

# Variational Flow Graphical Model

Anonymous Author(s)\*

## ABSTRACT

This paper introduces a novel approach to embed flow-based models with hierarchical structures. The proposed model learns the representation of high dimensional data via a message-passing scheme by carefully integrating flow-based functions through variational inference. Meanwhile, our model produces 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. With the proposed aggregation nodes, the model provides a new approach for distribution modeling and graphical inference on datasets. Multiple experiments on synthetic and real datasets show the benefits of our proposed method.

## CCS CONCEPTS

• **Computer systems organization** → **Embedded systems**; *Redundancy*; Robotics; • **Networks** → Network reliability;

## KEYWORDS

datasets, neural networks, gaze detection, text tagging

### ACM Reference Format:

Anonymous Author(s). 2018. Variational Flow Graphical Model. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, ?? pages. <https://doi.org/XXXXXXX.XXXXXXX>

## 1 INTRODUCTION

There are two general approaches for graphical inference: *exact inference* and *approximate inference*. Exact inference [?] resorts to an exact numerical calculation procedure of the quantity of interest. However, in most cases, exact inference is either *computationally involved* or simply *intractable*. Variational Inference (VI) is computationally efficient and has been applied to tackle the large scale inference problem [?]. In Variational Inference, mean-field approximation [?] and variational message passing [?] are two common approaches for graphical models. These approximation 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)*. These models include probabilistic decision graphs [?], arithmetic circuits [?], and-or search spaces [?], multi-valued

decision diagrams [?], cutset networks [?], sum-product nets [?], probabilistic sentential decision diagrams [?], and probabilistic circuits [?]. Probabilistic circuits (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.

Apart from probabilistic inference, generative models have been developed to model high dimensional datasets and to learn the 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) [?] and Generative Adversarial Networks (GAN) [?] are widely applied to different categories of datasets. Flow-based models [?] leverage invertible neural networks and can estimate the density values of data samples as well. 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 lie in the stream of tractable neural networks that allow to perform inference on graphical structures. Sum-product networks [?] and probabilistic circuits [?] 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. VFGs rely on approximation and consistency of nodes in graphical structures for tractable inference. Our contributions are summarized as follows.

**Contributions.** Dealing with high dimensional data using graphical models exacerbates this systemic inability to model the latent structure of the data efficiently, i.e., to sample from the posterior distribution of the latent variables given the observed data. To overcome these significant limitations, we propose a new framework, a variational hierarchical graphical flow model:

- **Hierarchical and Flow-Based:** Introducing the Variational Flow Graphical (VFG) model, we propose a novel graph 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.
- **Numerical Inference:** Aggregation nodes are introduced into our model in order to integrate hierarchical information through a variational forward-backward message passing scheme. We highlight the benefits of our VFG model on numerical applications: the missing values imputation problem and the numerical inference on graphical datasets. VFGs can achieve improved evidence lower bound (ELBO) and likelihood values due to the implicitly invertible flow-based model structure.
- **Representation Learning:** We specifically show in our experiments that our model achieves to disentangle the factors

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

Conference acronym 'XX, June 03–05, 2018, Woodstock, NY

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/XXXXXXX.XXXXXXX>

of variation underlying the high dimensional data given as input.

**Roadmap:** Section 2 presents important concepts used in the paper. Section 3 introduces the Variational Flow Graphical (VFG) model to tackle the posterior distribution learning of high dimensional data. Section ?? provides the algorithms used to train VFG models. Section ?? discusses how to perform inference with a trained VFG model. Section ?? showcases the advantages of VFG on various tasks. Section ?? and Section ?? provide a discussion and conclusion of the paper.

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

## 2 PRELIMINARIES

We introduce the general principles and notations of variational inference and flow-based models then explain how they can naturally be embedded with graphical models.

**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  $\theta$ , we maximize the following marginal log-likelihood  $\log p(\mathbf{x}) = \log \int p(\mathbf{z}) p_{\theta}(\mathbf{x}|\mathbf{z}) d\mathbf{z}$ . Direct optimization of the log-likelihood is usually not an option due to the intractable latent structure. Instead 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 goal of VI is to minimize the distance, in terms of Kullback-Leibler (KL) divergence, between the variational candidate and the true posterior  $\text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}|\mathbf{x}))$ . This optimization problem can be shown to be equivalent to maximizing the following evidence lower bound (ELBO) objective, noted  $\mathcal{L}(\mathbf{x}; \theta, \phi)$ :

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

In Variational Auto-Encoders (VAEs, [? ? ]), 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 [? ? ? ?] correspond to a probability distribution transformation using a sequence of invertible and differentiable mappings, noted  $\mathbf{f} : \mathcal{Z} \rightarrow \mathcal{X}$  between two random variables  $\mathbf{z} \in \mathcal{Z}$  of density  $p(\mathbf{z})$  and  $\mathbf{x} \in \mathcal{X}$ . The observed variable  $\mathbf{x} \sim p_{\theta}(\mathbf{x})$  is assumed to be distributed according to an unknown distribution  $p_{\theta}(\mathbf{x})$  parameterized by some parameter  $\theta$ . 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:

$$\log p_{\theta}(\mathbf{x}) = \log p(\mathbf{z}) + \log \left| \det \left( \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right) \right| \quad (3)$$

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

where we have  $\mathbf{h}^0 = \mathbf{x}$  and  $\mathbf{h}^L \stackrel{\text{def}}{=} \mathbf{z}$  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 (by sampling from the initial density and applying the transformations). [?] 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 negative EBLO is given by

$$\begin{aligned} \mathcal{F}(\theta, \phi) &= \mathbb{E}_{q_{\phi}} [\log q_{\phi}(\mathbf{z}|\mathbf{x}) - \log p_{\theta}(\mathbf{x}, \mathbf{z})] \\ &= \mathbb{E}_{q_0} [\log q_0(\mathbf{z}_0|\mathbf{x}) - \log p_{\theta}(\mathbf{x}, \mathbf{z})] \\ &\quad - \mathbb{E}_{q_0} \left[ \sum_{k=1}^K \log \left| \det \left( \frac{\partial \mathbf{f}_k(\mathbf{z}_k; \psi_k)}{\partial \mathbf{z}_k} \right) \right| \right]. \end{aligned} \quad (4)$$

Here  $\mathbf{f}_k$  is the  $k$ th flow with parameter  $\psi_k$ , i.e.  $\mathbf{z}_K = \mathbf{f}_K \circ \dots \circ \mathbf{f}_2 \circ \mathbf{f}_1(\mathbf{z}_0)$ . The parameters of the flows are considered as functions of data sample  $\mathbf{x}$ , and they determine the final distribution in amortized inference.

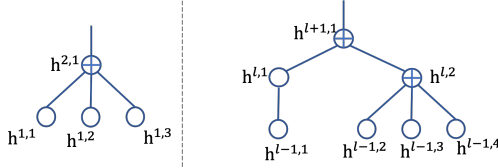
Several recent models have been proposed by leveraging the invertible flow-based models. Graphical normalizing flow [?] learns a DAG structure from the input data under sparse penalty and maximum likelihood estimation. The bivariate causal discovery method proposed in [?] relies on autoregressive structure of flow-based models and the log-likelihood ratio asymmetry for causal-effect inference. Here, we propose a framework that generalizes flow-based models [? ? ? ?] to graphical variable inference.

## 3 VARIATIONAL FLOW GRAPHICAL MODEL

Assume that there exist a relation between each data sample components and their corresponding variable in the latent space. 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 the model. We call this model a *Variational Flow Graphical Model* (VFG) and introduce it in the sequel.

