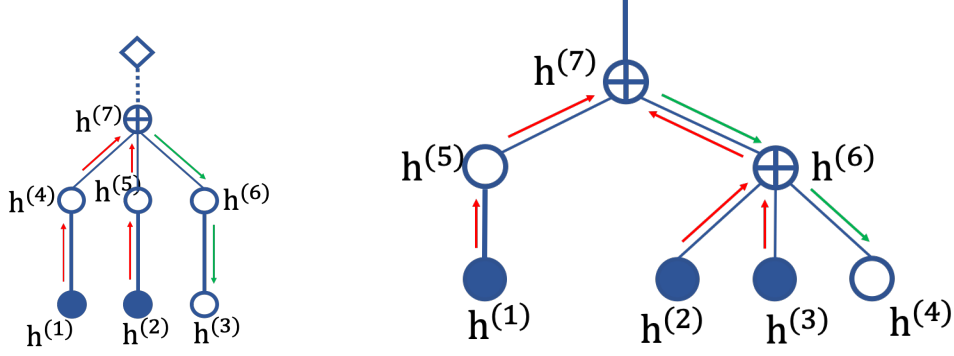


Figure 4: (Left) Inference on a VFG with single aggregation node. Node 7 aggregates information from node 1 and 2, and passes down the update to node 3 for prediction. (Right) Inference on a tree VFG. Observed node states are gathered at node 7 to predict the state of node 4. Red and green lines are forward and backward messages, respectively.

6 Inference on VFGs

With a VFG, we aim to infer node states given observed ones. The hidden state of a parent node j in $l = 1$ can be computed with the observed children as follows:

$$\mathbf{h}^{(j)} = \frac{1}{|ch(j) \cap O|} \sum_{i \in ch(j) \cap O} \mathbf{h}^{(i)}, \quad (13)$$

where O is the set of observed leaf nodes, see Figure 4-left for an illustration. Observe that for either a tree or a DAG, the state of any hidden node is updated via messages received from its children. After reaching the root node, we can update any nodes with backward message passing. Figure 4 illustrates this inference mechanism for trees in which the structure enables us to perform message passing among the nodes. We derive the following lemma establishing the relation between two leaf nodes.

Lemma 6.1. *Let \mathbb{G} be a tree VFG with L layers, and i and j are two leaf nodes with a as the closest common ancestor node. Given observed value at node i , the value of node j can be approximated by $\hat{\mathbf{x}}^j = \mathbf{f}_{(a,j)}(\mathbf{f}_{(i,a)}(\mathbf{x}^{(i)}))$. Here $\mathbf{f}_{(i,a)}$ is the flow function path from node i to node a .*

Proof. According to the aggregation operation (8) discussed in Section 3.2, at an aggregation node a , the reconstruction state of a child node j is the mean reconstruction state averaging the backward messages from the parent nodes. The reconstruction of the child node j can be calculated with the average reconstruction state regarding its parent node. Apply it sequentially, we have $\hat{\mathbf{x}}^{(j)} = \mathbf{f}_{(a,j)}(\hat{\mathbf{h}}^a)$. The forward state of node a can be computed by sequentially applying forward aggregating starting from its observed descendent i , i.e., $\mathbf{h}^{(a)} = \mathbf{f}_{(i,a)}(\mathbf{x}^{(i)})$. As there are no other observations, with forward and backward message passing to and from the root node, at node a , we have $\mathbf{h}^{(a)} = \hat{\mathbf{h}}^a$. Therefore, we have $\hat{\mathbf{x}}^{(j)} = \mathbf{f}_{(a,j)}(\mathbf{f}_{(i,a)}(\mathbf{x}^{(i)}))$. \square

Considering the flow-based model (3), we have the following identity for each node of the graph structure:

$$\begin{aligned} p(\mathbf{h}^{(i)}|\mathbf{h}^{pa(i)}) &= p(\mathbf{h}^{pa(i)}) \left| \det\left(\frac{\partial \mathbf{h}^{pa(i)}}{\partial \mathbf{h}^{(i)}}\right) \right| \\ &= p(\mathbf{h}^{pa(i)}) \left| \det(\mathbf{J}_{\mathbf{h}^{pa(i)}}(\mathbf{h}^{(i)})) \right|. \end{aligned}$$

Lemma 6.1 provides an approach to conduct inference on a tree and impute missing values in the data. It is easy to extend the inference method to DAG VFGs.

7 Numerical Experiments

In this section, we provide several studies to validate the proposed VFG models. The first application we present is missing value imputation. We compare our method with different baseline models on several datasets. The second set of experiments is to evaluate VFG models on three different datasets, i.e., MNIST, Caltech101, and Omniglot, with ELBO and likelihoods as the score. The third application we present here is the task of learning posterior distribution of the latent variables corresponding to the hidden explanatory factors of variations in the data (Bengio et al., 2013). For that latter application, the model is trained and evaluated on the MNIST handwritten digits dataset.

In this paper, we would rather assume the VFG graph structures are given and fixed. In the following experiments, the VFG structures are given in the dataset or designed heuristically (as other neural networks) for the sake of numerical illustrations. Learning the structure of VFG is an interesting research problem and is left for future works. A simple approach for VFG structure learning is to regularize the graph with the DAG structure penalty (Zheng et al., 2018; Wehenkel and Louppe, 2021).

All the experiments are conducted on NVIDIA-TITAN X (Pascal) GPUs. In the experiments, we use the same coupling block (Dinh et al., 2017) to construct different flow functions. The coupling block consists of three fully connected layers (of dimension 64) separated by two RELU layers along with the coupling trick. Each flow function has block number $\mathcal{B} > 3$.

7.1 Evaluation on Inference with Missing Entries Imputation

We now focus on the task of imputing missing entries in a graph structure. For all the following experiments, the models are trained on the training set and are used to infer the missing entries of samples in the testing set. We first study the proposed VFGs on two datasets without given graph structures, and we compare VFGs with several conventional methods that do not require the graph structures in the data. We then compare VFGs with graphical models that can perform inference on explicit graphs.

7.1.1 Synthetic Dataset

In this set of experiments, we study different methods with synthetic datasets. The baselines for this set of experiments include mean value method (Means), iterative imputation (Iterative) (Buck, 1960), and multivariate imputation by chained equation (MICE) (Van Buuren and Groothuis-Oudshoorn, 2011). Mean Squared Error as the metric of reference in order to compare the different methods for the imputation task. We use the baseline implementations in Pedregosa et al. (2011) in the experiments.

