

# VFG: VARIATIONAL FLOW GRAPHICAL MODEL WITH HIERARCHICAL LATENT STRUCTURE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

This paper introduces a novel approach to embed flow-based models with hierarchical latent data structures. Through a careful integration of normalizing flows in variational graphs, the proposed model uncovers the latent relational structures of high dimensional data via a message-passing scheme. Meanwhile, the model can generate data representations with reduced latent dimensions, and thus overcomes the drawbacks of many flow-based models usually requiring a high dimensional latent space involving many trivial variables. Theoretical analysis and numerical experiments on synthetic and real datasets show the benefits and broad potentials of our proposed method.

## 1 INTRODUCTION

Graphical models (Madigan et al., 1995; Hruschka et al., 2007) are powerful tools to combine the particular structure of a graph and probabilistic modeling, which provides a probabilistic (and hierarchical) characterization of variables. Due to their flexibility and ability to effectively learn and perform inference in large networks (Koller et al., 2007), they have attracted lots of interest and have been applied in many fields, *e.g.* speech recognition (Bilmes & Bartels, 2005), Quick Medical Reference (QMR) model (Shwe et al., 1990) and energy-based model (Jordan, 1999). The quantity of interest in such models is the marginal distribution of the observed data, also known as the incomplete likelihood, noted  $p(\mathbf{x})$ . Most statistical learning tasks involve a parameterized model and their training procedure involves computing the maximum likelihood estimate defined as  $\theta^* := \arg \max_{\theta \in \mathbb{R}^d} p_{\theta}(\mathbf{x})$ . A direct consequence of Bayes rule, which reads  $p_{\theta}(\mathbf{x}|\mathbf{z}) = p_{\theta}(\mathbf{z}, \mathbf{x})/p_{\theta}(\mathbf{x})$ , is that the maximization of such likelihood  $p_{\theta}(\mathbf{x})$  in a parameterized model is closely related to the inference of the density  $p_{\theta}(\mathbf{x}|\mathbf{z})$ , as a subroutine in the training procedure. Note that in the above,  $\mathbf{z}$  is the latent variable and  $p(\mathbf{x}, \mathbf{z})$  is the joint distribution of the complete data comprised of the observations  $x$  and  $z$ .

The focus of this paper is mostly on this graphical inference subroutine. There are two general approaches for this task: *exact inference* and *approximate inference*. (i) Exact inference, *e.g.* ELIMINATION ALGORITHM (Sanner & Abbasnejad, 2012) and JUNCTION TREE ALGORITHM (Kahle et al., 2008), resorts to an exact numerical calculation procedure of the quantity of interest. However, in most cases, exactly inferring from  $p_{\theta}(\mathbf{x}|\mathbf{z})$  is either *computationally involved* or simply *intractable*. It is the case for modern graphical models modeling complex tasks via deep neural networks. However it can be empirically observed that the distribution can be well determined by a small cluster of nodes in the network, see Jordan et al. (1999). There exist a trade-off between exact inference and light computations. (ii) In contrast, approximate inference, *e.g.* variational inference, yields an approximation procedures that generally provide bounds on the pdfs of interest without never attaining them. Despite such approximation and considering slow convergence issues of stochastic MCMC procedure (Salimans et al., 2015), we opt for the deterministic Variational Inference (VI) approach to tackle the graphical inference problem. VI is computationally efficient using off-the-shelf optimization techniques and easily applicable to large datasets (Hoffman et al., 2013; Kingma & Welling, 2013; Liu & Wang, 2016). In Variational Inference, mean-field approximation (Xing et al., 2012) and variational message passing (Winn & Bishop, 2005) are two common approaches for graphical models. Those methods leverage families of simple and tractable distributions to approximate the intractable posterior  $p(\mathbf{z}|\mathbf{x})$ . However, such approximation is limited by the choice of distributions that are inherently unable to recover the true posterior, often leading to a loose lower bound and they also often lack a flexible structure to learn the intrinsic disentangled latent representation.

Dealing with high dimensional data using graphical models exacerbates this systemic inability to efficiently model the latent structure of the data. Motivated by these major limitations, we propose a new framework, a variational hierarchical graphical flow model and list our contributions as follows:

- **Normalizing Flows:** A normalizing flow is introduced in the variational inference task on the original hierarchical latent data model. The result is a richer and tractable posterior distribution used as an approximation of the true posterior.
- **Hierarchical and Flow-Based:** Introducing the VARIATIONAL FLOW GRAPHICAL (VFG) model, we propose a novel graph architecture borrowing ideas from the *hierarchical latent data* modeling and *normalizing flows* concept to uncover the underlying complex structure of high dimensional data.
- **Numerical Applications:** We highlight the benefits of our VFG model on two main applications: – the graph missing entries imputation problem and – the disentanglement learning task where we specifically demonstrate that our model achieves to disentangle the factors of variation underlying the high dimensional data given as input.

Section 2 presents important concepts such as normalizing flows, VI and variational graphical models. Section 3 introduces the Variational Flow Graphical Model (VFG) model to tackle the latent relational structure learning of high dimensional data. Section 4 corresponds to theoretical findings of our model. Section 5 showcases the advantage of VFG on two different tasks: missing values imputation on both synthetic and real dataset and disentanglement learning.

## 2 PRELIMINARIES

In this section, we first introduce the general principles and notations of normalizing flows and variational inference. Then, we explain how they can naturally be embedded with graphical models.

**Normalizing flows:** Normalizing flows (Kingma & Dhariwal, 2018; Rezende & Mohamed, 2015) is a transformation of a simple probability distribution into a more complex distribution by a sequence of invertible and differentiable mappings, noted  $\mathbf{f} : \mathcal{Z} \rightarrow \mathcal{X}$  between two random variables  $z \in \mathcal{Z}$  of density  $p(\mathbf{z})$  and  $x \in \mathcal{X}$ . Firstly introduced by Tabak et al. (2010) for single maps, it has been popularized in Dinh et al. (2016); Rippel & Adams (2013) with deep neural networks for variational inference (Rezende & Mohamed, 2015). Flow-based models (Dinh et al., 2016; 2014; De Cao et al., 2020; Ho et al., 2019; Papamakarios et al., 2019) are attractive approaches for density estimation as they result in better performance enjoying the exact inference capability at a *low computational cost*. 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 a user-designed model  $\theta$ . We focus on a finite sequence of transformations  $\mathbf{f} := \mathbf{f}_1 \circ \mathbf{f}_2 \circ \dots \circ \mathbf{f}_L$  such that :

$$\mathbf{x} = \mathbf{f}(\mathbf{z}), \quad \mathbf{z} = \mathbf{f}^{-1}(\mathbf{x}) \quad \text{and} \quad \mathbf{z} \xleftrightarrow[\mathbf{f}_1^{-1}]{\mathbf{f}_1} \mathbf{h}^1 \xleftrightarrow[\mathbf{f}_2^{-1}]{\mathbf{f}_2} \mathbf{h}^2 \dots \xleftrightarrow[\mathbf{f}_L^{-1}]{\mathbf{f}_L} \mathbf{x}.$$

By defining the aforementioned invertible maps  $\{f_\ell\}_{\ell=1}^L$ , 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:

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

where we have  $\mathbf{h}^0 = \mathbf{x}$  and  $\mathbf{h}^L = \mathbf{z}$  for conciseness. The scalar value  $\log |\det(\frac{\partial \mathbf{h}^i}{\partial \mathbf{h}^{i-1}})|$  is the logarithm of the absolute value of the determinant of the Jacobian matrix  $\frac{\partial \mathbf{h}^i}{\partial \mathbf{h}^{i-1}}$ , also called the log-determinant. Identity (1) yields an easy 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).