### 3.1 Evidence Lower Bound of Variational Flow Graphical Models

A VFG model  $\mathbb{G} = \{\mathcal{V}, \mathbf{f}\}$  consists of a set of nodes ( $\mathcal{V}$ ) and a set of edges ( $\mathbf{f}$ ). Figure ??-Right gives an illustration of a tree structure induced by a VFG model. 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 other nodes with a single flow function or an identity function. An aggregation node has multiple children, and it connects with each of them with an identity function. Unlike classical graphical models, a node in a VFG model may represent multiple variables. A node only belongs to one layer, and one layer can have



**Figure 1: (Left) Node  $h^{2,1}$  connects its children with invertible functions. Messages from the children are aggregated at the parent node,  $h^{2,1}$ ;  $\oplus$  is an aggregation node, and circles stand for non-aggregation nodes. (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.**

multiple nodes. Moreover, each latent variable belongs to only one node in a VFG.

In the following sections of this paper, identity function is considered as a special case of flow functions. We apply variational inference to assemble the layers of a VFG and learn parameters 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.

The hierarchical generative network comprises  $L$  layers,  $h^l$  denotes the latent variable in layer  $l$ , and  $\theta$  is the vector of model parameters. We use  $h^{l,i}$  to denote the  $i$ th node's latent variable in layer  $l$ , and  $h^{(j)}$  to represent node  $j$ 's latent variable without specification of the layer number, and  $j$  is the node index on a tree or graph. The joint distribution of the hierarchical model is given by:

where  $\mathbf{h} = \{h^1, \dots, h^L\}$  denotes the set of latent variables of the model. The hierarchical generative model is given by factorization  $p(\mathbf{x}|\mathbf{h}^L) = \prod_{l=1}^{L-1} p(h^l|h^{l+1})p(\mathbf{x}|h^1)$ , and the prior distribution is  $p(h^L)$ . Note that only the root nodes have prior distributions.

The probability density function  $p(h^{l-1}|h^l)$  in the generative model is parameterized with one or multiple invertible flow-based functions. We follow the variational inference to approximate the posterior distribution of latent variables. The hierarchical posterior (recognition model) is factorized as

Evaluation of the posterior (recognition model) (??) 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 given by

$$\log p_{\theta}(\mathbf{x}) \geq \mathcal{L}(\mathbf{x}; \theta) \quad (6)$$

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

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

$$\text{KL}^l = \mathbb{E}_{q(h^{1:L}|\mathbf{x})} [\log q(h^l|h^{l-1}) - \log p(h^l|h^{l+1})]. \quad (7)$$

When  $l = L$ ,  $\text{KL}^L = \mathbb{E}_{q(h^{1:L}|\mathbf{x})} [\log q(h^L|h^{L-1}) - \log p(h^L)]$ . It is easy to extend the computation of the ELBO (??) 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:

$$\mathcal{L}(\mathbf{x}; \theta) = \mathbb{E}_{q(\mathbf{h}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{h})] - \sum_{i \in \mathcal{V} \setminus \mathcal{R}_{\mathbb{G}}} \text{KL}^{(i)} \quad (8)$$

$$- \sum_{i \in \mathcal{R}_{\mathbb{G}}} \text{KL}(q(h^{(i)}|h^{ch(i)})||p(h^{(i)})).$$

Here  $\text{KL}^{(i)} = \mathbb{E}_{q(\mathbf{h}|\mathbf{x})} [\log q(h^{(i)}|h^{ch(i)}) - \log p(h^{(i)}|h^{pa(i)})]$ .  $\mathcal{R}_{\mathbb{G}}$  is the set of root or source 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  $\mathbf{x} = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)}]$ , then the hidden variables in both (??) and (??) are computed with forward and backward message passing.



**Figure 2: Recognition model consists of forward message from data to approximate the posterior distributions; the generative model is realized by backward message from the root node and generates the samples or reconstructions in each layer.**

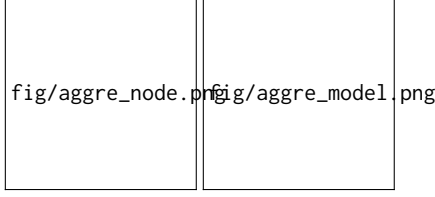
### 3.2 Aggregation Nodes

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 ??-Left 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, we observe that at node  $i$

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

$$\widehat{\mathbf{h}}^{(i)} = \frac{1}{|pa(i)|} \sum_{j \in pa(i)} \mathbf{f}_{(i,j)}^{-1}(\widehat{\mathbf{h}}^{(j)}).$$

Notice that the above two equations hold even when node  $i$  has only one child or parent. In the sequel, we consider that all latent variables, noted  $h^{l,i}$ , for all  $l \in [L]$  and  $i \in \mathbb{N}$ , are distributed according to Laplace distributions. With the identity function between the parent and its children, there are *node consistency rules* regarding an *average aggregation node*: (a) the parent value is the mean of its children, i.e.,  $\mathbf{h}^{(i)} = \frac{1}{|ch(i)|} \sum_{j \in ch(i)} \mathbf{h}^{(j)}$ ; (b) the child node have the same reconstruction value with its parent, i.e.,  $\widehat{\mathbf{h}}^{(j)} = \widehat{\mathbf{h}}^{(i)}, \forall j \in ch(i)$ .



**Figure 3: (Left) Aggregation node  $h^{l+1,1}$  has three children,  $h^{l,1}$ ,  $h^{l,2}$ , and  $h^{l,3}$ . (Right) A VFG model with one aggregation node,  $h^{(r)}$ . Solid circles are nodes with observed values, and the diamond is the prior for the root node.**

We use Figure ??-Right as an example to illustrate a simple model that has only one average aggregation node. Then (??) yields the ELBO,

$$\begin{aligned} \log p(\mathbf{x}) &\geq \mathcal{L}(\mathbf{x}; \theta) \\ &= \mathbb{E}_{q(\mathbf{h}|\mathbf{x})} [\log p(\mathbf{x}|\hat{\mathbf{h}}^1)] - \mathbb{E}_{q(\mathbf{h}|\mathbf{x})} [\log q(\mathbf{h}^1|\mathbf{x}) \\ &\quad - \log p(\mathbf{h}^1|\hat{\mathbf{h}}^2)] - \text{KL}(q(\mathbf{h}^2|\mathbf{h}^1)|p(\mathbf{h}^2)). \end{aligned} \quad (9)$$

From Figure ??-Right,  $h^{(r)}$  is the root, and has  $k$  children,  $h^{(t)}$ ,  $t = 1, \dots, k$ , and  $k = 3$ . With  $\mathbf{f}_t$  as the flow function connecting  $h^{(t)}$  and  $\mathbf{x}^{(t)}$ , according to the aggregation rules, we get:

$$\begin{aligned} \mathbf{h}^{(t)} &= \mathbf{f}_t(\mathbf{x}^{(t)}), \quad \hat{\mathbf{h}}^{(r)} = \mathbf{h}^{(r)} = \frac{1}{k} \sum_{t=1}^k \mathbf{h}^{(t)}, \\ \hat{\mathbf{h}}^{(t)} &= \hat{\mathbf{h}}^{(r)}, \quad t = 1, \dots, k. \end{aligned} \quad (10)$$