We generate 10 synthetic datasets (using different seeds) of 1,300 data points, 1,000 for the training phase of the model, 300 for imputation testing. Each data sample has 8 dimensions with 2 latent variables. Let $z_1 \sim \mathcal{N}(0, 1.0^2)$ and $z_2 \sim \mathcal{N}(1.0, 2.0^2)$ be the latent variables. For a sample \mathbf{x} , we have $x_1 = x_2 = z_1$, $x_3 = x_4 = 2\sin(z_1)$, $x_5 = x_6 = z_2$, and $x_7 = x_8 = z_2^2$. In the testing dataset, x_3 , x_4 , x_7 , and x_8 are missing. We use a VFG model with a single average aggregation node that has four children, and each child connects the parent with a flow function consisting of 3 coupling layers (Dinh et al., 2017). Each child takes 2 variables as input data section, and the latent dimension of the VFG is 2. We compare, in Figure 5, our VFG method with the baselines described above using boxplots on obtained MSE values for those 10 simulated datasets. We can see that the proposed VFG model performs much better than mean value, iterative, and MICE methods. Figure 5 shows that VFGs also demonstrates more performance robustness compared against other methods.

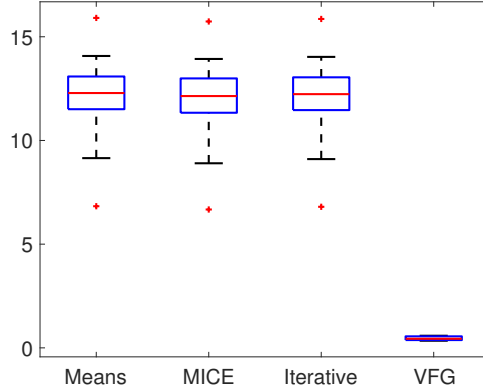


Figure 5: Synthetic datasets: MSE boxplots of VFG and baseline methods.

7.1.2 California Housing Dataset

We further investigate the method on a real dataset. The California Housing dataset has 8 feature entries and 20,640 data samples. We use the first 20,000 samples for training and 100 of the rest for testing. We get 4 data sections, and each section contains 2 variables. In the testing set, the second section is assumed missing for illustration purposes, as the goal is to impute this missing section. In addition to the three baselines in introduced the main file, we also compared with KNN (k-nearest neighbor) method. Again, we use the implementations from Pedregosa et al. (2011) for the baselines in this set of experiments.

The VFG structure is designed heuristically. We construct a tree structure VFG with 2 layers. The first layer has two aggregation nodes, and each of them has two children. The second layer consists of one aggregation node that has two children connecting with the first layer. Each flow function has $\mathcal{B} = 4$ coupling blocks. Table 1 shows that our model yields significantly better results than any other method in terms of prediction error. It indicates that with the help of universal approximation power of neural networks, VFGs have superior inference capability.

Table 1: California Housing dataset: Imputation Mean Squared Error (MSE) results.

<i>Methods</i>	<i>Imputation MSE</i>
Mean Value	1.993
MICE	1.951
Iterative Imputation	1.966
KNN (k=5)	1.969
VFG	1.356

7.1.3 Comparison with Graphical Models

In this set of experiments, we use a synthetic Gaussian graphical model dataset from the bnlearn package (Scutari, 2009) to evaluate the proposed model. The data graph structure is given. The dataset consists of 7 variables and 5,000 samples. Sample values at each node are generated according to a structured causal model with a diagram given by Figure 6.

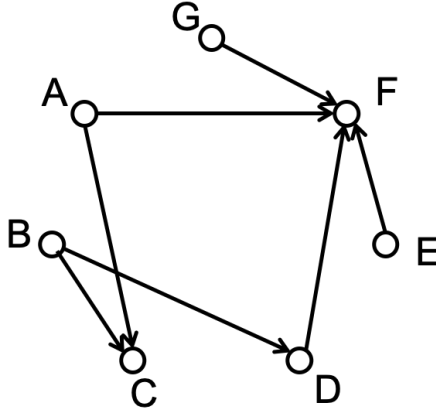


Figure 6: Graph structure for Gaussian graphical model dataset.

In Figure 6, each node represents a variable generated with a function of its parent nodes. For instance, node V is generated with $V = \mathbf{f}(pa(V), N_V)$. Here $pa(V)$ is the set of V 's parents, and N_V is a noise term for V . A node without any parent is determined only by the noise term. $\mathbf{f}()$ is V 's generating function, and only linear functions are used in this dataset. All the noise terms are Normal distributions.

We take Bayesian network implementation (Scutari, 2009) and sum-product network (SPN) package (Molina et al., 2019; Poon and Domingos, 2011) as experimental baselines. 4500 samples are used for training, and the rest 500 samples are for testing. The structure of VFG is designed based on the directed graph given by Figure 6. In the imputation task, we take Node 'F' as the missing entry, and use the values of other node to impute the missing entry. Table 2 gives the imputation results from the three methods. We can see that VFG achieves the smallest prediction error. Besides the imputation MSE, Table 2 also gives the prediction error variance. Compared against Bayesian net and SPN, VFG achieves much smaller performance variance. It means VFGs are much more stable in this set of experiments.

Table 2: Gaussian graphical model dataset: Imputation Mean Squared Error (MSE) and Variance results.

<i>Methods</i>	Bayesian Net	SPN	VFG
Imputation MSE	1.059	0.402	0.104
Imputation Variance	2.171	0.401	0.012

7.2 ELBO and Likelihood

We further qualitatively compare our VFG model with existing methods on data distribution learning and variational inference using three standard datasets. The baselines we compare in this experiment are VAE (Kingma and Welling, 2014), Planer (Rezende and Mohamed, 2015), IAF (Kingma et al., 2016), and SNF (van den Berg et al., 2018). The evaluation datasets and setup are following two standard flow-based variational models, Sylvester Normalizing Flows (van den Berg et al., 2018) and (Rezende and Mohamed, 2015). We use a tree VFG with structure as shown in Figure 7 for three datasets.