**Variational inference:** Following the setting discussed above, the functional mapping  $\mathbf{f}: \mathbf{x} \rightarrow \mathbf{z}$  can be viewed as an encoding process and the mapping  $\mathbf{f}^{-1}: \mathbf{z} \rightarrow \mathbf{x}$  as a decoding one:  $\mathbf{z} \sim p(\mathbf{z})$ ,  $\mathbf{x} \sim p_\theta(\mathbf{x}|\mathbf{z})$ . To learn the vector of parameters  $\theta$ , we maximize the following marginal log-likelihood  $\log p_\theta(\mathbf{x}) = \log \int p(\mathbf{z})p_\theta(\mathbf{x}|\mathbf{z})d\mathbf{z}$ . 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  $p_\theta(\mathbf{z}|\mathbf{x}) \propto p(\mathbf{z})p_\theta(\mathbf{x}|\mathbf{z})$ . The goal of VI is to minimize the distance, in terms of Kullback-Leibler (KL), between the variational candidate and the true posterior  $\mathbf{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x}))$ . This optimization problem can be shown to be equivalent to maximizing the following evidence lower bound (ELBO) objective, noted  $\mathcal{L}(\mathbf{x}; \theta)$ :

$$\log p_\theta(\mathbf{x}) \geq \mathcal{L}(\mathbf{x}; \theta) = E_{p_\theta(\mathbf{x})} \{ E_{q_\phi(\mathbf{z}|\mathbf{x})} \log p_\theta(\mathbf{x}|\mathbf{z}) - \mathbf{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) \}. \quad (2)$$

**Variational graphical models:** In Directed Acyclic Graph (DAG) models, each node  $\mathbf{v}$  corresponds to a random variable, *e.g.*  $\mathbf{v}$  include the latent variables  $\mathbf{z}$  and observed variables  $\mathbf{x}$  in the variational framework. The edges represent the statistical dependencies between the variables, it can be for instance a function  $\mathbf{f}_\theta$  parameterized by  $\theta$  which serves as a link function between two variables. The joint distribution of the model is given by  $p_\theta(\mathbf{v}) = \prod_{\mathbf{v} \in \mathcal{V}} p_\theta(\mathbf{v}|pa(\mathbf{v}))$ , where  $\mathbf{v} = (\mathbf{z}, \mathbf{x})$ ,  $\mathcal{V}$  is a sample space for all graph variables and  $pa(\mathbf{v})$  denotes the parent node of  $\mathbf{v}$ . The goal of variational Bayesian networks, as a special instance of variational graphical models, is to find a variational distribution, noted  $q(\mathbf{z}|\mathbf{x})$ , approximating  $p(\mathbf{z}|\mathbf{x})$ . In this paper, we focus on the factorization of the independent and disjoint latent variables (Bishop et al., 2003) such that  $q(\mathbf{z}|\mathbf{x}) = \prod_i q_i(\mathbf{z}_i)$ , where  $\mathbf{z}_i$  is the latent variable at node  $i$  of the graph, assuming that  $\mathbf{x} = pa(\mathbf{z}_i)$ .

### 3 VARIATIONAL FLOW GRAPHICAL MODEL

Assume that there exist a sequence of variables that maps the latent and the observation sets. 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).

#### 3.1 THE EVIDENCE LOWER BOUND OF VARIATIONAL FLOW GRAPHICAL MODELS

Figure 1 illustrates the tree structure induced by variational flows. The hierarchical generative network is comprised of  $L$  layers,  $\mathbf{h}^l$  denotes the latent variable in layer  $l$  and  $\theta$  is the vector of parameters of the model. The hierarchical generative process of the model is defined as:

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

The probability density function  $p_{\theta_f}(\mathbf{h}^{l-1}|\mathbf{h}^l)$  is modeled with an invertible normalizing flow function. The hierarchical recognition network is factorized as

$$q_{\theta_f}(\mathbf{h}|\mathbf{x}) = q_{\theta_f}(\mathbf{h}^1|\mathbf{x}) q_{\theta_f}(\mathbf{h}^2|\mathbf{h}^1) \cdots q_{\theta_f}(\mathbf{h}^L|\mathbf{h}^{L-1}),$$

where  $\mathbf{h} = \{\mathbf{h}^1, \dots, \mathbf{h}^L\}$  denotes all latent variables of the model. At node  $i$ , the invertible function  $\mathbf{h}^{(i)}$  is used as the forward evidence message received from its children, and  $\hat{\mathbf{h}}^{(i)}$  as the reconstruction of  $\mathbf{h}^{(i)}$  with backward message received from the root. We denote by  $ch(i)$  and  $pa(i)$ , the node  $i$ 's child set and parent, respectively. 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

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

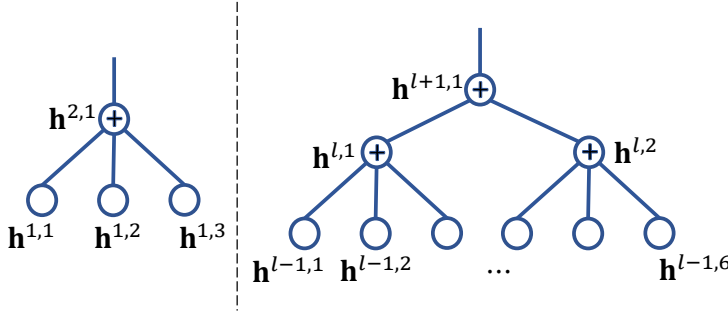
The inference procedure includes forward and backward message passing corresponding to the encoding and decoding steps, respectively. With  $\mathbf{h}^0 = \mathbf{x}$ , the layer-wise ELBO (for latent states in each layer) can be derived as

$$\mathcal{L}(\mathbf{x}; \theta_f) = \sum_{l=0}^{L-1} \mathbb{E}_{q(\mathbf{h}^{l+1}|\mathbf{h}^l)} \left[ \log p(\mathbf{h}^l|\hat{\mathbf{h}}^{l+1}) \right] + \sum_{l=1}^{L-1} \mathbf{H}(\mathbf{h}^l|\mathbf{h}^{l-1}) - \mathbf{KL}(q(\mathbf{h}^L|\mathbf{h}^{L-1})|p(\mathbf{h}^L)). \quad (3)$$

The details of the derivation of the ELBO can be found in the Appendix. The first term of ELBO is the reconstruction term for each layer:  $\mathbf{x}$  and the latent representations  $\mathbf{h}^1, \dots, \mathbf{h}^{L-1}$  where the

model pushes the variational distribution to fit the observed data. At layer  $l$ , the reconstruction  $\hat{\mathbf{h}}^l$  is generated based on  $\hat{\mathbf{h}}^{l+1}$ . Optimizing the reconstruction term  $\log p(\mathbf{h}^l | \hat{\mathbf{h}}^{l+1})$  is equivalent to force latent value  $\mathbf{h}^l$  close to its reconstruction  $\hat{\mathbf{h}}^l$ , i.e.,  $\mathbf{h}^l = \hat{\mathbf{h}}^l$ . At the root node, we have  $\hat{\mathbf{h}}^L = \mathbf{h}^L$ .

The second and third terms are some regularizations term for the latent representation where the negated **KL** term in the third position appears to keep the model near the prior distribution of the nodes. A trade-off is thus performed via such a loss function. Invertible functions are employed to connect the nodes of the studied graph as in flow-based models (Dinh et al., 2016) in order to achieve tractable message passing.



**Figure 1:** (Left) The structure of one node. Node  $\mathbf{h}^{2,1}$  connects with its children with invertible functions. The messages from its children are aggregated at  $\mathbf{h}^{2,1}$ . (Right) An illustration of the latent structure from layer  $l-1$  to  $l+1$ .  $\mathbf{h}^{l,i}$  means the  $i$ th latent variable in layer  $l$ .