Assume the data to be normally distributed, then

$$\begin{aligned} \log p(\mathbf{x}|\hat{\mathbf{h}}^1) &= - \sum_{t=1}^k \left\{ \frac{1}{2\sigma_{\mathbf{x}}^2} \|\mathbf{x}^{(t)} - \mathbf{f}_t^{-1}(\hat{\mathbf{h}}^{(r)})\|_2^2 \right\} + C \\ &\quad \text{By } \hat{\mathbf{x}}^{(t)} = \mathbf{f}_t^{-1}(\hat{\mathbf{h}}^{(t)}) = \mathbf{f}_t^{-1}(\hat{\mathbf{h}}^{(r)}) \\ \log p(\mathbf{h}^1|\hat{\mathbf{h}}^2) &= - \sum_{t=1}^k \left\{ \|\mathbf{f}_t(\mathbf{x}^{(t)}) - \hat{\mathbf{h}}^{(r)}\|_1 \right\} + C. \end{aligned} \quad (11)$$

By  $\hat{\mathbf{h}}^2 = \hat{\mathbf{h}}^{(r)}$ ,  $\mathbf{h}^{(t)} = \mathbf{f}_t(\mathbf{x}^{(t)})$

Here  $\sigma_{\mathbf{x}}$  is a constant standard deviation. We notice that maximizing the ELBO, or minimizing the KL term of an aggregation node implies the aggregation rule (b).

### 3.3 Universal Approximation Property

Coupling layer based flows have universal approximation power [? ]. Following the analysis for flows [? ], we can prove that coupling layer based VFGs have universal approximation as well.

## 4 ALGORITHM AND IMPLEMENTATION

In this section, we develop the training algorithm, see Algorithm ??, that outputs the fitted vector of parameters resulting from the maximization of the ELBO objective function (??) or (??) depending on what graph structure is used. In Algorithm ??, the inference of the latent variables is performed via forwarding message passing, cf. Line 6, and their reconstructions are computed in backward message passing, cf. Line 11. Unlike Variational Autoencoders (VAE), the variance of latent variables in a VFG is approximated rather