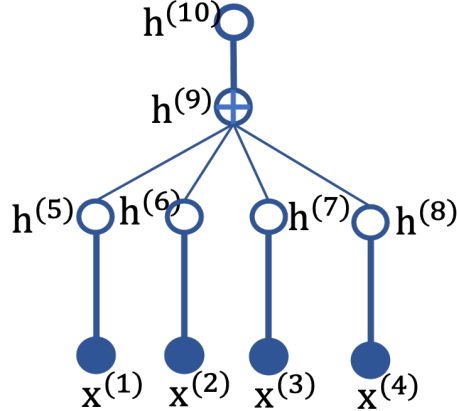


Figure 7: MIST Tree structure.

We train the tree VFG with the following ELBO objective that incorporate a β coefficient for the \mathbf{KL} terms. Empirically, a small β yields better ELBO and NLL values, and we set β around 0.1 in the experiments. Recall that

$$\text{ELBO} = \mathcal{L}(\mathbf{x}; \theta) = \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{h}^{1:L})] - \beta \sum_{l=1}^L \mathbf{KL}^l.$$

Table 3 presents the negative evidence lower bound (-ELBO) and the estimated negative likelihood (NLL) for all methods on three datasets: MNIST, Caltech101, and Omniglot. The baseline methods are VAE based methods enhanced with normalizing flows. They use 16 flows to improve the posterior estimation. SNF is orthogonal Sylvester flow method with a bottleneck of $M = 32$. We set the VFG coupling block (Dinh et al., 2017) number with $\mathcal{B} = 4$, and following (van den Berg et al., 2018) we run multiple times to get the mean and standard derivation as well. VFG can achieve superior EBLO as well as NLL values on all three datasets compared against the baselines as given in Table 3. VFGs can achieve better variational inference and data distribution modeling results (ELBOs and NLLs) in Table 3 in part due to VFGs’ universal approximation power as given

Table 3: Numerical values of negative log-likelihood and free energy (negative evidence lower bound) for static MNIST, Caltech101, and Omniglot datasets.

Model	MNIST		Caltech101		Omniglot	
	-ELBO	NLL	-ELBO	NLL	-ELBO	NLL
VAE (Kingma and Welling, 2014)	86.55 \pm 0.06	82.14 \pm 0.07	110.80 \pm 0.46	99.62 \pm 0.74	104.28 \pm 0.39	97.25 \pm 0.23
Planer (Rezende and Mohamed, 2015)	86.06 \pm 0.31	81.91 \pm 0.22	109.66 \pm 0.42	98.53 \pm 0.68	102.65 \pm 0.42	96.04 \pm 0.28
IAF (Kingma et al., 2016)	84.20 \pm 0.17	80.79 \pm 0.12	111.58 \pm 0.38	99.92 \pm 0.30	102.41 \pm 0.04	96.08 \pm 0.16
SNF (van den Berg et al., 2018)	83.32 \pm 0.06	80.22 \pm 0.03	104.62 \pm 0.29	93.82 \pm 0.62	99.00 \pm 0.04	93.77 \pm 0.03
VFG (ours)	80.80 \pm 0.76	63.66 \pm 0.14	67.26 \pm 0.53	65.74 \pm 0.84	80.16 \pm 0.73	78.65 \pm 0.66

in Theorem 4.1. Also, the intrinsic approximate invertible property of VFGs ensures the decoder or generative model in a VFG to achieve smaller reconstruction errors for data samples and hence smaller NLL values.

7.3 Latent Representation Learning on MNIST

In this set of experiments, we evaluate VFGs on latent representation learning of the MNIST dataset (LeCun et al.). We construct a tree VFG model depicted in Figure 7. In the first layer, there are 4 flow functions, and each of them takes 14×14 image blocks as the input. Thus a 28×28 input image is divided into four 14×14 blocks as the input of VFG model. We use $\mathcal{B} = 4$ for all the flows. The latent dimension for this model is $m = 196$. Following Sorrenson et al. (2020), the VFG model is trained with image labels to learn the latent representation of the input data. We set the parameters of \mathbf{h}^L 's prior distribution as a function of image label, i.e., $\lambda^L(u)$, where u denotes the image label. In practice, we use 10 trainable λ^L s regarding the 10 digits. The images in the second

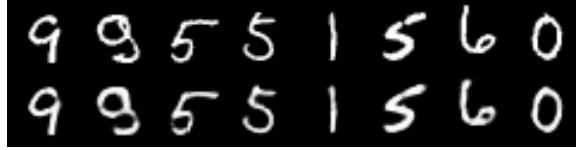


Figure 8: (Top) original MNIST digits. (Bottom) reconstructed images using VFG.

row of Figure 8 are reconstructions of MNIST samples extracted from the testing set, displayed in the first row of the same Figure, using our proposed VFG model.

Figure 9-Left shows t-distributed stochastic neighbor embedding (t-SNE) (van der Maaten and Hinton, 2008) plot of 2,000 testing images' latent variables learned with our model, and 200 for each digit. Figure 9-Left illustrates that VFG can learn separated latent representations to distinguish different hand-written numbers. For comparison, we also present the results of a baseline model. The baseline model (coupling-based flow) is constructed using the same coupling block and similar number of parameters as VFGs but with 28×28 as the input and latent dimension. Figure 9-Right gives the baseline coupling-layer-based flow training and testing with the same procedures. These show that coupling-based flow cannot give a clear division between some digits, e.g., 1 and 2, 7 and 9 due to the bias introduced by the high-dimensional redundant latent variables.

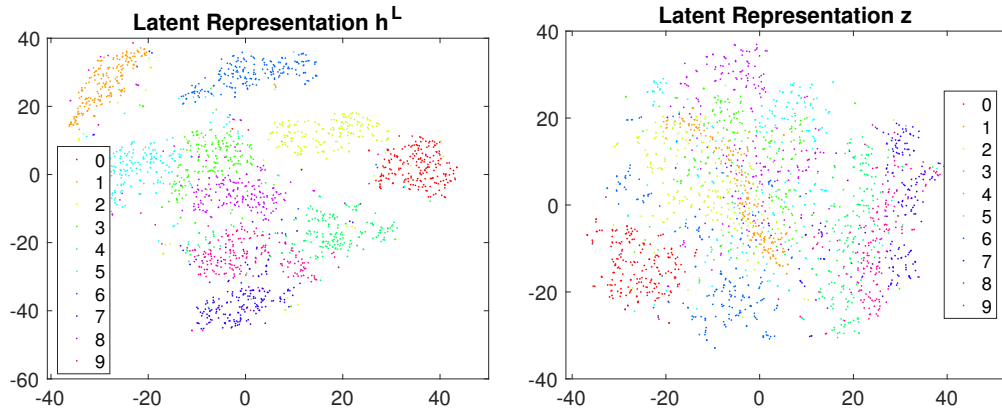


Figure 9: t-SNE of latent variables for VFG (Left) and coupling-layer based flow (Right) on MNIST.