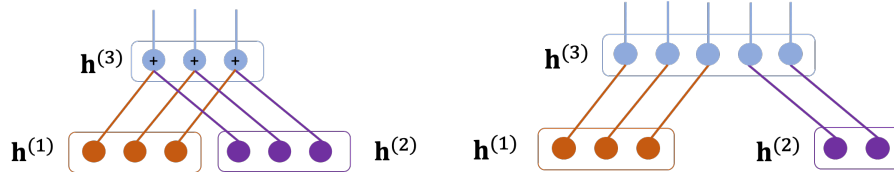
As shown in Figure 1- (Left), a node in a flow-graph can have multiple children and multiple parents. Each node has the forward messages from the input data samples and the backward messages from the root. If all the nodes have only one parent, then the structure becomes a tree. If there several nodes have multiple parents, the graph will be a DAG (directed acyclic graph). It is easy to extend the computation of

the ELBO (3) to DAGs with topology ordering of the nodes and thus the layer number. We develop the ELBO for a DAG structure as follows:

$$\begin{aligned} \log p(\mathbf{x}) \geq \mathcal{L}(\mathbf{x}; \theta_{\mathbf{f}}) = & \sum_{i \in \mathcal{G} \setminus \mathcal{R}_{\mathcal{G}}} \mathbb{E}_{q(\mathbf{h}^{pa(i)} | \mathbf{h}^{ch(pa(i))})} \left[ \log p(\mathbf{h}^{(i)} | \hat{\mathbf{h}}^{pa(i)}) \right] \\ & + \sum_{i \in \mathcal{G} \setminus \mathcal{R}_{\mathcal{G}}} \mathbf{H}(\mathbf{h}^{(i)} | \mathbf{h}^{ch(i)}) - \sum_{i \in \mathcal{R}_{\mathcal{G}}} \mathbf{KL}(q(\mathbf{h}^{(i)} | \mathbf{h}^{ch(i)}) | p(\mathbf{h}^{(i)})). \end{aligned} \quad (4)$$

Here  $\mathcal{G}$  stands for the node set of the GAG, and  $\mathcal{R}_{\mathcal{G}}$  is the set of root, or source, nodes.

Assume there are  $k$  leaf nodes on a tree or a DAG model, and they correspond to  $k$  sections of the input sample  $\mathbf{x} = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)}]$ , then the hidden variables in both (3) and (4) are computed with forward and backward message passing. We provide more details about the nodes in next subsection.



**Figure 2:** (Left) Aggregation with average. (Right) Aggregation with concatenation.

### 3.2 NODE AGGREGATION

In the sequel, we consider that all nodes latent variables, noted  $\mathbf{h}^{l,i}$ , for all  $l[L]$  and  $i \in \mathbb{N}$ , admit Gaussian distributions as prior distribution. There are two approaches to aggregate signals from different nodes: – Average-based and – Concatenation-based aggregation, see Figure 2 for an illustrative scheme. While concatenation-based aggregation is simple and straightforward, we rather focus on Average-based aggregation for the purpose of this paper. We assume each entry of a hidden node follows a normal distribution, i.e.,  $\mathbf{h}_j^{(i)} \sim \mathcal{N}(\mu_j^{(i)}, \sigma^2)$  for node  $i$ 's  $j$ th entry. To avoid cumbersome notations, we use the same standard deviation  $\sigma$  across all nodes. Extending to different

values for each node does not affect the rest of the paper. Assume a model only has one average aggregation node as shown in Figure 2. According to (3), we have

$$\begin{aligned} \log p(\mathbf{x}) \geq \mathcal{L}(\mathbf{x}; \theta_f) = & \mathbb{E}_{q(\mathbf{h}^1|\mathbf{x})} [\log p(\mathbf{x}|\hat{\mathbf{h}}^1)] + \mathbf{H}(\mathbf{h}^1|\mathbf{x}) \\ & + \mathbb{E}_{q(\mathbf{h}^2|\mathbf{h}^1)} [\log p(\mathbf{h}^1|\hat{\mathbf{h}}^2)] - \mathbf{KL}(q(\mathbf{h}^2|\mathbf{h}^1)|p(\mathbf{h}^2)). \end{aligned} \quad (5)$$

Note that at an average-based aggregation node  $i$ , we have two aggregation rules: (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 children have the same reconstruction value with its parent, i.e.,  $\hat{\mathbf{h}}^{(j)} = \hat{\mathbf{h}}^{(i)}, \forall j \in ch(i)$ . Here we give an example with a model that has only one aggregation node. Let  $\mathbf{h}^{(r)}$  be the root, and it has  $k$  children,  $\mathbf{h}^{(t)}, t = 1, \dots, k$ . With  $\mathbf{f}_t$  as the flow function connecting  $\mathbf{h}^{(t)}$  and  $\mathbf{x}^{(t)}$ , according to the aggregation rules,

$$\mathbf{h}^{(t)} = \mathbf{f}_t(\mathbf{x}^{(t)}), \quad \hat{\mathbf{h}}^{(r)} = \mathbf{h}^{(r)} = \frac{1}{k} \sum_{t=1}^k \mathbf{h}^{(t)}, \quad \text{and} \quad \hat{\mathbf{h}}^{(t)} = \hat{\mathbf{h}}^{(r)}, \quad t = 1, \dots, k.$$

In Figure 2-left,  $k = 3$ .

Given one sample  $\mathbf{x}$  and since  $\hat{\mathbf{h}}^2 = \hat{\mathbf{h}}^{(r)}$  and  $\hat{\mathbf{h}}^{(t)} = \hat{\mathbf{h}}^{(r)}$ , the reconstruction terms in ELBO (5) are computed with

$$\log p(\mathbf{x}|\hat{\mathbf{h}}^1) + \log p(\mathbf{h}^1|\hat{\mathbf{h}}^2) = - \sum_{t=1}^k \left\{ \underbrace{\frac{1}{2\sigma_{\mathbf{x}}^2} \|\mathbf{x}^{(t)} - \mathbf{f}_t^{-1}(\hat{\mathbf{h}}^{(r)})\|^2}_{\text{By } \hat{\mathbf{x}}^{(t)} = \mathbf{f}_t^{-1}(\hat{\mathbf{h}}^{(t)}) = \mathbf{f}_t^{-1}(\hat{\mathbf{h}}^{(r)})} + \underbrace{\frac{1}{2\sigma^2} \|\mathbf{f}_t(\mathbf{x}^{(t)}) - \hat{\mathbf{h}}^{(r)}\|^2}_{\text{By } \hat{\mathbf{h}}^2 = \hat{\mathbf{h}}^{(r)}} \right\} + C$$

Here  $C = -dk \ln(2\pi) - \frac{dk}{2} \ln(\sigma_{\mathbf{x}}^2) - \frac{dk}{2} \ln(\sigma^2)$ . We use constant values for both  $\sigma_{\mathbf{x}}^2$  and  $\sigma^2$ , hence the value of  $C$  is constant as well. We use the latent variables from a batch of training samples to approximate the entry  $\mathbf{H}$  and  $\mathbf{KL}$  terms in (5). We can see that maximizing the ELBO will force the average aggregation node to satisfy aggregation rule (b). We take the parent and children involved an aggregation operation as one node in the graphical figures, e.g., Figure 1.

### 3.3 INFERENCE ON SUB-GRAPHS

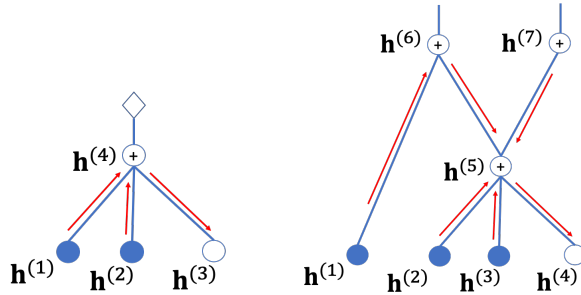
Given a trained VFG model, we can infer the state of a node given the observed nodes. Relations between variables at different nodes can also be inferred via our flow-based graphical model.