than parameterized with neural networks. A VFG is a deterministic network passing latent variable values between nodes. The reconstruction (likelihood) terms in each layer are computed with forward and backward node states. Ignoring explicit neural network parameterized variances for all latent nodes enables us to use flow-based models as both the encoders and decoders. We thus obtain a deterministic ELBO objective (??)- (??) 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
```

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 (??) or (??) with the above algorithm, the *node consistency rules* in Section ?? are expected to be satisfied. Inference on sub-graphs can be further improved with random masking strategy, and more details are given in the supplemental file.

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  $q(\mathbf{h}|\mathbf{x})$ .

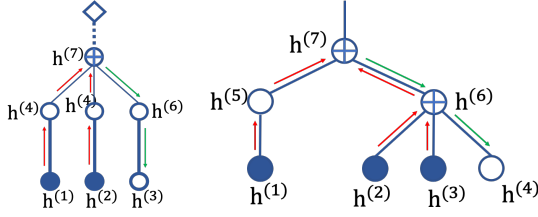
## 5 INFERENCE ON VFG MODELS

Given a VFG model, we aim to infer the state of any node given observed ones. Relations between variables at different nodes can also be inferred via our flow-based graphical model.

The hidden state of the parent node  $j$  in an aggregation model can be computed by the observed children as follows:

$$\mathbf{h}^{(j)} = \frac{1}{|ch(j) \cap O|} \sum_{i \in ch(j) \cap O} \mathbf{h}^{(i)}, \quad (12)$$





**Figure 4: (Left) Inference on model with single aggregation node. Node 7 aggregates information from node 1 and 2, and pass down the updated state to node 3 for prediction. (Right) Inference on a tree model. 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.**

where  $O$  is the set of observed leaf nodes, see Figure ??-left for an illustration. Observe that for either a tree or a DAG, the state of any given node is updated via messages received from its children. Figure ?? 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 5.1.** *Let  $\mathbb{G}$  be a variational flow graphical model (with a tree structure) with  $L$  layers, and  $i$  and  $j$  are two leaf nodes with  $a$  as the closest common ancestor. Given observed value at node  $i$ , the value of node  $j$  can be approximated by  $\hat{x}^j = f_{(a,j)}(f_{(i,a)}(x^{(i)}))$ . Here  $f_{(i,a)}$  is the flow function path from node  $i$  to node  $a$ .*

**PROOF.** Without loss generality, we assume that there are relationships among different data sections, and the value of one section can be imputed by other sections. According to the *node consistency rule* (b) discussed in section ??, at an aggregation node  $a$ , the latent value of a child node  $j$  has the same reconstruction value as the parent node. The reconstruction of the child node  $j$  can be calculated with the reconstruction of the parent node, i.e.,  $\hat{h}^{(j)} = f_{(a,j)}(\hat{h}^{(a)})$ . Recalling the reconstruction term in the ELBO (??), at each node we have  $h^{(a)} = \hat{h}^{(a)}$ . Hence for node  $a$ 's descendent  $j$ , we have  $\hat{h}^{(j)} = f_{(a,j)}(h^{(a)})$ , and  $f_{(a,j)}$  is the flow function path from  $a$  to  $j$ . The value of node  $a$  can be computed by the value of its descendent  $i$ , i.e.,  $h^{(a)} = f_{(i,a)}(h^{(i)})$ . Hence, we have  $\hat{x}^{(j)} = f_{(a,j)}(f_{(i,a)}(x^{(i)}))$ .  $\square$

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

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

Let  $O$  be the set of observed leaf nodes,  $j$  be an unobserved node, and  $a$  the closest ancestor of  $O \cup \{a\}$ . Then the state of  $j$  can be imputed with  $\hat{x}^j = f_{(a,j)}(f_{(O,a)}(x^{(i)}))$ , where  $f_{(O,a)}$  is the flow function path from all nodes in  $O$  to  $a$ . Lemma ?? provides an approach to conduct inference on a tree and impute missing values in the dataset. It is easy to extend the inference method to DAG-VFG models.

## 6 NUMERICAL EXPERIMENTS

The first application we present is the imputation of missing values. We compare our method with several baseline models on a synthetic

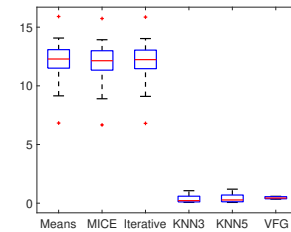
dataset. The second set of experiments is to evaluate VFG models on three different datasets, i.e., MNIST, Caltech101, and Omniglot, with ELBO and likelihoods as the score. The third application we present is the task of learning posterior distribution of the latent variables corresponding to the hidden explanatory factors of variations in the data [? ]. 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 experiments, the VFG structures are 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 [? ].

### 6.1 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. The baselines for this set of experiments include mean value method (Means), iterative imputation (Iterative) [? ], and multivariate imputation by chained equation (MICE) [? ]. Mean Squared Error as the metric of reference in order to compare the different methods for the imputation task.

**Synthetic dataset:** In this set of experiments, we study the proposed model with synthetic datasets. 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  $x$ , 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 [? ]. Each child takes 2 variables as input data section, and the latent dimension of the VFG is 2. We compare, Figure ??, 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: Synthetic datasets: MSE boxplots of VFG and baseline methods.**

**Table 1: Negative log-likelihood and free energy (negative evidence lower bound) for static MNIST, Caltech101, and Omniglot.**

Model	MNIST		Caltech101		Omniglot	
	-ELBO	NLL	-ELBO	NLL	-ELBO	NLL
VAE [? ]	86.55 ± 0.06	82.14 ± 0.07	110.80 ± 0.46	99.62 ± 0.74	104.28 ± 0.39	97.25 ± 0.23
Planer [? ]	86.06 ± 0.31	81.91 ± 0.22	109.66 ± 0.42	98.53 ± 0.68	102.65 ± 0.42	96.04 ± 0.28
IAF [? ]	84.20 ± 0.17	80.79 ± 0.12	111.58 ± 0.38	99.92 ± 0.30	102.41 ± 0.04	96.08 ± 0.16
SNF [? ]	83.32 ± 0.06	80.22 ± 0.03	104.62 ± 0.29	93.82 ± 0.62	99.00 ± 0.04	93.77 ± 0.03
VFG	<b>80.80 ± 0.76</b>	<b>63.66 ± 0.14</b>	<b>67.26 ± 0.53</b>	<b>65.74 ± 0.84</b>	<b>80.16 ± 0.73</b>	<b>78.65 ± 0.66</b>

## 6.2 ELBO and Likelihood

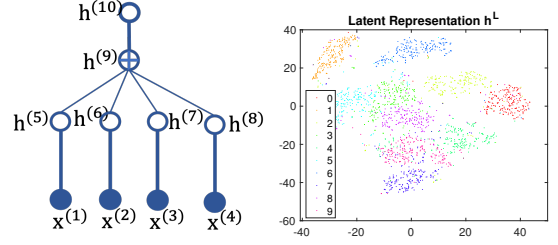
We further qualitatively compared VFG with existing methods on variational inference using three standard datasets. The evaluation datasets and setup are exactly following two standard flow-based variational models, sylvester normalizing flows [? ] and [? ]. We use the tree VFG structure as shown in Figure ??-Left for three datasets. We train the tree-VFG model using the following ELBO objective. Empirically, a small  $\beta$  yields better ELBO and NLL values, and we set  $\beta$  around 0.1 in the experiments.

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

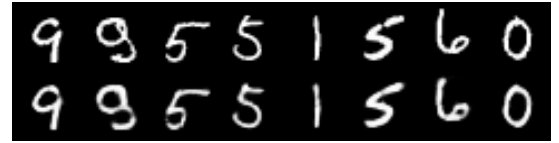
In Table ??, we provide the negative evidence lower bound (-ELBO) and the estimated negative likelihood (NLL) for baseline methods on three datasets, MNIST, Caltech101, and Omniglot. The results of the baseline methods are from [? ]. These methods are VAE based methods enhanced with normalizing flows. 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[? ] number with  $\mathcal{B} = 4$ , and following [? ] we run multiple times to get the mean and standard derivation as well. VFG can achieve superior EBLO and NLL values as given in Table ?. The main reason why VFGs can achieve superior variational inference results (ELBOs and NLLs) in Table ?? is due to VFGs's approximately invertible flow-based model structure. The intrinsic invertibility ensures the decoder or generative model in the VFG to achieve smaller reconstruction errors for data samples and hence smaller NLL values.

## 6.3 Latent Representation Learning on MNIST

In this set of experiments, we evaluate Variational Flow Graphical Models on latent representation learning of the MNIST dataset [? ]. We construct a tree VFG model depicted in Figure ??-Left. In the first layer, there are 4 flow functions, and each of them takes  $14 \times 14$  image blocks as the input. Thus a  $28 \times 28$  input image is divided into four  $14 \times 14$  blocks as the input of VFG model. The latent dimension for this model is 196. Following [? ], 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  $\lambda^L$ 's regarding the 10 digits. The images in the second row of Figure ?? 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 6: (Left) Tree structure for MNIST; (Right) MNIST: t-SNE plot of latent variables from VFG learned with labels.**



**Figure 7: (Top) original MNIST digits. (Bottom) reconstructed images using VFG.**

Figure ??-Right shows t-distributed stochastic neighbor embedding (t-SNE) [? ] plot of 2,000 testing images's latent variables learned with our model, and 200 for each digit. Figure ??-Right illustrates that VFG can learn separated latent representations to distinguish different hand-written numbers.

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 ?? 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  $d = 196$  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.

## 7 DISCUSSION

One of the motivations for VFG is to develop a tractable model that can be used for inference and sampling on datasets. As long as the variable states in the aggregation nodes are consistent, we can always apply VFGs to missing value inference. We provide more discussion on the structures of VFGs below.



Figure 8: MNIST: Increasing each latent variable from a small value to a larger one.