To provide a description of the learned latent representation, we first obtain the root latent variables of a set of images through forward message passing. Each latent variable's values are changed increasingly within a range centered at the value of the latent variable obtained from last step. This perturbation is performed for each image in the set. Figure 10 shows the change of images by increasing one latent variable from a small value to a larger one. The figure presents some of the latent variables that have obvious effects on images, and most of the $m = 196$ variables do not impact the generation significantly. Latent variables $i = 6$ and $i = 60$ control the digit width. Variable $i = 19$ affects the brightness. $i = 92, i = 157$ and some of the variables not displayed here control the style of the generated digits.

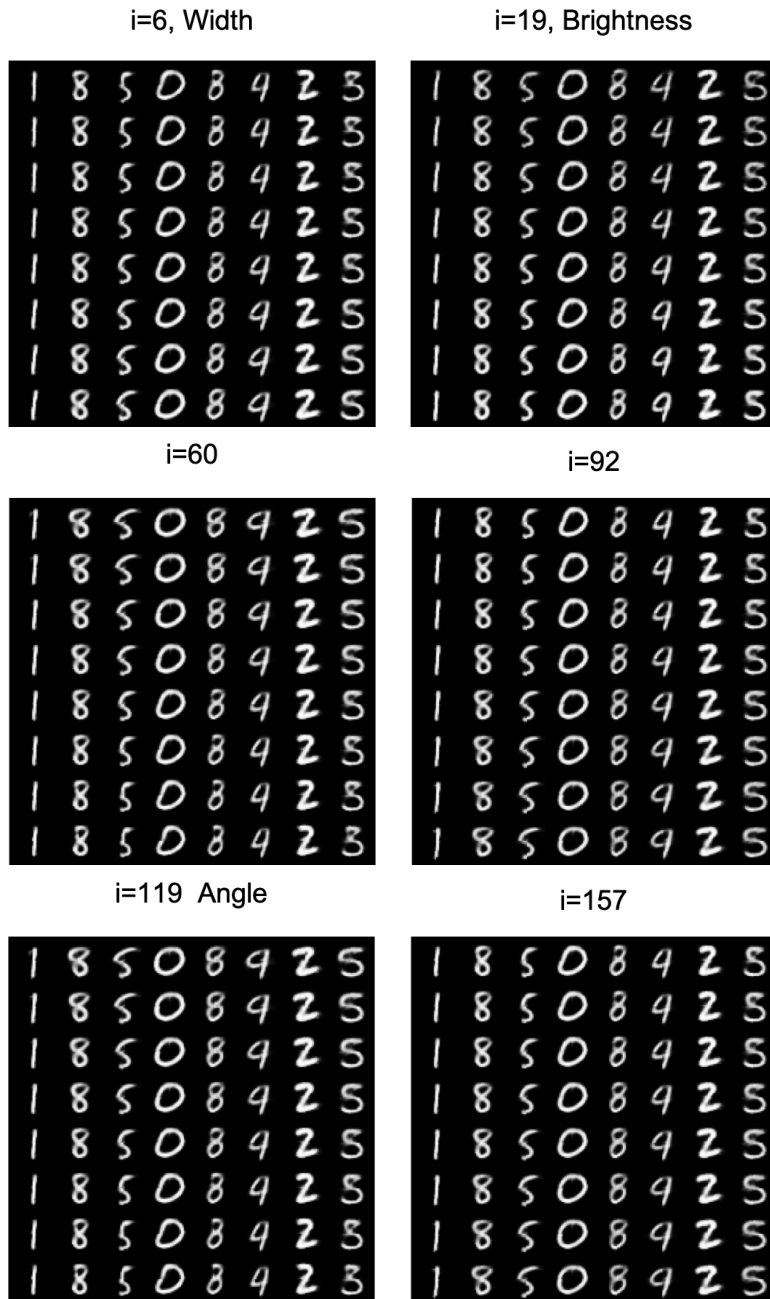


Figure 10: MNIST: Increasing each latent variable from a small value to a larger one.

8 Discussion

One of the motivations for proposing our VFG algorithm is to develop a tractable model that can be used for distribution learning and posterior inference. As long as the node states in the aggregation nodes are consistent, we can always apply VFGs in order to infer missing values. We provide more discussion on the structures of VFGs in the sequel.

8.1 Benefits of Encoder-decoder Parameter Sharing

There are several advantages for the encoder and decoder to share parameters. Firstly, it makes the network’s structure simple. Secondly, the training and inference can be simplified with concise and simple graph structures. Thirdly, by leveraging invertible flow-based functions, VFGs obtain tighter ELBOs in comparison with VAE based models. The intrinsic invertibility introduced by flow functions ensures the decoder or generative model in a VFG achieves smaller reconstruction errors for data samples and hence smaller NLL values and tighter ELBO. Whereas without the intrinsic constraint of invertibility or any help or regularization from the encoder, VAE-based models have to learn an unassisted mapping function (decoder) to reconstruct all data samples with the latent variables, and there are always some discrepancy errors in the reconstruction that lead to relatively larger NLL values and hence inferior ELBOs.

8.2 Structures of VFGs

In the experiments, the model structures have been chosen heuristically and for the sake of numerical illustrations. A tree VFG model can be taken as a dimension reduction model that is available for missing value imputation as well. Variants of those structures will lead to different numerical results and at this point, we can not claim any generalization regarding the impact of the VFG structure on the outputs. Meanwhile, learning the structure of VFG is an interesting research problem and is left for future works. VFG structures could be learned through the regularization of DAG structures (Zheng et al., 2018; Wehenkel and Louppe, 2021).

VFGs rely on minimizing the KL term to achieve *consistency* in aggregation nodes. As long as the aggregation nodes retain consistency, the model always has a tight ELBO and can be applied to tractable posterior inference. According to Teshima et al. (2020), coupling-based flows are endowed with the universal approximation power. Hence, we believe that the consistency of aggregation nodes on a VFG can be attained with a tight ELBO.