The hidden state of the parent node  $s$  in a single aggregation model can be approximated by the observed children,  $\mathbf{h}^{(s)} = \frac{1}{|ch(s)|} \sum_{i \in ch(s) \cap O} \mathbf{h}^{(i)}$ . Here  $O$  is the set of observed leaf nodes. Figure 3-left illustrates one example of this case.

Observe that for either a tree or a DAG, the state of any given node is updated via messages received from its children. The message passing firstly occurs from the children to the parent with updating and then pass it back to the children without updating. Figure 3 illustrates this inference mechanism for both trees and DAGS. The tree and DAG structures enable the model to perform message passing among the nodes. We now establish the following Lemma regarding the relation between two leaf nodes:

**Lemma 1.** Let  $\mathcal{G}$  be a trained tree structured variational flow graphical model 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  $\hat{\mathbf{x}}^j \approx \mathbf{f}_{(a,j)}(\mathbf{f}_{(i,a)}(\mathbf{x}^{(i)}))$ . Here  $\mathbf{f}_{(i,a)}$  is the flow function path from node  $i$  to node  $a$ . The conditional density of  $\mathbf{x}^{(j)}$  given  $\mathbf{x}^{(i)}$  can be approximated by:

$$\log p(\mathbf{x}^{(j)}|\mathbf{x}^{(i)}) \approx \log p(\hat{\mathbf{h}}^L) - \frac{1}{2} \log (\det (\mathbf{J}_{\hat{\mathbf{x}}^{(j)}}(\hat{\mathbf{h}}^L)^\top \mathbf{J}_{\hat{\mathbf{x}}^{(i)}}(\hat{\mathbf{h}}^L))). \quad (6)$$



**Figure 3:** (Left) Inference of single aggregation node model. Node 4 aggregates from node 1 and 2, and pass the updated state to node 3 for prediction. (Right) Inference on a DAG model. Observed node states are gathered in node 5 to predict the state of node 4.

where we recall that using the normalizing flow equation (1), we have the following identity for each node of the graph structure:

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

The proof of Lemma 1 can be found in the appendix.

**Remark 1.** Let  $O$  be the set of observed leaf nodes,  $j$  be an unobserved node, and  $a$  is the closest ancestor of  $O \cup \{a\}$ . Then the state of  $j$  can be imputed with  $\hat{\mathbf{x}}^j \approx \mathbf{f}_{(a,j)}(\mathbf{f}_{(O,a)}(\mathbf{x}^{(i)}))$ . We denote  $\mathbf{f}_{(O,a)}$  as the flow function path from all nodes in  $O$  to  $a$ , and approximation (7) still holds for  $p(\mathbf{x}^{(j)}|\mathbf{x}^O)$ .

We note that these results can be easily extended to DAG models.

### 3.4 ALGORITHM AND IMPLEMENTATION

In this section, we develop the training algorithm, see Algorithm 1, that outputs the fitted vector of parameters resulting from the maximization of the ELBO objective function (equation 3) or (equation 4) depending on what graph structure is used. In Algorithm 1, the inference of the latent variables is performed via forward message passing, cf. Line 5, and their reconstructions are computed in backward message passing, cf. Line 6.

Different from VAE, the variance of latent variables in a VFG are set with a fixed value 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. We use the empirical variance in a batch of training samples to approximate the entropy and **KL** terms. Ignoring explicit neural network parameterized variances for all latent nodes enable us to use flow-based models as the encoders as well as the decoders. It also helps the model to get rid of sampling steps and to obtain a deterministic ELBO objective equation 3 and equation 4 that can be efficiently optimized with standard stochastic solvers.

---

#### Algorithm 1 Inference model parameters with forward and backward message propagation

---

```

1: Input: Data distribution  $\mathcal{D}$ ,  $\mathcal{G} = \{\mathcal{V}, \mathbf{f}\}$ 
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:      $\mathbf{h}^{(i)} = \frac{1}{|ch(i)|} \sum_{j \in ch(i)} \mathbf{f}_{(j,i)}(\mathbf{h}^{(j)})$ ; // forward message passing
6:      $\hat{\mathbf{h}}^{(i)} = \frac{1}{|pa(i)|} \sum_{j \in pa(i)} \mathbf{f}_{(i,j)}^{-1}(\hat{\mathbf{h}}^{(j)})$ ; // backward message passing
7:   end for
8:    $\mathbf{h} = \{\mathbf{h}^{(1)}, \dots, \mathbf{h}^{(|\mathcal{V}|)}\}$  and  $\hat{\mathbf{h}} = \{\hat{\mathbf{h}}^{(1)}, \dots, \hat{\mathbf{h}}^{(|\mathcal{V}|)}\}$ ;
9:   Approximate the entry H and KL terms in ELBO for each layer with  $b$  samples;
10:  Updating flow-graph  $\mathcal{G}$  using SGD:  $\theta^{(k+1)} = \theta^{(k)} - \nabla_{\theta} \frac{1}{b} \sum_{i=1}^b \mathcal{L}(\mathbf{x}_b; \theta^{(k)})$ .
11: end for
```

---

## 4 THEORY

The proposed VFG models provide approaches to integrate multi-modal data or data sets from different sources. With invertible flow functions, we try to analyse the identifiability (Khemakhem et al., 2020; Sorrenson et al., 2020) of the VFG in this subsection. We assume each data sample has  $k$  section, and  $\mathbf{h}^{(t)}$  is the latent variable for section  $t$ ,  $\mathbf{x}^{(t)}$ . Suppose the distribution regarding latent variable  $\mathbf{h}^{(t)}$  is a factorial member of the exponential family with  $m$  sufficient statistics, conditioned on  $\mathbf{u}$ . Here  $\mathbf{u}$  is additional observed variable. The general form of the distribution can be written as

$$p_{\mathbf{h}^{(t)}}(\mathbf{h}^{(t)}|\mathbf{u}) = \Pi_{i=1}^d \frac{Q_i(h^{(t,i)})}{Z_i(\mathbf{u})} \exp \left[ \sum_{j=1}^m T_{i,j}(h^{(t,i)}) \lambda_{i,j}(\mathbf{u}) \right]. \quad (7)$$

Here  $Q_i$  is the base measure,  $Z_i(\mathbf{u})$  is the normalizing constant,  $T_{i,j}$  are the component of the sufficient statistic and  $\lambda_{i,j}$  the corresponding parameters, depending on  $\mathbf{u}$ . Variable  $\mathbf{x}^{(t)}$  is generated



with some complex, invertible, and deterministic function from the latent space:  $\mathbf{x}^{(t)} = \mathbf{f}_t^{-1}(\mathbf{h}^{(t)}, \epsilon)$ . Let  $\mathbf{T} = [\mathbf{T}_1, \dots, \mathbf{T}_l]$ , and  $\lambda = [\lambda_1, \dots, \lambda_l]$ . We define the domain of  $\mathbf{f}_t^{-1}$  as  $\mathcal{H} = \mathcal{H}_1 \times \dots \times \mathcal{H}_l$ . Here  $Q_i$  is the base measure,  $Z_i(\mathbf{u})$  is the normalizing constant,  $T_{i,j}$  are the component of the sufficient statistic and  $\lambda_{i,j}$  the corresponding parameters, depending on  $\mathbf{u}$ . Variable  $\mathbf{x}^{(t)}$  is generated with some complex, invertible, and deterministic function from the latent space:  $\mathbf{x}^{(t)} = \mathbf{f}_t^{-1}(\mathbf{h}^{(t)}, \epsilon)$ . Let  $\mathbf{T} = [\mathbf{T}_1, \dots, \mathbf{T}_l]$ ,  $\lambda = [\lambda_1, \dots, \lambda_l]$ , and  $\Theta = \{\theta := (\mathbf{T}, \lambda, \mathbf{f}_k^{-1})\}$ . We define the domain of  $\mathbf{f}_t^{-1}$  as  $\mathcal{H} = \mathcal{H}_1 \times \dots \times \mathcal{H}_l$ . We use  $\hat{\Theta} = \{\hat{\theta} := (\hat{\mathbf{T}}, \hat{\lambda}, \mathbf{g})\}$  to represent the model learned by a piratical algorithm. In the limit of infinite data and perfect convergence, we have the following theorem regarding the identifiability of  $\mathbf{T}$ .