## 7.1 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 achieve 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.

## 7.2 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 can be used 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. 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 [? ? ].

VFGs rely on minimizing the KL term to achieve *consistency* among the nodes in an aggregation node. As long as the aggregation nodes retain consistency, the model always has a tight ELBO and can be applied for tractable posterior inference. According to the recent theoretical study [? ], 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 training algorithm and thus a tight ELBO.

## 8 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 the latent representation of the input data through message passing between nodes in the model structure. The posterior inference, of the latent nodes given input observations, is facilitated by the careful embedding of flow functions in the general graph structure. Experiments on different datasets illustrate the effectiveness of the model. Future work includes applying our VFG model to fine grained data relational structure learning and reasoning.

The proposed variational flow graphical models assemble flow functions with tree or DAG structures via variational inference on aggregation nodes. In this supplemental file, we first present additional results, then we give more details on the methodology of VFG models. Section ?? provides additional experiments; universal approximation property of VFGs is given in Section ??; Section ?? adds details on ELBO computation; Section ?? provides more details on the training algorithm; Section ?? presents other details on the aggregation node; ELBO derivation can be found in Section ??.

## A ADDITIONAL RESULTS

All the experiments are conducted on NVIDIA-TITAN X (Pascal) GPUs. In the experiments, we use the same coupling block [?] 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  $\mathcal{B} \geq 2$ . All latent variables,  $\mathbf{h}^i, i \in \mathcal{V}$  are forced to be non-negative via Sigmoid or RELU functions. Non-negativeness can help the model to identify sparse structures of the latent space.

### A.1 Imputation on 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. We use the baseline implementations in [?] or [?] in the experiments.

Here, 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 4 coupling blocks. We can see Table ?? that our model yields significantly better results than any other method in terms of prediction error.

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

**Table 2: California Housing dataset: Imputation Mean Squared Error (MSE) results.**

## B UNIVERSAL APPROXIMATION

In [?], the authors analysis the approximation power of coupling-based flows. Under mild conditions, coupling-based flows are universal distributional approximators. Following [?], we show that VFGs are also universal approximators. We first give several definitions regarding universal approximation, then we prove that VFGs have universal approximation.

For a measurable mapping  $\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^n$  and a subset  $K \subset \mathbb{R}^m$ , we have the following definition,

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

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

**Definition B.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  $\|\cdot\|_{\text{sup},K}$ .

**Definition B.2.** (Distributional universality) Let  $\mathcal{M}$  be a model which is a set of measurable mappings from  $\mathbb{R}^m$  to  $\mathbb{R}^n$ . We say that a model  $\mathcal{M}$  is a distributional universal approximator or has the distributional universal approximation property if, for any absolutely continuous probability measure  $\mu$  on  $\mathbb{R}^m$  and any probability measure  $\nu$  on  $\mathbb{R}^n$ , there exists a sequence  $\{\mathbf{f}_i\}_{i=1}^{\infty} \subset \mathcal{M}$  such that  $(\mathbf{f}_i)_* \mu$  converges to  $\nu$  in distribution as  $i \rightarrow \infty$ , where  $(\mathbf{f}_i)_* \mu := \mu \circ \mathbf{f}_i^{-1}$ .

**Definition B.3.** (Immersion and submanifold)  $\mathbf{g} : \mathcal{M} \rightarrow \mathcal{N}$  is said to be an immersion if  $\text{rank}(\mathbf{g}) = m = \dim(\mathcal{M})$  everywhere. If  $\mathbf{g}$  is injective (one-to-one) immersion, then  $\mathbf{g}$  establish an one-to-one correspondence of  $\mathcal{M}$  and the subset  $\mathcal{N} = \mathbf{g}(\mathcal{M})$  of  $\mathcal{N}$ . If we use this correspondence to endow  $\mathcal{M}$  with a topology and  $C^\infty$  structure, then  $\mathcal{M}$  will be called a submanifold (or immersed submanifold) and  $\mathbf{g} : \mathcal{M} \rightarrow \mathcal{N}$  is a diffeomorphism.

**Definition B.4.** ( $C^r$ -diffeomorphisms for Submanifold:  $Q^r$ ). We define  $Q^r$  as the set of all  $C^r$ -diffeomorphisms  $\mathbf{g} : U_{\mathbf{g}} \rightarrow \mathcal{U}$ , where  $U_{\mathbf{g}} \subset \mathbb{R}^m$  is an open set  $C^r$ -diffeomorphic to  $\mathcal{U}$ , which may depend on  $\mathbf{g}$ , and  $\mathcal{U}$  is a submanifold of  $\mathbb{R}^n$ .

We use  $m$  to represent the root dimension of a VFG, and  $n$  to denote the dimension of data samples. VFGs learn the data manifold embedding the  $\mathbb{R}^n$ . We define  $C_c^\infty(\mathbb{R}^{m-1})$  as the set of all compactly-supported  $C^\infty$  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 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 B.5.** ( $L^p$ -universality) Let  $p \in [0, \infty)$ . Assume  $\mathcal{H}$  is a sup-universal approximator for  $C_c^\infty(\mathbb{R}^{m-1})$ , and that it consists of  $C^1$ -functions. Then  $\text{VFG}_{\mathcal{H}\text{-ACF}}$  is an  $L^p$ -universal approximator for  $Q_c^0$ .

**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 an one-layer VFG tree model  $\mathbf{f}$  with the root as an aggregation node. The children divide the  $n$  input entries into  $\tau = n/m$  sections, and each section connects the aggregation node with a flow function.

Given an injective immersion  $\mathbf{g} : \mathcal{M} \rightarrow \mathcal{N}$ , function  $\mathbf{g}$  can be represented with the concatenation of a set of functions, i.e.  $\mathbf{g} = [\mathbf{g}_1, \dots, \mathbf{g}_\tau]^\top$ . According to the function decomposition theory [?],



the inverse can be represent as the summation of functions  $\mathbf{g}_i^{-1}$ ,  $1 \leq i \leq \tau$ , i.e.  $\mathbf{g}^{-1} = \frac{1}{\tau} \sum_{i=1}^{\tau} \mathbf{g}_i^{-1}$ . Each  $\mathbf{g}_i \in \mathcal{D}^2$ , and  $\mathfrak{M}_i = \mathbf{g}_i(\mathfrak{M})$  is diffeomorphic to  $\mathfrak{M}$ . According to Theorem 2 in [?],  $\mathcal{H} - \text{ACF}$  is an universal approximator for each  $\mathbf{g}_i$ . 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  $\tau$ th section and the remaining  $r$  entries into two equal small sections that are denoted with  $\tau$  and  $\tau + 1$ . Sections  $\tau$  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  $\mathbf{g}_\tau$  in  $\mathbb{R}^m$  is  $I_\tau = [1, 2, \dots, \lceil (m+r)/2 \rceil]$ , and the out put index of  $\mathbf{g}_\tau$  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 out put index of  $\mathbf{g}_{\tau+1}$  in  $\mathbb{R}^n$  is  $I_{\tau+1} + \gamma$ . We can see that the  $m$  dimension is 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  $\mathbf{g} : \mathfrak{M} \rightarrow \mathfrak{N}$  is 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  $\mathbf{g}_\tau^{-1}$  or  $\mathbf{g}_{\tau+1}^{-1}$  is omitted due to the missing of entry  $j$  in the function output.

The mapping  $\mathbf{g}_\tau$  is a diffeomorphism from manifold  $\mathfrak{M}_\tau$  ( $\mathfrak{M}_\tau \subset \mathfrak{M}$ ) to sub-manifold  $\mathfrak{M}_\tau$  of  $\mathfrak{N}$ . Similarly  $\mathbf{g}_{\tau+1}$  is a diffeomorphism from  $\mathfrak{M}_{\tau+1}$  to manifold  $\mathfrak{M}_{\tau+1}$ . For each  $\mathbf{g}_i$ ,  $i : 1 \leq i \leq \tau + 1$ , it can universally approximated with a function in  $\mathcal{H} - \text{ACF}$ . Here we can construct a VFG with universal approximation for  $\mathbf{g}$ .

□

## C ADDITIONAL DETAILS ON ELBO

This section presents more details on ELBO calculation. The recognition model in a VFG is the neural network (encoder) used to approximate the posterior of latent variables. With invertible neural networks (flows), the recognition model and the generative model in a VFG share the same structure and parameters. As shown in Figure ??, the recognition and generative models are realized with forward and backward message passing, respectively.

Maximize the ELBOs (??,??) requires evaluation of both the reconstruction and the KL terms. It involves node state samples from the posterior in both forward and backward messages.

### C.1 ELBO Calculation

Maximizing the ELBO (??) equals to optimizing the parameters of the flows,  $\theta$ . 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 ???. VFGs are hierarchical aggregation model that requires information from all the leaf nodes aggregates at the root nodes.

For the following sections, we use  $\mathbf{h}^l$  to represent the node state of layer  $l$  in the forward message, and  $\hat{\mathbf{h}}^l$  for the node state in the backward message. For all layer, both  $\mathbf{h}^l$  and  $\hat{\mathbf{h}}^l$  are sampled from the posterior. At the rood nodes, we have  $\hat{\mathbf{h}}^R = \mathbf{h}^R$ . The calculation of the data reconstruction term in (??) requires samples of  $\mathbf{h}^l$  and

$\hat{\mathbf{h}}^l$  from the posterior. It corresponds to the encoding and decoding procedures in VAE model [?] as given by Eq. (2). Specifically, we have

$$\mathbf{h}^l \sim q(\cdot | \mathbf{h}^{l-1}) \quad \text{where} \quad 1 \leq l \leq L, \quad (13)$$

$$\hat{\mathbf{h}}^l \sim p(\cdot | \hat{\mathbf{h}}^{l+1}) \quad \text{where} \quad 0 \leq l \leq L-1. \quad (14)$$