9 Conclusion

In this paper, we propose VFG, a variational flow graphical model that aims at bridging the gap between flow-based models and the paradigm of graphical models. Our VFG model learns data distribution and latent representation through message passing between nodes in the model structure. We leverage the power of invertible flow functions in any general graph structure to simplify the inference step of the latent nodes given some input observations. We illustrate the effectiveness of our variational model through experiments. Future work includes applying our VFG model to relational data structure learning and reasoning.

Appendix

A ELBO of Tree VFGs

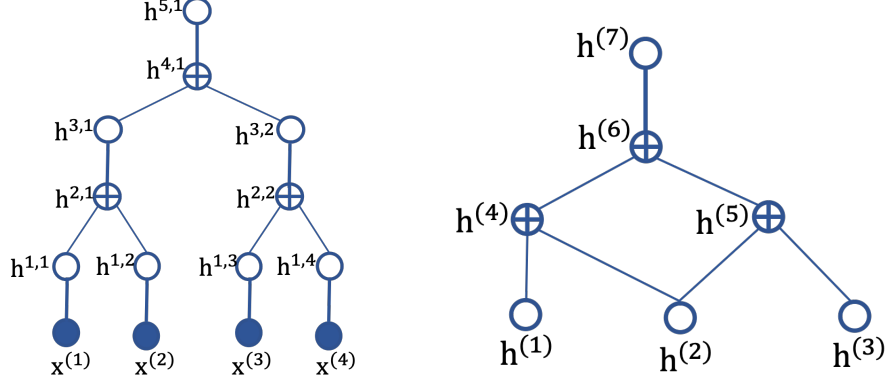


Figure 11: (Left) A tree VFG with $L = 5$ and three aggregation nodes. (Right) A DAG with inverse topology order $\{ \{1,2,3\}, \{4,5\}, \{6\}, \{7\} \}$, and they correspond to layers 0 to 3.

Let each data sample has k sections, i.e., $\mathbf{x} = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)}]$. VFGs are graphical models that can integrate different sections or components of the dataset. We assume that for each pair of connected nodes, the edge is an invertible flow function. The vector of parameters for all the edges is denoted by θ . The forward message passing starts from \mathbf{x} and ends at \mathbf{h}^L , and backward message passing in the reverse direction. We start with the hierarchical generative tree network structure illustrated by an example in Figure 11-Left. Then the marginal likelihood term of the data reads

$$p(\mathbf{x}|\theta) = \sum_{\mathbf{h}^1, \dots, \mathbf{h}^L} p(\mathbf{h}^L|\theta) p(\mathbf{h}^{L-1}|\mathbf{h}^L, \theta) \dots p(\mathbf{x}|\mathbf{h}^1, \theta).$$

The hierarchical generative model is given by factorization

$$p(\mathbf{h}) = p(\mathbf{h}^L) \prod_{l=1}^{L-1} p(\mathbf{h}^l|\mathbf{h}^{l+1}). \quad (14)$$

The probability density function $p(\mathbf{h}^{l-1}|\mathbf{h}^l)$ in the generative model is modeled with one or multiple invertible normalizing flow functions. The hierarchical posterior (recognition network) is factorized as

$$q_\theta(\mathbf{h}|\mathbf{x}) = q(\mathbf{h}^1|\mathbf{x}) q(\mathbf{h}^2|\mathbf{h}^1) \dots q(\mathbf{h}^L|\mathbf{h}^{L-1}). \quad (15)$$

Draw samples from the generative model (14) involves sequential conditional sampling from the top of the tree to the bottom, and computation of the recognition model (15) takes the reverse direction. Notice that

$$q(\mathbf{h}|\mathbf{x}) = q(\mathbf{h}^1|\mathbf{x}) q(\mathbf{h}^{2:L}|\mathbf{h}^1).$$

With the hierarchical structure of a tree, we further have

$$q(\mathbf{h}^{l:L}|\mathbf{h}^{l-1}) = q(\mathbf{h}^l|\mathbf{h}^{l-1}) q(\mathbf{h}^{l+1:L}|\mathbf{h}^l) = q(\mathbf{h}^l|\mathbf{h}^{l-1}) q(\mathbf{h}^{l+1:L}|\mathbf{h}^l) \quad (16)$$

$$p(\mathbf{h}^{l:L}) = p(\mathbf{h}^l|\mathbf{h}^{l+1:L}) p(\mathbf{h}^{l+1:L}) = p(\mathbf{h}^l|\mathbf{h}^{l+1}) p(\mathbf{h}^{l+1:L}) \quad (17)$$

By leveraging the conditional independence in the chain structures of both recognition and generative models, the derivation of trees' ELBO becomes easier.

$$\begin{aligned}
\log p(\mathbf{x}) &= \log \int p(\mathbf{x}|\mathbf{h})p(\mathbf{h})d\mathbf{h} \\
&= \log \int \frac{q(\mathbf{h}|\mathbf{x})}{q(\mathbf{h}|\mathbf{x})} p(\mathbf{x}|\mathbf{h})p(\mathbf{h})d\mathbf{h} \\
&\geq \mathbb{E}_{q(\mathbf{h}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{h}) - \log q(\mathbf{h}|\mathbf{x}) + \log p(\mathbf{h})] \\
&= \mathcal{L}(x; \theta).
\end{aligned}$$

The last step is due to the Jensen inequality. With $\mathbf{h} = \mathbf{h}^{1:L}$,

$$\begin{aligned}
\log p(\mathbf{x}) &\geq \mathcal{L}(x; \theta) \\
&= \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{h}^{1:L}) - \log q(\mathbf{h}^{1:L}|\mathbf{x}) + \log p(\mathbf{h}^{1:L})] \\
&= \underbrace{\mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{h}^{1:L})]}_{\text{Reconstruction of data}} - \underbrace{\mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log q(\mathbf{h}^{1:L}|\mathbf{x}) - \log p(\mathbf{h}^{1:L})]}_{\mathbf{KL}^{1:L}}
\end{aligned} \tag{18}$$

With conditional independence in the hierarchical structure, we have

$$q(\mathbf{h}^{1:L}|\mathbf{x}) = q(\mathbf{h}^{2:L}|\mathbf{h}^1\mathbf{x})q(\mathbf{h}^1|\mathbf{x}) = q(\mathbf{h}^{2:L}|\mathbf{h}^1)q(\mathbf{h}^1|\mathbf{x}).$$