**Theorem 1.** Assume we observe data distributed according to the generative model given by equation 8 and  $\mathbf{x}^{(t)} = \mathbf{f}_t^{-1}(\mathbf{h}^{(t)}, \epsilon)$ , we further have the following assumptions,

- (a) The sufficient statistics  $T_{ij}(h)$  are differentiable almost everywhere and their derivatives  $\frac{dT_{i,j}}{dh}$  are nonzero almost surely for all  $h \in \mathcal{H}_i$  and all  $1 \leq i \leq d$  and  $1 \leq j \leq m$ .
- (b) There exist  $dm + 1$  distinct conditions  $\mathbf{u}^{(0)}, \dots, \mathbf{u}^{(dm)}$  such that the matrix

$$\mathbf{L} = [\lambda(\mathbf{u}^{(1)}) - \lambda(\mathbf{u}^{(0)}), \dots, \lambda(\mathbf{u}^{(dm)}) - \lambda(\mathbf{u}^{(0)})]$$

of size  $dm \times dm$  is invertible. Then the model parameters  $\mathbf{T}(\mathbf{h}_k) = \mathbf{A}\hat{\mathbf{T}}(\mathbf{h}_k) + \mathbf{c}$ . Here  $\mathbf{A}$  is an  $dm \times dm$  invertible matrix and  $\mathbf{c}$  is a vector of size  $dm$ .

The proof of Theorem 1 and further analysis can be found in the appendix.

## 5 NUMERICAL EXPERIMENTS

We present in this section several numerical experiments to highlight the benefits of our VFG model. The first main application we present consists in missing values imputation. Several baseline models are compared with our newly introduced one on both synthetic and real datasets. The second application we present is to learn the disentangled latent representations that separate the explanatory factors of variations in the data, see Bengio et al. (2013). For that latter application, the model is trained and evaluated on the MNIST handwritten digits dataset.

### 5.1 MISSING ENTRIES IMPUTATION

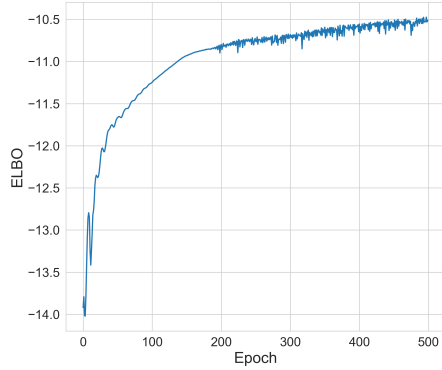
We now focus on the task of imputing missing entries in a graph structure. For all the following experiments, the models are trained on the training set and are used to infer the missing entries of samples in the testing set.

**Baseline Methods:** We use the following baselines for data imputation:

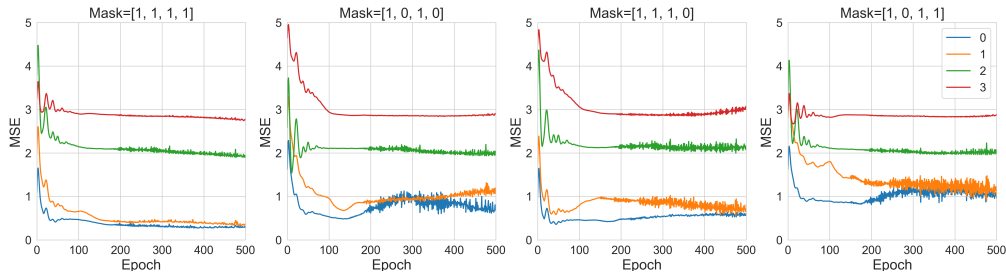
- *Mean Value:* We can directly use the mean values in the corresponding position of training set to replace the missing entries in the testing set.
- *Iterative Imputation:* A strategy for imputing missing values by modeling each feature with missing values as a function of other features in a Round-Robin fashion. We choose the KNeighborRegressor as the specific function (Pedregosa et al., 2011).
- *KNN:* To use K-Nearest Neighbor for data imputation, we compare the non-missing entries of each sample to the training set and use the average of top  $k$  samples to impute the missing entries.
- *Multivariate Imputation by Chained Equation (MICE):* This method impute the missing entries with multiple rounds of inference. This method can handle different type of data.

**Evaluation with Synthetic Data:** In this set of experiments, we study the proposed model with synthetic datasets. We generate a synthetic dataset of 1 000 data points for the training phase of the model, where each data sample has 8 dimensions with 2 latent variables. The relation between the latent variables and the [BK: to complete]

Figure 4 gives the ELBO values of the proposed method.

**Figure 4:** ELBO on the synthetic data

Methods	Imputation MSE
Mean Value	8.43
MICE	8.38
Iterative Imputation	2.64
KNN (k=3)	0.14
KNN (k=5)	0.18
Proposed	1.45

**Table 1:** Imputation Results on Synthetic Data.**Figure 5:** Imputation with Mask on the child nodes [0, 1, 2, 3] indicated by colored legends.

**Arrhythmia Dataset:** We further investigate the method on a tabular dataset. The Arrhythmia (Dua & Graff, 2017) dataset is obtained from the ODDS repository. The smallest classes, including 3, 4, 5, 7, 8, 9, 14, and 15, are combined to form the anomaly class, and the rest of the classes are combined to form the normal class. Table ?? shows the anomaly detection results with different methods.

## 5.2 DISENTANGLEMENT ON MNIST

In this set of experiments, we evaluate disentanglement of VFGs on MNIST data. Disentanglement evaluation metrics have been proposed by Karaletsos et al. (2015); Glorot et al. (2011); Higgins et al. (2017); Kim & Mnih (2018); Eastwood & Williams (2018). The authors of Chen et al. (2018) proposed to use the mutual information gap (MIG) between latent variables and the ground truth factors to measure disentanglement. In this manuscript, we utilize the approach proposed in Eastwood & Williams (2018) for the following experiments. Eastwood & Williams (2018) is a regression-based approach dividing the latent space data into training, evaluation, and testing. The disentanglement score is obtained based on the performance of the learned regression model.