At the root node, we have  $\hat{\mathbf{h}}^L = \mathbf{h}^L \sim q(\cdot | \mathbf{h}^{L-1})$ . The computation of  $\mathbf{h}^l$ s and  $\hat{\mathbf{h}}^l$ s with (??) and (??) requires message passing via the flow functions shown in Figure ?. Latent node states in both (??) and (??) are directly obtained with flow functions.

The reconstruction term in ELBO (??) can be computed with the backward message from the generative model  $p(\mathbf{x} | \hat{\mathbf{h}}^1)$ , i.e.,

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

For a VFG model, we set  $U = 1$ . In the last term,  $p(\mathbf{x} | \hat{\mathbf{h}}^1)$  is either Gaussian or Binary distribution parameterized with  $\hat{\mathbf{x}}$  generated via the flow function with  $\hat{\mathbf{h}}^1$  as the input. For any  $l \in [L]$ , the calculation of the  $\text{KL}^l$  term is done in a similar manner, and (??) requires both forward message from the posterior and backward message from the generative model, i.e.,

$$\begin{aligned} \text{KL}^l &= \mathbb{E}_{q(\mathbf{h}^{1:L} | \mathbf{x})} [\log q(\mathbf{h}^l | \mathbf{h}^{l-1}) - \log p(\mathbf{h}^l | \hat{\mathbf{h}}^{l+1})] \\ &\simeq \log q(\mathbf{h}^l | \mathbf{h}^{l-1}) - \log p(\mathbf{h}^l | \hat{\mathbf{h}}^{l+1}). \end{aligned} \quad (15)$$

In a tree structured VFG, each layer may consist of multiple nodes. In layer  $l$ , each node's KL value can be computed individually, thus  $\text{KL}^l = \sum_{i \in l} \text{KL}^{(i)}$ , and here each  $i$  represents a node. For a DAG structure, we can individually evaluate each node's KL term.

### C.2 Distributions of Latent Variables

We first discuss additional details on the distributions of latent variables.

**C.2.1 Generative Model.** In a tree VFG, the sample reconstruction in the generative model consists of layer-wise backward message passing, i.e., latent variable generation in each layer. For any  $l$ ,  $0 \leq l \leq L-1$ , latent variable backward state (reconstruction)  $\hat{\mathbf{h}}^l$  is propagated from layer  $l+1$  via the flow function  $\mathbf{f}_l$  between the two layers with  $\hat{\mathbf{h}}^l = \mathbf{f}_l^{-1}(\hat{\mathbf{h}}^{l+1})$ . The prior  $p(\mathbf{h}^L)$  for the root latent variable  $\mathbf{h}^L$  is Laplace(0,1). With a sample  $\hat{\mathbf{h}}^L$  from the posterior, i.e.,  $\hat{\mathbf{h}}^L = \mathbf{h}^L \sim q(\cdot | \mathbf{h}^{L-1})$ , the conditional distribution for latent variable in layer  $l$  is  $p(\cdot | \hat{\mathbf{h}}^{l+1}) := \text{Laplace}(\hat{\mathbf{h}}^l, 1)$ . Here the location parameter is generated from layer  $l+1$ , i.e.,  $\hat{\mathbf{h}}^l = \mathbf{f}_l^{-1}(\hat{\mathbf{h}}^{l+1})$ . For a latent variable  $\mathbf{h}^l$  sampling from the posterior, its log-likelihood regarding  $p(\cdot | \hat{\mathbf{h}}^{l+1})$  in (??) is given by

$$\log p(\mathbf{h}^l | \hat{\mathbf{h}}^{l+1}) = -\|\mathbf{h}^l - \hat{\mathbf{h}}^l\|_1 - d \cdot \log 2.$$

Here  $d = \dim(\mathbf{h}^l)$ . Hence, minimizing KLs is to minimize the  $\ell_1$  distance between latent variables and their reconstructions.

**C.2.2 Recognition Model.** The forward message passing in the recognition model consists of layer-wise sample generation. In layer  $l$ ,  $1 \leq l \leq L$ , latent variable forward state  $\mathbf{h}^l$  is propagated from

layer  $l - 1$  via the flow function  $f_{l-1}$  between the two layers with  $\mathbf{h}^l = f_{l-1}(\mathbf{h}^{l-1})$ .

We assume each entry of hidden variable  $\mathbf{h}^l$  follows a Laplace distribution, i.e.,  $\mathbf{h}_j^l \sim \text{Laplace}(\mu_j^l, s_j^l)$  for layer  $l$ 's  $j$ th entry. Here  $\mu_j^l$  is the location and  $s_j^l$  is the scale. Compared with other distributions, Laplace can introduce sparsity to the model and it works well in practice. At level  $l \in [L]$ , we set  $q(\cdot|\mathbf{h}^{l-1}) := \text{Laplace}(\mu^l, \mathbf{s}^l)$  with

$$\mu^l = \text{median}(H), \quad \mathbf{s}^l = \frac{1}{B} \sum_{b=1}^B |\mathbf{h}^l(\mathbf{x}_b) - \mu^l|. \quad (16)$$

Here  $H = \{\mathbf{h}^l(\mathbf{x}_b) | 1 \leq b \leq B\}$  is a batch of latent values generated from a batch of data samples with size  $B$ , i.e.,  $X_B = \{\mathbf{x}_b | 1 \leq b \leq B\}$ . The median operation is performed element-wisely. For each  $\mathbf{x}_b$ ,  $\mathbf{h}^l(\mathbf{x}_b) = f^{l-1}(\mathbf{h}^{l-1}(\mathbf{x}_b))$ .

In summary, we use Laplace distribution to model latent distributions.  $q(\cdot|\mathbf{h}^{l-1})$  is a Laplace with location and scale equal to the median and scale of a batch of  $\mathbf{h}^l$ , respectively;  $p(\cdot|\widehat{\mathbf{h}}^{l+1})$  is a Laplace parameterized with  $(\widehat{\mathbf{h}}^l, 1.0)$  as the location and scale parameters. Hence with  $\mathbf{h}^l$  we can compute the log-likelihoods on RHS of (??) and thus the  $\text{KL}^l$  value. We use Laplace distribution for latent variables because it may introduce sparsity in latent variables. We also tried Gaussian and Gamma distributions, and they also work well in most tasks. VFGs leverage the universal approximation power [?] of coupling flows, and tractable inference could be achieved with different latent distributions.

### C.3 KL Term

For any  $l$ ,  $1 \leq l \leq L - 1$ , the calculation of the  $\text{KL}^l$  term (??) requires message passing and samples from both recognition and generative models, i.e.,

$$\begin{aligned} \text{KL}^l &= \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log q(\mathbf{h}^l|\mathbf{h}^{l-1}) - \log p(\mathbf{h}^l|\widehat{\mathbf{h}}^{l+1})] \\ &\simeq \log q(\mathbf{h}^l|\mathbf{h}^{l-1}) - \log p(\mathbf{h}^l|\widehat{\mathbf{h}}^{l+1}). \end{aligned} \quad (17)$$

Here  $q(\cdot|\mathbf{h}^{l-1})$  is a Laplace with location and scale equal to the median and scale defined in (??);  $p(\cdot|\widehat{\mathbf{h}}^{l+1})$  is a Laplace parameterized with  $(\widehat{\mathbf{h}}^l, 1.0)$  as discussed in ?? . Hence with  $\mathbf{h}^l$  we can compute the log-likelihoods on RHS of (??) and thus the  $\text{KL}^l$  value. When  $l = L$ ,  $\text{KL}^L = \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log q(\mathbf{h}^L|\mathbf{h}^{L-1}) - \log p(\mathbf{h}^L)] \simeq \log q(\mathbf{h}^L|\mathbf{h}^{L-1}) - \log p(\mathbf{h}^L)$ .

Assume  $k$  leaf nodes on a tree or a DAG model, corresponding to  $k$  sections of input sample  $\mathbf{x} = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)}]$ . Then the hidden variables in both (??) and (??) are computed with forward and backward message passing.

In practice, we set  $U = 1$  for efficiency. With a batch of training samples, the structure of flow functions make the forward and backward message passing very efficient, and thus the estimation of the ELBO.

### C.4 Reconstruction Term

The reconstruction term in ELBO (??) can be computed with the backward message from the generative model  $p(\mathbf{x}|\widehat{\mathbf{h}}^1)$ , i.e.,

$$\begin{aligned} \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{h}^{1:L})] &= \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} [\log p(\mathbf{x}|\widehat{\mathbf{h}}^{1:L})] \\ &\simeq \frac{1}{U} \sum_{u=1}^U \log p(\mathbf{x}|\widehat{\mathbf{h}}_u^{1:L}) = \frac{1}{U} \sum_{u=1}^U \log p(\mathbf{x}|\widehat{\mathbf{h}}_m^1) \\ &\simeq \log p(\mathbf{x}|\widehat{\mathbf{h}}^1). \end{aligned}$$

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

## D ADDITIONAL DETAILS ON TRAINING

One import motivation of VFG is that we aim to develop a model that is tractable of inference on datasets. A very simple application case with inference is to impute the missing values of data samples. It motivates us to use the imputation loss as the objective along with the KL terms from the ELBO as regularization terms.

Inference on a VFG model requires the aggregation node's state to be imputed from observed children's state, as shown in (??). Then, unobserved children's state can be inferred from their parent. The inference ability of VFG via imputation 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

$$\begin{aligned} \mathcal{L}(\mathbf{x}, O_{\mathbf{x}}; \theta) &= \sum_{t: 1 \leq t \leq k, t \notin O} \mathbb{E}_{q(\mathbf{h}^1|\mathbf{x}^{O_{\mathbf{x}}})} [\log p(\mathbf{x}^{(t)}|\widehat{\mathbf{h}}^1)] \\ &\quad - \sum_{l=1}^{L-1} \mathbb{E}_{q(\mathbf{h}|\mathbf{x})} [\log q(\mathbf{h}^l|\mathbf{h}^{l-1}) - \log p(\mathbf{h}^l|\widehat{\mathbf{h}}^{l+1})] \\ &\quad - \text{KL}(q(\mathbf{h}^L|\mathbf{h}^{L-1})|p(\mathbf{h}^L)). \end{aligned} \quad (18)$$

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. Considering a minibatch of training samples, the objectives function in (??) and (??) can thus be optimized sequentially. The training with *random masking* is described in Algorithm ??.

**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 (??) with Line 4 to Line 15 in Algorithm ??;
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:    $\widehat{\mathbf{h}}^{(i)} = \mathbf{h}^{(i)}$  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:      $\widehat{\mathbf{h}}^{(i)} = \frac{1}{|pa(i)|} \sum_{j \in pa(i)} \mathbf{f}_{(i,j)}^{-1}(\mathbf{h}^{(j)})$ ;
15:   end for
16:    $\mathbf{h} = \{\mathbf{h}^{(i)} | i \in \mathcal{V} \cap O\}$ ,  $\widehat{\mathbf{h}} = \{\widehat{\mathbf{h}}^{(i)} | i \in \mathcal{V}\}$ ;
17:   Approximate the KL terms in ELBO for each layer with  $b$ 
   samples;
18:   Updating VFG with gradient of (??):  $\theta_{\mathbf{f}}^{(s+1)} = \theta_{\mathbf{f}}^{(s)} +$ 
    $\nabla_{\theta_{\mathbf{f}}} \frac{1}{b} \sum_{i=1}^b \mathcal{L}(\mathbf{x}_b, O_{\mathbf{x}}; \theta_{\mathbf{f}}^{(s)})$ ,
19: end for