The second term of (18) can be further expanded as

$$\begin{aligned}
\mathbf{KL}^{1:L} &= \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log q(\mathbf{h}^1|\mathbf{x}) + \log q(\mathbf{h}^{2:L}|\mathbf{h}^1) \\
&\quad - \log p(\mathbf{h}^1|\mathbf{h}^{2:L}) - \log p(\mathbf{h}^{2:L})].
\end{aligned} \tag{19}$$

Similarly, with conditional independence of the hierarchical latent variables, $p(\mathbf{h}^1|\mathbf{h}^{2:L}) = p(\mathbf{h}^1|\mathbf{h}^2)$. Thus

$$\begin{aligned}
\mathbf{KL}^{1:L} &= \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log q(\mathbf{h}^1|\mathbf{x}) - \log p(\mathbf{h}^1|\mathbf{h}^2) \\
&\quad + \log q(\mathbf{h}^{2:L}|\mathbf{h}^1) - \log p(\mathbf{h}^{2:L})] \\
&= \underbrace{\mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log q(\mathbf{h}^1|\mathbf{x}) - \log p(\mathbf{h}^1|\mathbf{h}^2)]}_{\mathbf{KL}^1} \\
&\quad + \underbrace{\mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log q(\mathbf{h}^{2:L}|\mathbf{h}^1) - \log p(\mathbf{h}^{2:L})]}_{\mathbf{KL}^{2:L}}.
\end{aligned}$$

We can further expand the $\mathbf{KL}^{2:L}$ term following similar conditional independent rules regarding the tree structure. At level l , we get

$$\mathbf{KL}^{l:L} = \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log q(\mathbf{h}^{l:L}|\mathbf{h}^{l-1}) - \log p(\mathbf{h}^{l:L})].$$

With (16) and (17), it is easy to show that

$$\begin{aligned}
\mathbf{KL}^{l:L} &= \underbrace{\mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log q(\mathbf{h}^l|\mathbf{h}^{l-1}) - \log p(\mathbf{h}^l|\mathbf{h}^{l+1})]}_{\mathbf{KL}^l} \\
&\quad + \underbrace{\mathbb{E}_{q(\mathbf{h}^{l:L}|\mathbf{x})} [\log q(\mathbf{h}^{l+1:L}|\mathbf{h}^l) - \log p(\mathbf{h}^{l+1:L})]}_{\mathbf{KL}^{l+1:L}}.
\end{aligned} \tag{20}$$

The ELBO (18) can be written as

$$\mathcal{L}(\mathbf{x}; \theta) = \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{h}^{1:L})] - \sum_{l=1}^{L-1} \mathbf{KL}^l - \mathbf{KL}^L. \quad (21)$$

When $1 \leq l \leq L-1$

$$\mathbf{KL}^l = \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log q(\mathbf{h}^l|\mathbf{h}^{l-1}) - \log p(\mathbf{h}^l|\mathbf{h}^{l+1})]. \quad (22)$$

According to conditional independence, the expectation regarding variational distribution layer l just depends on layer $l-1$. We can simplify the expectation each term of (21) with the default assumption that all latent variables are generated regarding data sample \mathbf{x} . Therefore the ELBO (21) can be simplified as

$$\mathcal{L}(\mathbf{x}; \theta) = \mathbb{E}_{q(\mathbf{h}^1|\mathbf{x})} [\log p(\mathbf{x}|\hat{\mathbf{h}}^1)] - \sum_{l=1}^L \mathbf{KL}^l. \quad (23)$$

The \mathbf{KL} term (22) becomes

$$\mathbf{KL}^l = \mathbb{E}_{q(\mathbf{h}^l|\mathbf{h}^{l-1})} [\log q(\mathbf{h}^l|\mathbf{h}^{l-1}) - \log p(\mathbf{h}^l|\hat{\mathbf{h}}^{l+1})].$$

When $l = L$,

$$\mathbf{KL}^L = \mathbb{E}_{q(\mathbf{h}^L|\mathbf{h}^{L-1})} [\log q(\mathbf{h}^L|\mathbf{h}^{L-1}) - \log p(\mathbf{h}^L)].$$

B ELBO of DAG VFGs

Note that if we reverse the edge directions in a DAG, the resulting graph is still a DAG graph. The nodes can be listed in a topological order regarding the DAG structure as shown in Figure 11-Right.

By taking the topology order as the layers in tree structures, we can derive the ELBO for DAG structures. Assume the DAG structure has L layers, and the root nodes are in layer L . We denote by \mathbf{h} the vector of latent variables, then following (18) we develop the ELBO as

$$\begin{aligned} \log p(\mathbf{x}) &\geq \mathcal{L}(\mathbf{x}; \theta) \\ &= \mathbb{E}_{q(\mathbf{h}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{h})}{q(\mathbf{h}|\mathbf{x})} \right] \\ &= \underbrace{\mathbb{E}_{q(\mathbf{h}|\mathbf{x})} \left[\log p(\mathbf{x}|\mathbf{h}) \right]}_{\text{Reconstruction of the data}} - \underbrace{\mathbb{E}_{q(\mathbf{h}|\mathbf{x})} \left[\log q(\mathbf{h}|\mathbf{x}) - \log p(\mathbf{h}) \right]}_{\mathbf{KL}}. \end{aligned} \quad (24)$$

Similarly the KL term can be expanded as in the tree structures. For nodes in layer l

$$\mathbf{KL}^{l:L} = \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log q(\mathbf{h}^{l:L}|\mathbf{h}^{1:l-1}) - \log p(\mathbf{h}^{l:L})].$$

Note that $ch(l)$ may include nodes from layers lower than $l-1$, and $pa(l)$ may include nodes from layers higher than l . Some nodes in l may not have parent. Based on conditional independence with the topology order of a DAG, we have

$$q(\mathbf{h}^{l:L}|\mathbf{h}^{1:l-1}) \quad (25)$$

$$\begin{aligned} &= q(\mathbf{h}^l|\mathbf{h}^{1:l-1})q(\mathbf{h}^{l+1:L}|\mathbf{h}^l) \\ &= q(\mathbf{h}^l|\mathbf{h}^{1:l-1})q(\mathbf{h}^{l+1:L}|\mathbf{h}^{1:l})p(\mathbf{h}^{l:L}) \\ &= p(\mathbf{h}^l|\mathbf{h}^{l+1:L})p(\mathbf{h}^{l+1:L}) \end{aligned} \quad (26)$$