## 6 CONCLUSION



## REFERENCES

- Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013.
- Jeff A Bilmes and Chris Bartels. Graphical model architectures for speech recognition. *IEEE signal processing magazine*, 22(5):89–100, 2005.
- Christopher M Bishop, David Spiegelhalter, and John Winn. Vibes: A variational inference engine for bayesian networks. In *Advances in neural information processing systems*, pp. 793–800, 2003.
- Ricky T. Q. Chen, Xuechen Li, Roger Grosse, and David Duvenaud. Isolating sources of disentanglement in vaes. In *NeurIPS*, 2018.
- Nicola De Cao, Wilker Aziz, and Ivan Titov. Block neural autoregressive flow. In *Uncertainty in Artificial Intelligence*, pp. 1263–1273. PMLR, 2020.
- Laurent Dinh, David Krueger, and Yoshua Bengio. Nice: Non-linear independent components estimation. *arXiv preprint arXiv:1410.8516*, 2014.
- Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp. *ArXiv*, abs/1605.08803, 2016.
- Dheeru Dua and Casey Graff. UCI machine learning repository, 2017. URL <http://archive.ics.uci.edu/ml>.
- Cian Eastwood and Christopher K. I. Williams. A framework for the quantitative evaluation of disentangled representations. In *ICLR*, 2018.
- Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Domain adaptation for large-scale sentiment classification: A deep learning approach. In *ICML*, 2011.
- Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. In *ICLR*, 2017.
- Jonathan Ho, Xi Chen, Aravind Srinivas, Yan Duan, and Pieter Abbeel. Flow++: Improving flow-based generative models with variational dequantization and architecture design. *arXiv preprint arXiv:1902.00275*, 2019.
- Matthew D Hoffman, David M Blei, Chong Wang, and John Paisley. Stochastic variational inference. *The Journal of Machine Learning Research*, 14(1):1303–1347, 2013.
- Estevam R Hruschka, Eduardo R Hruschka, and Nelson FF Ebecken. Bayesian networks for imputation in classification problems. *Journal of Intelligent Information Systems*, 29(3):231–252, 2007.
- Michael I. Jordan (ed.). *Learning in Graphical Models*. MIT Press, Cambridge, MA, USA, 1999. ISBN 0262600323.
- Michael I Jordan, Zoubin Ghahramani, Tommi S Jaakkola, and Lawrence K Saul. An introduction to variational methods for graphical models. *Machine learning*, 37(2):183–233, 1999.
- David Kahle, Terrance Savitsky, Stephen Schnelle, and Volkan Cevher. Junction tree algorithm. *Stat*, 631, 2008.
- Theofanis Karaletsos, Serge J. Belongie, and Gunnar Ratsch. Bayesian representation learning with oracle constraints. In *ICLR 2016*, 2015.
- Ilyes Khemakhem, Diederik Kingma, Ricardo Monti, and Aapo Hyvarinen. Variational autoencoders and nonlinear ica: A unifying framework. volume 108 of *Proceedings of Machine Learning Research*, pp. 2207–2217, Online, 26–28 Aug 2020. PMLR.
- Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In *ICML*, 2018.

- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In *Advances in Neural Information Processing Systems*, pp. 10215–10224, 2018.
- Daphne Koller, Nir Friedman, Lise Getoor, and Ben Taskar. Graphical models in a nutshell. *Introduction to statistical relational learning*, 43, 2007.
- Qiang Liu and Dilin Wang. Stein variational gradient descent: A general purpose bayesian inference algorithm. In *Advances in neural information processing systems*, pp. 2378–2386, 2016.
- David Madigan, Jeremy York, and Denis Allard. Bayesian graphical models for discrete data. *International Statistical Review/Revue Internationale de Statistique*, pp. 215–232, 1995.
- George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *arXiv preprint arXiv:1912.02762*, 2019.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows. *arXiv preprint arXiv:1505.05770*, 2015.
- Oren Rippel and Ryan Prescott Adams. High-dimensional probability estimation with deep density models. *arXiv preprint arXiv:1302.5125*, 2013.
- Tim Salimans, Diederik Kingma, and Max Welling. Markov chain monte carlo and variational inference: Bridging the gap. In *International Conference on Machine Learning*, pp. 1218–1226, 2015.
- Scott Sanner and Ehsan Abbasnejad. Symbolic variable elimination for discrete and continuous graphical models. In *Twenty-Sixth AAAI Conference on Artificial Intelligence*, 2012.
- Michael Shwe, Blackford Middleton, David Heckerman, Max Henrion, Eric Horvitz, Harold Lehmann, and Gregory Cooper. A probabilistic reformulation of the quick medical reference system. In *Proceedings of the Annual Symposium on Computer Application in Medical Care*, pp. 790. American Medical Informatics Association, 1990.
- Peter Sorrenson, Carsten Rother, and Ullrich Köthe. Disentanglement by nonlinear ica with general incompressible-flow networks (gin). In *Ninth International Conference on Learning Representations*, 2020.
- Esteban G Tabak, Eric Vanden-Eijnden, et al. Density estimation by dual ascent of the log-likelihood. *Communications in Mathematical Sciences*, 8(1):217–233, 2010.
- John Winn and Christopher M Bishop. Variational message passing. *Journal of Machine Learning Research*, 6(Apr):661–694, 2005.
- Eric P Xing, Michael I Jordan, and Stuart Russell. A generalized mean field algorithm for variational inference in exponential families. *arXiv preprint arXiv:1212.2512*, 2012.

## APPENDIX A. ELBO OF TREE MODELS

The hierarchy generative network as given in Figure 6. For each pair of connected nodes, the edge is linked with an invertible function. We use  $\theta$  to represent the parameters for all the edges. The forward message passing starts from  $\mathbf{x}$  and ends at  $\mathbf{h}^L$ , and backward message passing is in the reverse direction. Then the likelihood for the data is given by

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

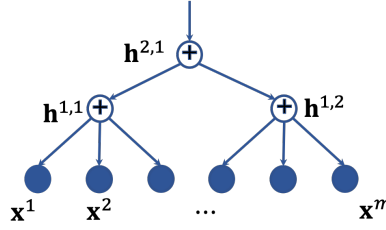


Figure 6: Tree structure.

With the flow-based ensemble model, each edge is invertible. The hierarchy of recognition network is the procedure from top to down of the structure as shown in Figure 6. Similarly, with the Markov property of the structure, the posterior density of the latent variables is given by

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

which can be simplified as

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

Note that we also have

$$q(\mathbf{h}|\mathbf{x}) = q(\mathbf{h}^1|\mathbf{x}) q(\mathbf{h}^{2:L}|\mathbf{h}^1). \quad (8)$$

To derive the ELBO of a hierarchy model, we take all layers of latent variables as the latent vector in conventional VAE, and we have

$$\begin{aligned} \log p(\mathbf{x}) &= \mathbb{E}_{q(\mathbf{h}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{h})}{p(\mathbf{h}|\mathbf{x})} \right] \\ &= \mathbb{E}_{q(\mathbf{h}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{h})}{q(\mathbf{h}|\mathbf{x})} \frac{q(\mathbf{x}, \mathbf{h})}{p(\mathbf{h}|\mathbf{x})} \right] \\ &= \underbrace{\mathbb{E}_{q(\mathbf{h}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{h})}{q(\mathbf{h}|\mathbf{x})} \right]}_{\mathcal{L}_\theta(\mathbf{x}) \text{ (ELBO)}} + \underbrace{\mathbb{E}_{q(\mathbf{h}|\mathbf{x})} \left[ \log \frac{q(\mathbf{h}|\mathbf{x})}{p(\mathbf{h}|\mathbf{x})} \right]}_{\text{KL}(q(\mathbf{h}|\mathbf{x})|p(\mathbf{h}|\mathbf{x}))}. \end{aligned}$$

Since  $\mathbf{KL}(q(\mathbf{h}|\mathbf{x})|p(\mathbf{h}|\mathbf{x})) \geq 0$  as a distance between two distributions, we obtain

$$\log p(\mathbf{x}) \geq \mathcal{L}_\theta(x) \quad (9)$$

$$\begin{aligned} &= \mathbb{E}_{q(\mathbf{h}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{h})}{q(\mathbf{h}|\mathbf{x})} \right] \\ &= \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}|\mathbf{h}^{1:L})p(\mathbf{h}^{1:L})}{q(\mathbf{h}^{1:L}|\mathbf{x})} \right] \\ &= \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} \left[ \log p(\mathbf{x}|\mathbf{h}^{1:L}) \right] + \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} \left[ \log \frac{p(\mathbf{h}^{1:L})}{q(\mathbf{h}^{1:L}|\mathbf{x})} \right] \\ &= \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} \left[ \log p(\mathbf{x}|\mathbf{h}^1) \right] + \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} \left[ \log \frac{p(\mathbf{h}^{1:L})}{q(\mathbf{h}^{1:L}|\mathbf{x})} \right] \end{aligned} \quad (10)$$

$$\begin{aligned} &= \underbrace{\mathbb{E}_{q(\mathbf{h}^1|\mathbf{x})} \left[ \log p(\mathbf{x}|\mathbf{h}^1) \right]}_{\text{Reconstruction of the data given hidden layer 1}} + \underbrace{\mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} \left[ \log \frac{p(\mathbf{h}^{1:L})}{q(\mathbf{h}^{1:L}|\mathbf{x})} \right]}_{-\mathbf{KL}^{1:L}}. \end{aligned} \quad (11)$$