```

---

In practice, we use Algorithm 2 along with Algorithm 1 to train a VFG model. Training with Algorithm 1 is to improve the data fitting and the *consistency* of aggregation nodes. Geometrically speaking, training with Algorithm 2 is to enhance the learned relational knowledge among leaf nodes. Though Algorithm 2 improves the model's inference capability in a heuristic way, it works in practice.

## E AGGREGATION NODE

We present additional details on aggregation nodes here. 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, at an aggregation node  $i$  that has multiple ( $|ch(i)|$ ) children, its latent variable in forward message passing is the mean of all children's output, i.e.,

$$\mathbf{h}^{(i)} = \frac{1}{|ch(i)|} \sum_{j \in ch(i)} \mathbf{f}_{(j,i)}(\mathbf{h}^{(j)}). \quad (19)$$

On the other hand, if node  $i$  in a DAG has multiple parents, the reconstruction of its latent variable is the mean of all parents' output, i.e.,

$$\widehat{\mathbf{h}}^{(i)} = \frac{1}{|pa(i)|} \sum_{j \in pa(i)} \mathbf{f}_{(i,j)}^{-1}(\widehat{\mathbf{h}}^{(j)}). \quad (20)$$

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

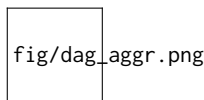


Figure 9: Aggregation node on a DAG.

Besides averaging, the aggregation nodes also ensure the latent variable on the two ends of an identity function are *consistent*. We use node  $i$  in the DAG presented in Figure ?? as an example. 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 (??) and (??), we have  $\mathbf{h}^{(i)} = (\mathbf{h}^{(d)} + \mathbf{h}^{(e)})/2$  and  $\widehat{\mathbf{h}}^{(i)} = (\widehat{\mathbf{h}}^{(u)} + \widehat{\mathbf{h}}^{(v)})/2$ . Here aggregation *consistency* means, for  $i$ 's children, their forward state should be consistent with  $i$ 's backward state, i.e.,

$$\mathbf{h}^{(d)} = \mathbf{h}^{(e)} = \widehat{\mathbf{h}}^{(i)}. \quad (21)$$

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

$$\widehat{\mathbf{h}}^{(u)} = \widehat{\mathbf{h}}^{(v)} = \mathbf{h}^{(i)}. \quad (22)$$

We utilize the KL term in the ELBOs to ensure (??) and (??) can be satisfied during parameter updating. The KL term regarding node  $i$  is

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

Here

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

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

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

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

With

$$\log p(\mathbf{h}^{(d)}|\widehat{\mathbf{h}}^{pa(d)}) = \log p(\mathbf{h}^{(d)}|\widehat{\mathbf{h}}^{(i)}) = -\|\mathbf{h}^{(d)} - \widehat{\mathbf{h}}^{(i)}\|_1 - d \cdot \log 2,$$

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

## F DERIVATION OF THE ELBOS FOR TREES AND DAGS

### F.1 ELBO of Tree Models

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  $\theta$ . 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 ?. Then the marginal likelihood term of the data reads

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



**Figure 10: A tree VFG with  $L = 5$  and three aggregation nodes.**

The hierarchical generative model is given by factorization

$$p(\mathbf{h}) = p(\mathbf{h}^L) \prod_{l=1}^{L-1} p(\mathbf{h}^l | \mathbf{h}^{l+1}). \quad (23)$$

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

$$q_{\theta}(\mathbf{h} | \mathbf{x}) = q(\mathbf{h}^1 | \mathbf{x}) q(\mathbf{h}^2 | \mathbf{h}^1) \cdots q(\mathbf{h}^L | \mathbf{h}^{L-1}). \quad (24)$$

Draw samples from the generative model (??) involves sequential conditional sampling from the top of the tree to the bottom, and computation of the recognition model (??) 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

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

$$p(\mathbf{h}^{l:L}) = p(\mathbf{h}^l | \mathbf{h}^{l+1:L}) p(\mathbf{h}^{l+1:L}) = p(\mathbf{h}^l | \mathbf{h}^{l+1}) p(\mathbf{h}^{l+1:L}) \quad (26)$$