Following (20) and with (25-26), we have

$$\begin{aligned} \mathbf{KL}^{l:L} = & \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log q(\mathbf{h}^l|\mathbf{h}^{1:l-1}) - \log p(\mathbf{h}^l|\mathbf{h}^{l+1:L})] \\ & + \underbrace{\mathbb{E}_{q(\mathbf{h}^{l:L}|\mathbf{x})} [\log q(\mathbf{h}^{l+1:L}|\mathbf{h}^{1:l}) - \log p(\mathbf{h}^{l+1:L})]}_{\mathbf{KL}^{l+1:L}}. \end{aligned}$$

Furthermore,

$$q(\mathbf{h}^l|\mathbf{h}^{1:l-1}) = q(\mathbf{h}^l|\mathbf{h}^{ch(l)}), \quad p(\mathbf{h}^l|\mathbf{h}^{l+1:L}) = p(\mathbf{h}^l|\mathbf{h}^{pa(l)}).$$

Hence,

$$\mathbf{KL}^{l:L} = \underbrace{\mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log q(\mathbf{h}^l|\mathbf{h}^{ch(l)}) - \log p(\mathbf{h}^l|\mathbf{h}^{pa(l)})]}_{\mathbf{KL}^l} + \mathbf{KL}^{l+1:L} \quad (27)$$

For nodes in layer l ,

$$\mathbf{KL}^l = \sum_{i \in l} \underbrace{\mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log q(\mathbf{h}^{(i)}|\mathbf{h}^{ch(i)}) - \log p(\mathbf{h}^{(i)}|\mathbf{h}^{pa(i)})]}_{\mathbf{KL}^{(i)}}.$$

Recurrently applying (27) to (24) yields

$$\begin{aligned} \mathcal{L}(\mathbf{x}; \theta) = & \mathbb{E}_{q(\mathbf{h}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{h})] - \sum_{i \in \mathcal{V} \setminus \mathcal{R}_G} \mathbf{KL}^{(i)} \\ & - \sum_{i \in \mathcal{R}_G} \mathbf{KL}(q(\mathbf{h}^{(i)}|\mathbf{h}^{ch(i)})||p(\mathbf{h}^{(i)})). \end{aligned}$$

For node i ,

$$\mathbf{KL}^{(i)} = \mathbb{E}_{q(\mathbf{h}|\mathbf{x})} [\log q(\mathbf{h}^{(i)}|\mathbf{h}^{ch(i)}) - \log p(\mathbf{h}^{(i)}|\mathbf{h}^{pa(i)})].$$

References

- James R Anderson and Carsten Peterson. A mean field theory learning algorithm for neural networks. *Complex Systems*, 1:995–1019, 1987.
- Martin Arjovsky and Léon Bottou. Towards principled methods for training generative adversarial networks. *arXiv preprint arXiv:1701.04862*, 2017.
- Yoshua Bengio, Aaron C. Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(8):1798–1828, 2013.
- Yoshua Bengio, Tristan Deleu, Edward J Hu, Salem Lahlou, Mo Tiwari, and Emmanuel Bengio. Gflownet foundations. *arXiv preprint arXiv:2111.09266*, 2021.
- Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*, volume 4. Springer, 2006.
- Christopher M Bishop, David Spiegelhalter, and John Winn. Vibes: A variational inference engine for bayesian networks. In *NeurIPS*, pages 793–800, 2003.
- David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017.
- Samuel F Buck. A method of estimation of missing values in multivariate data suitable for use with an electronic computer. *Journal of the Royal Statistical Society: Series B (Methodological)*, 22(2):302–306, 1960.
- YooJung Choi, Antonio Vergari, and Guy Van den Broeck. Probabilistic circuits: A unifying framework for tractable probabilistic models. Technical report, 2020.
- Adnan Darwiche. A differential approach to inference in bayesian networks. *J. ACM*, 50(3):280–305, 2003.
- Rina Dechter and Robert Mateescu. AND/OR search spaces for graphical models. *Artif. Intell.*, 171(2-3):73–106, 2007.
- Laurent Dinh, David Krueger, and Yoshua Bengio. NICE: non-linear independent components estimation. In *Proceedings of the 3rd International Conference on Learning Representations (ICLR Workshop)*, San Diego, CA, 2015.
- Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real NVP. In *Proceedings of the 5th International Conference on Learning Representations (ICLR)*, Toulon, France, 2017.
- Guanhua Fang and Ping Li. On variational inference in biclustering models. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*, pages 3111–3121, Virtual Event, 2021.
- Zoubin Ghahramani and Matthew J. Beal. Variational inference for bayesian mixtures of factor analysers. In *Advances in Neural Information Processing Systems (NIPS)*, pages 449–455, Denver, CO, 1999.

- Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2672–2680, Montreal, Canada, 2014.
- Geoffrey E Hinton. A practical guide to training restricted boltzmann machines. In *Neural networks: Tricks of the trade*, pages 599–619. Springer, 2012.
- Geoffrey E. Hinton and Drew van Camp. Keeping the neural networks simple by minimizing the description length of the weights. In Lenny Pitt, editor, *Proceedings of the Sixth Annual ACM Conference on Computational Learning Theory (COLT)*, pages 5–13, Santa Cruz, CA, 1993.
- Matthew D. Hoffman, David M. Blei, Chong Wang, and John W. Paisley. Stochastic variational inference. *J. Mach. Learn. Res.*, 14(1):1303–1347, 2013.
- Manfred Jaeger, Jens Dalgaard Nielsen, and Tomi Silander. Learning probabilistic decision graphs. *Int. J. Approx. Reason.*, 42(1-2):84–100, 2006.
- Michael I. Jordan, Zoubin Ghahramani, Tommi S. Jaakkola, and Lawrence K. Saul. An introduction to variational methods for graphical models. *Mach. Learn.*, 37(2):183–233, 1999.
- Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019.
- Ilyes Khemakhem, Ricardo Pio Monti, Robert Leech, and Aapo Hyvärinen. Causal autoregressive flows. In *Proceedings of the 24th International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 3520–3528, Virtual Event, 2021.
- Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In *Proceedings of the 2nd International Conference on Learning Representations (ICLR)*, Banff, Canada, 2014.
- Diederik P. Kingma, Tim Salimans, Rafal Józefowicz, Xi Chen, Ilya Sutskever, and Max Welling. Improved variational autoencoders with inverse autoregressive flow. In *Advances in Neural Information Processing Systems (NIPS)*, pages 4736–4744, Barcelona, Spain, 2016.
- Doga Kisa, Guy Van den Broeck, Arthur Choi, and Adnan Darwiche. Probabilistic sentential decision diagrams. In *Proceedings of the Fourteenth International Conference on Principles of Knowledge Representation and Reasoning (KR)*, Vienna, Austria, 2014.
- Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- Frances Y. Kuo, Ian H. Sloan, Grzegorz W. Wasilkowski, and Henryk Wozniakowski. On decompositions of multivariate functions. *Math. Comput.*, 79(270):953–966, 2010.
- Yann LeCun, Corinna Cortes, and Christopher J.C. Burges. MNIST handwritten digit database. URL <http://yann.lecun.com/exdb/mnist/>.
- Yann LeCun, Sumit Chopra, Raia Hadsell, M Ranzato, and F Huang. A tutorial on energy-based learning. *Predicting structured data*, 1(0), 2006.
- Radu Marinescu and Rina Dechter. AND/OR branch-and-bound for graphical models. In *Proceedings of the Nineteenth International Joint Conference on Artificial Intelligence (IJCAI)*, pages 224–229, Edinburgh, Scotland, UK, 2005.