The first term in equation 11 is due to  $p(\mathbf{x}|\mathbf{h}^{1:L}) = p(\mathbf{x}|\mathbf{h}^1)$ . The first term in equation 12 is due to that the expectation is regarding  $\mathbf{h}^1$ . The hidden variables  $\mathbf{h}^{l+1:L}$  can be taken as the parameters for  $\mathbf{h}^l$ 's prior distribution. We expand the minus KL term in equation 12 as follows

$$\begin{aligned} -\mathbf{KL}^{1:L} &= \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} \left[ \log \frac{p(\mathbf{h}^{1:L})}{q(\mathbf{h}^{1:L}|\mathbf{x})} \right] \\ &= \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} \left[ \log \frac{p(\mathbf{h}^1|\mathbf{h}^{2:L})p(\mathbf{h}^{2:L})}{q(\mathbf{h}^1|\mathbf{x})q(\mathbf{h}^{2:L}|\mathbf{h}^1)} \right] \\ &\quad \text{Due to equation 9} \\ &= \underbrace{\mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} \left[ \log \frac{p(\mathbf{h}^1|\mathbf{h}^{2:L})p(\mathbf{h}^{2:L})}{q(\mathbf{h}^{2:L}|\mathbf{h}^1)} \right]}_{(a)} + \underbrace{\mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} \left[ \log \frac{1}{q(\mathbf{h}^1|\mathbf{x})} \right]}_{(b)}. \end{aligned} \quad (12)$$

Given a batch of data, we take the inference in each layer as encoding and decoding procedures. In forward message passing, the hidden layer  $\mathbf{h}^l$  only depends on its previous layer  $l-1$ . The logarithm term in (a) only relates to hidden states  $\mathbf{h}^{1:L}$ . With equation 9, given the hidden states  $\mathbf{h}^1$  samples from layer 0, we have

$$(a) = \mathbb{E}_{q(\mathbf{h}^1|\mathbf{x})} \left[ \mathbb{E}_{q(\mathbf{h}^{2:L}|\mathbf{h}^1)} \left[ \log \frac{p(\mathbf{h}^1|\mathbf{h}^{2:L})p(\mathbf{h}^{2:L})}{q(\mathbf{h}^{2:L}|\mathbf{h}^1)} \right] \right]. \quad (13)$$

The inner expectation is actually the ELBO for layer hidden variable  $\mathbf{h}^1$ . Hence

$$\begin{aligned} &\mathbb{E}_{q(\mathbf{h}^{2:L}|\mathbf{h}^1)} \left[ \log \frac{p(\mathbf{h}^1|\mathbf{h}^{2:L})p(\mathbf{h}^{2:L})}{q(\mathbf{h}^{2:L}|\mathbf{h}^1)} \right] \\ &= \mathbb{E}_{q(\mathbf{h}^{2:L}|\mathbf{h}^1)} \left[ \log p(\mathbf{h}^1|\mathbf{h}^{2:L}) \right] + \mathbb{E}_{q(\mathbf{h}^{2:L}|\mathbf{h}^1)} \left[ \log \frac{p(\mathbf{h}^{2:L})}{q(\mathbf{h}^{2:L}|\mathbf{h}^1)} \right] \\ &= \mathbb{E}_{q(\mathbf{h}^2|\mathbf{h}^1)} \left[ \log p(\mathbf{h}^1|\mathbf{h}^2) \right] + \mathbb{E}_{q(\mathbf{h}^{2:L}|\mathbf{h}^1)} \left[ \log \frac{p(\mathbf{h}^{2:L})}{q(\mathbf{h}^{2:L}|\mathbf{h}^1)} \right] \\ &= \mathbb{E}_{q(\mathbf{h}^2|\mathbf{h}^1)} \left[ \log p(\mathbf{h}^1|\mathbf{h}^2) \right] - \mathbf{KL}^{2:L}. \end{aligned} \quad (14)$$

For the term (b) we develop as follows:

$$(b) = \mathbb{E}_{q(\mathbf{h}^{1:L}|\mathbf{x})} \left[ \log \frac{1}{q(\mathbf{h}^1|\mathbf{x})} \right] = \mathbb{E}_{q(\mathbf{h}^1|\mathbf{x})} \left[ \log \frac{1}{q(\mathbf{h}^1|\mathbf{x})} \right] = \mathbf{H}(\mathbf{h}^1|\mathbf{x}). \quad (15)$$

With equation 13 equation 14 equation 15 equation 16,

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

Similarly, for layer  $l$ , we have

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

Given a batch of samples, we compute and store the forward message and the backward message for each node in the forward and backward message passing procedures. The above KL term can be simplified as

$$-\mathbf{KL}^{l:L} = \mathbb{E}_{q(\mathbf{h}^{l+1}|\mathbf{h}^l)} [\log p(\mathbf{h}^l|\mathbf{h}^{l+1})] + \mathbf{H}(\mathbf{h}^l|\mathbf{h}^{l-1}) - \mathbf{KL}^{l+1:L}. \quad (16)$$

For a hierarchy model with  $L$  layers, we can recursively expand the KL term regarding the ELBO for each layer. Thus

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

With  $\mathbf{h}^0 = \mathbf{x}$ , with the ELBO can be written as

$$\log p(\mathbf{x}) \geq \sum_{l=0}^{L-1} \mathbb{E}_{q(\mathbf{h}^{l+1}|\mathbf{h}^l)} \left[ \log p(\mathbf{h}^l|\mathbf{h}^{l+1}) \right] + \sum_{l=1}^{L-1} \mathbf{H}(\mathbf{h}^l|\mathbf{h}^{l-1}) - \mathbf{KL}(q(\mathbf{h}^L|\mathbf{h}^{L-1})|p(\mathbf{h}^L)).$$

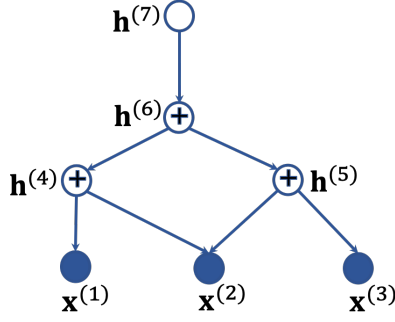
The hidden variables are computed with forward message passing with encoders  $q(\mathbf{h}^l|\mathbf{h}^{l-1})$ ,  $l = 1, \dots, L$ . The reconstructed hidden variables are computed with decoders  $p(\mathbf{h}^l|\mathbf{h}^{l+1})$ ,  $l = L-1, \dots, 0$ . We use  $\hat{\mathbf{h}}^l$  to represent the reconstruction of  $\mathbf{h}^l$ . Only at the root level  $L$ , we have  $\hat{\mathbf{h}}^L = \mathbf{h}^L$ . Each latent variable is reconstructed with messages from higher layer. Hence the ELBO can be rewritten as

$$\log p(\mathbf{x}) \geq \sum_{l=0}^{L-1} \mathbb{E}_{q(\mathbf{h}^{l+1}|\mathbf{h}^l)} \left[ \log p(\mathbf{h}^l|\hat{\mathbf{h}}^{l+1}) \right] + \sum_{l=1}^{L-1} \mathbf{H}(\mathbf{h}^l|\mathbf{h}^{l-1}) - \mathbf{KL}(q(\mathbf{h}^L|\mathbf{h}^{L-1})|p(\mathbf{h}^L)).$$