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

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

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

$$\begin{aligned} \log p(\mathbf{x}) &\geq \mathcal{L}(\mathbf{x}; \theta) \\ &= \mathbb{E}_{q(\mathbf{h}^{1:L} | \mathbf{x})} [\log p(\mathbf{x} | \mathbf{h}^{1:L}) - \log q(\mathbf{h}^{1:L} | \mathbf{x}) + \log p(\mathbf{h}^{1:L})] \\ &= \underbrace{\mathbb{E}_{q(\mathbf{h}^{1:L} | \mathbf{x})} [\log p(\mathbf{x} | \mathbf{h}^{1:L})]}_{\text{(a) Reconstruction of the data}} \\ &\quad - \underbrace{\mathbb{E}_{q(\mathbf{h}^{1:L} | \mathbf{x})} [\log q(\mathbf{h}^{1:L} | \mathbf{x}) - \log p(\mathbf{h}^{1:L})]}_{\text{KL}^{1:L}} \end{aligned} \quad (27)$$

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 (??) can be further expanded as

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

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

$$\begin{aligned} \text{KL}^{1:L} &= \mathbb{E}_{q(\mathbf{h}^{1:L} | \mathbf{x})} [\log q(\mathbf{h}^1 | \mathbf{x}) - \log p(\mathbf{h}^1 | \mathbf{h}^2)] \\ &\quad + \log q(\mathbf{h}^{2:L} | \mathbf{h}^1) - \log p(\mathbf{h}^{2:L}) \\ &= \underbrace{\mathbb{E}_{q(\mathbf{h}^{1:L} | \mathbf{x})} [\log q(\mathbf{h}^1 | \mathbf{x}) - \log p(\mathbf{h}^1 | \mathbf{h}^2)]}_{\text{KL}^1} + \underbrace{\mathbb{E}_{q(\mathbf{h}^{1:L} | \mathbf{x})} [\log q(\mathbf{h}^{2:L} | \mathbf{h}^1) - \log p(\mathbf{h}^{2:L})]}_{\text{KL}^{2:L}} \end{aligned} \quad (28)$$

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

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

With (??) and (??), it is easy to show that

$$\text{KL}^{l:L} = \underbrace{\mathbb{E}_{q(\mathbf{h}^{1:L} | \mathbf{x})} [\log q(\mathbf{h}^l | \mathbf{h}^{l-1}) - \log p(\mathbf{h}^l | \mathbf{h}^{l+1})]}_{\text{KL}^l} + \underbrace{\mathbb{E}_{q(\mathbf{h}^{1:L} | \mathbf{x})} [\log q(\mathbf{h}^{l+1:L} | \mathbf{h}^l) - \log p(\mathbf{h}^{l+1:L})]}_{\text{KL}^{l+1:L}} \quad (29)$$

The ELBO (??) can be written as

$$\mathcal{L}(\mathbf{x}; \theta) = \mathbb{E}_{q(\mathbf{h}^{1:L} | \mathbf{x})} [\log p(\mathbf{x} | \mathbf{h}^{1:L})] - \sum_{l=1}^{L-1} \text{KL}^l - \text{KL}^L. \quad (30)$$

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

$$\text{KL}^l = \mathbb{E}_{q(\mathbf{h}^{1:L} | \mathbf{x})} [\log q(\mathbf{h}^l | \mathbf{h}^{l-1}) - \log p(\mathbf{h}^l | \mathbf{h}^{l+1})]. \quad (31)$$

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 (??) with the default assumption that all latent variables are generated regarding data sample  $\mathbf{x}$ . Therefore the ELBO (??) can be simplified as

$$\mathcal{L}(\mathbf{x}; \theta) = \mathbb{E}_{q(\mathbf{h}^{1:L} | \mathbf{x})} [\log p(\mathbf{x} | \mathbf{h}^1)] - \sum_{l=1}^L \text{KL}^l. \quad (32)$$

The  $\text{KL}$  term (??) becomes

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

When  $l = L$ ,

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



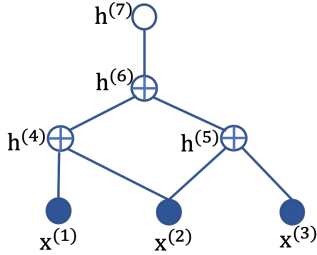
## F.2 Improve ELBO Estimation with Flows

In this paper we follow the approach in [??] using normalizing flows to further improve posterior estimation on a tree VFG model. At each layer, minimizing KL term is to optimize the parameters of the network so that the posterior is closer to the prior. As shown in Figure ??, for layer  $l$ , we can take the encoding-decoding procedures (discussed in section ??) as transformation of the posterior distribution from layer  $l$  to  $l+1$ , and then transform it back. By counting in the transformation difference [??], the KL at layer  $l$  becomes

$$\begin{aligned} \text{KL}^l &= \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} \left[ \log q(\mathbf{h}^l|\mathbf{h}^{l-1}) + \log \left| \det \frac{\partial \mathbf{h}^l}{\partial \mathbf{h}^{l+1}} \right| + \log \left| \det \frac{\partial \widehat{\mathbf{h}}^{l+1}}{\partial \widehat{\mathbf{h}}^l} \right| - \log p(\mathbf{h}^l|\widehat{\mathbf{h}}^{l+1}) \right] \\ &\approx \frac{1}{U} \sum_{u=1}^U \left[ \log q(\mathbf{h}_u^l|\mathbf{h}^{l-1}) + \log \left| \det \frac{\partial \mathbf{h}_u^l}{\partial \mathbf{h}_u^{l+1}} \right| + \log \left| \det \frac{\partial \widehat{\mathbf{h}}_u^{l+1}}{\partial \widehat{\mathbf{h}}_u^l} \right| - \log p(\mathbf{h}_u^l|\widehat{\mathbf{h}}_u^{l+1}) \right]. \end{aligned} \quad (36)$$

## F.3 ELBO of DAG Models

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 ??.



**Figure 11: A DAG with inverse topology order  $\{ \{1,2,3\}, \{4,5\}, \{6\}, \{7\} \}$ , and they correspond to layers 0 to 3.**

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 (??) we develop the ELBO as

$$\begin{aligned} \log p(\mathbf{x}) &\geq \mathcal{L}(\mathbf{x}; \theta) = \mathbb{E}_{q(\mathbf{h}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{h})}{q(\mathbf{h}|\mathbf{x})} \right] \\ &= \underbrace{\mathbb{E}_{q(\mathbf{h}|\mathbf{x})} \left[ \log p(\mathbf{x}|\mathbf{h}) \right]}_{\text{Reconstruction of the data}} - \underbrace{\mathbb{E}_{q(\mathbf{h}|\mathbf{x})} \left[ \log q(\mathbf{h}|\mathbf{x}) - \log p(\mathbf{h}) \right]}_{\text{KL}}. \end{aligned} \quad (33)$$

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

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

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

$$q(\mathbf{h}^{l:L}|\mathbf{h}^{1:l-1}) = q(\mathbf{h}^l|\mathbf{h}^{1:l-1})q(\mathbf{h}^{l+1:L}|\mathbf{h}^l) \quad (34)$$

$$= q(\mathbf{h}^l|\mathbf{h}^{1:l-1})q(\mathbf{h}^{l+1:L}|\mathbf{h}^{1:l})p(\mathbf{h}^{l:L}) = p(\mathbf{h}^l|\mathbf{h}^{l+1:L})p(\mathbf{h}^{l+1:L}) \quad (35)$$

Following (??) and with (??-??), we have

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

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,

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

For nodes in layer  $l$ ,

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

Recurrently applying (??) to (??) yields

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

For node  $i$ ,

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