- Alejandro Molina, Antonio Vergari, Karl Stelzner, Robert Peharz, Pranav Subramani, Nicola Di Mauro, Pascal Poupart, and Kristian Kersting. SPFlow: An easy and extensible library for deep probabilistic learning using sum-product networks. *arXiv preprint arXiv:1901.03704*, 2019.
- Erik Nijkamp, Mitch Hill, Song-Chun Zhu, and Ying Nian Wu. Learning non-convergent non-persistent short-run MCMC toward energy-based model. In *Advances in Neural Information Processing (NeurIPS)*, pages 5233–5243, Vancouver, Canada, 2019.
- Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake VanderPlas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Edouard Duchesnay. Scikit-learn: Machine learning in python. *J. Mach. Learn. Res.*, 12:2825–2830, 2011.
- Hoifung Poon and Pedro M. Domingos. Sum-product networks: A new deep architecture. In *Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 337–346, Barcelona, Spain, 2011.
- Shaogang Ren, Dingcheng Li, Zhixin Zhou, and Ping Li. Estimate the implicit likelihoods of gans with application to anomaly detection. In *Proceedings of the Web Conference (WWW)*, pages 2287–2297, Taipei, 2020.
- Shaogang Ren, Haiyan Yin, Mingming Sun, and Ping Li. Causal discovery with flow-based conditional density estimation. In *Proceedings of the IEEE International Conference on Data Mining (ICDM)*, pages 1300–1305, Auckland, New Zealand, 2021.
- Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *Proceedings of the 32nd International Conference on Machine Learning (ICML)*, pages 1530–1538, Lille, France, 2015.
- Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *Proceedings of the 31th International Conference on Machine Learning (ICML)*, pages 1278–1286, Beijing, China, 2014.
- Raquel Sánchez-Cauce, Iago París, and Francisco Javier Díez. Sum-product networks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence (early access)*, 2021.
- Marco Scutari. Learning bayesian networks with the bnlearn r package. *arXiv preprint arXiv:0908.3817*, 2009.
- Peter Sorrenson, Carsten Rother, and Ullrich Köthe. Disentanglement by nonlinear ICA with general incompressible-flow networks (GIN). In *Proceedings of the 8th International Conference on Learning Representations (ICLR)*, Addis Ababa, Ethiopia, 2020.
- Takeshi Teshima, Isao Ishikawa, Koichi Tojo, Kenta Oono, Masahiro Ikeda, and Masashi Sugiyama. Coupling-based invertible neural networks are universal diffeomorphism approximators. In *Advances in Neural Information Processing Systems (NeurIPS)*, virtual, 2020.
- Stef Van Buuren and Karin Groothuis-Oudshoorn. mice: Multivariate imputation by chained equations in R. *Journal of statistical software*, 45:1–67, 2011.
- Rianne van den Berg, Leonard Hasenclever, Jakub M. Tomczak, and Max Welling. Sylvester normalizing flows for variational inference. In *Proceedings of the Thirty-Fourth Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 393–402, Monterey, CA, 2018.

- Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *J. Mach. Learn. Res.*, 9:2579–2605, 2008.
- Martin J Wainwright and Michael Irwin Jordan. *Graphical models, exponential families, and variational inference*. Now Publishers Inc, 2008.
- Antoine Wehenkel and Gilles Louppe. Graphical normalizing flows. In *Proceedings of the 24th International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 37–45, Virtual Event, 2021.
- John M. Winn and Christopher M. Bishop. Variational message passing. *J. Mach. Learn. Res.*, 6: 661–694, 2005.
- Jianwen Xie, Yang Lu, Song-Chun Zhu, and Ying Nian Wu. A theory of generative convnet. In *Proceedings of the 33rd International Conference on Machine Learning (ICML)*, pages 2635–2644, New York City, NY, 2016.
- Eric P. Xing, Michael I. Jordan, and Stuart Russell. A generalized mean field algorithm for variational inference in exponential families. In *Proceedings of the 19th Conference in Uncertainty in Artificial Intelligence (UAI)*, pages 583–591, Acapulco, Mexico, 2003.
- Haiyan Yin, Dingcheng Li, Xu Li, and Ping Li. Meta-cotgan: A meta cooperative training paradigm for improving adversarial text generation. In *Proceedings of the Thirty-Fourth AAAI Conference on Artificial Intelligence (AAAI)*, pages 9466–9473, New York, NY, 2020.
- Yang Zhao, Jianwen Xie, and Ping Li. Learning energy-based generative models via coarse-to-fine expanding and sampling. In *Proceeding of the 9th International Conference on Learning Representations (ICLR)*, Virtual Event, 2021.
- Xun Zheng, Bryon Aragam, Pradeep Ravikumar, and Eric P. Xing. DAGs with NO TEARS: continuous optimization for structure learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 9492–9503, Montréal, Canada, 2018.
- Zilong Zheng, Jianwen Xie, and Ping Li. Patchwise generative convnet: Training energy-based models from a single natural image for internal learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2961–2970, virtual, 2021.
- Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Venice, Italy, 2017.
- Song-Chun Zhu, Ying Nian Wu, and David Mumford. Filters, random fields and maximum entropy (FRAME): towards a unified theory for texture modeling. *International Journal of Computer Vision (IJCV)*, 27(2):107–126, 1998.