## APPENDIX B. ELBO OF DAG MODELS

If we reverse the edge directions in a DAG, the result graph is still a DAG graph. The nodes can be listed in a topological order regarding the DAG structure as shown in Figure 7. By taking the topology order as the layers in tree structures, we can derive the ELBO for DAG structures. Let's assume the DAG structure has  $L$  layers, and the root nodes are in layer  $L$ . With  $\mathbf{h}$  to represent the whole latent variables, following equation 10 we have the ELBO for the log-likelihood of data

$$\begin{aligned} \log p(\mathbf{x}) &\geq \mathcal{L}_\theta(\mathbf{x}) = \mathbb{E}_{q(\mathbf{h}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{h})}{q(\mathbf{h}|\mathbf{x})} \right] \\ &= \underbrace{\mathbb{E}_{q(\mathbf{h}^{pa(\mathbf{x})}|\mathbf{x})} \left[ \log p(\mathbf{x}|\mathbf{h}^{pa(\mathbf{x})}) \right]}_{\text{Reconstruction of the data given the parent nodes of the data}} + \underbrace{\mathbb{E}_{q(\mathbf{h}|\mathbf{x})} \left[ \log \frac{p(\mathbf{h})}{q(\mathbf{h}|\mathbf{x})} \right]}_{-\mathbf{KL}}. \end{aligned} \quad (18)$$



**Figure 7:** DAG structure. The inverse topology order is  $\{ \{1,2,3\}, \{4,5\}, \{6\}, \{7\} \}$ , and it corresponds to layers 0 to 3.

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

$$-\mathbf{KL}^{l:L} = \mathbb{E}_{q(\mathbf{h}^{pa(l)}|\mathbf{h}^l)} [\log p(\mathbf{h}^l|\mathbf{h}^{pa(l)})] + \mathbf{H}(\mathbf{h}^l|\mathbf{h}^{ch(l)}) - \mathbf{KL}^{l+1:L}. \quad (19)$$

The forward and backward messages or latent state of a node are stored in the message passing procedures. They can be used by the node's parents and children to compute the ELBO. It enables the calculation even the parents or children are not in layer  $l+1$  or  $l-1$ . For the node  $i$  in layer  $l$ ,  $pa(i)$  may have children in layers below  $l$ . Some nodes in  $l$  may not have parent, and combining with the prior, the entropy term will become an KL term in this case. Thus, we have

$$\begin{aligned} -\mathbf{KL}^{l:L} = & \sum_{i:i \in l, i \notin \mathcal{R}_G} \left\{ \mathbb{E}_{q(\mathbf{h}^{pa(i)}|\mathbf{h}^{ch(pa(i))})} [\log p(\mathbf{h}^i|\mathbf{h}^{pa(i)})] + \mathbf{H}_q(\mathbf{h}^i|\mathbf{h}^{ch(i)}) \right\} \\ & - \sum_{i \in l \cap \mathcal{R}_G} \mathbf{KL}(q(\mathbf{h}^{(i)}|\mathbf{h}^{ch(i)})|p(\mathbf{h}^{(i)})) - \mathbf{KL}^{l+1:L}. \end{aligned} \quad (20)$$

Recurrently applying equation 21 yields

$$\begin{aligned} \mathbb{E}_{q(\mathbf{h}|\mathbf{x})} \left[ \log \frac{p(\mathbf{h})}{q(\mathbf{h}|\mathbf{x})} \right] = & \sum_{l=1}^{L-1} \sum_{i:i \in l, i \notin \mathcal{R}_G} \left\{ \mathbb{E}_{q(\mathbf{h}^{pa(i)}|\mathbf{h}^{(i)})} \left[ \log p(\mathbf{h}^{(i)}|\mathbf{h}^{pa(i)}) \right] + \mathbf{H}(\mathbf{h}^i|\mathbf{h}^{ch(i)}) \right\} \\ & - \sum_{l=1}^{L-1} \sum_{i \in l \cap \mathcal{R}_G} \mathbf{KL}(q(\mathbf{h}^{(i)}|\mathbf{h}^{ch(i)})|p(\mathbf{h}^{(i)})) - \mathbf{KL}(q(\mathbf{h}^L|\mathbf{h}^{L-1})|p(\mathbf{h}^L)). \end{aligned} \quad (21)$$

Since  $L \subseteq \mathcal{R}_G$ , with  $\mathbf{h}^{(0)} = \mathbf{x}$ , equation 19, and equation 22 we have

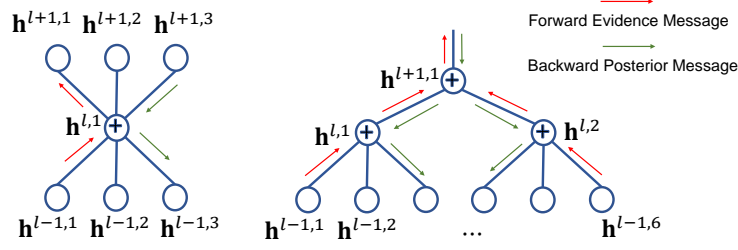
$$\begin{aligned} \log p(\mathbf{x}) \geq \mathcal{L}(\mathbf{x}; \theta) = & \sum_{i \in \mathcal{G} \setminus \mathcal{R}_G} \mathbb{E}_{q(\mathbf{h}^{pa(i)}|\mathbf{h}^{ch(pa(i))})} \left[ \log p(\mathbf{h}^{(i)}|\mathbf{h}^{pa(i)}) \right] \\ & + \sum_{i \in \mathcal{G} \setminus \mathcal{R}_G} \mathbf{H}(\mathbf{h}^{(i)}|\mathbf{h}^{ch(i)})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

## APPENDIX C. PROOF OF LEMMA 1

**Lemma 1.** Let  $\mathcal{G}$  be a well trained tree structured variational flow graphical model 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 with  $\hat{\mathbf{x}}^j \approx \mathbf{f}_{(a,j)}(\mathbf{f}_{(i,a)}(\mathbf{x}^{(i)}))$ . Here  $\mathbf{f}_{(i,a)}$  is the flow function path from node  $i$  to node  $a$ . The conditional density of  $\mathbf{x}^{(j)}$  given  $\mathbf{x}^{(i)}$  can be approximated with

$$\log p(\mathbf{x}^{(j)}|\mathbf{x}^{(i)}) \approx \log p(\hat{\mathbf{h}}^L) - \frac{1}{2} \log (\det (\mathbf{J}_{\hat{\mathbf{x}}^{(j)}}(\hat{\mathbf{h}}^L)^\top \mathbf{J}_{\hat{\mathbf{x}}^{(j)}}(\hat{\mathbf{h}}^L))).$$





**Figure 8:** (Left) Message passing in a node. (Right) Message passing in a tree.

## APPENDIX D. PROOF OF THEOREM 1

**Theorem 1.** Assume we observe data distributed according to the generative model given by equation 8 and  $\mathbf{x}^{(t)} = \mathbf{f}_t^{-1}(\mathbf{h}^{(t)}, \epsilon)$ , we further have the following assumptions,

- (a) The sufficient statistics  $T_{ij}(h)$  are differentiable almost everywhere and their derivatives  $\frac{dT_{i,j}}{dh}$  are nonzero almost surely for all  $h \in \mathcal{H}_i$  and all  $1 \leq i \leq d$  and  $1 \leq j \leq m$ .
- (b) There exist  $dm + 1$  distinct conditions  $\mathbf{u}^{(0)}, \dots, \mathbf{u}^{(dm)}$  such that the matrix

$$\mathbf{L} = [\lambda(\mathbf{u}^{(1)}) - \lambda(\mathbf{u}^{(0)}), \dots, \lambda(\mathbf{u}^{(dm)}) - \lambda(\mathbf{u}^{(0)})]$$

of size  $dm \times dm$  is invertible. Then the model parameters  $\mathbf{T}(\mathbf{h}_k) = \mathbf{A}\hat{\mathbf{T}}(\mathbf{h}_k) + \mathbf{c}$ . Here  $\mathbf{A}$  is an  $dm \times dm$  invertible matrix and  $\mathbf{c}$  is a vector of size  $dm$ .