

HIERARCHICALLY CLUSTERED REPRESENTATION LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

The joint optimization of representation learning and clustering in the embedding space has experienced a breakthrough in recent years. In spite of the advance, clustering with representation learning has been limited to flat-level categories, which oftentimes involves cohesive clustering with a focus on instance relations. To overcome the limitations of flat clustering, we introduce *hierarchically clustered* representation learning (HCRL), which simultaneously optimizes representation learning and hierarchical clustering in the embedding space. Specifically, we place a nonparametric Bayesian prior on embeddings to handle dynamic mixture hierarchies under the variational autoencoder framework, and to adopt the generative process of a hierarchical-versioned Gaussian mixture model. Compared with a few prior works focusing on unifying representation learning and hierarchical clustering, HCRL is the first model to consider a generation of deep embeddings from every component of the hierarchy, not just leaf components. This generation process enables more meaningful separations and mergers of clusters via branches in a hierarchy. In addition to obtaining hierarchically clustered embeddings, we can reconstruct data by the various abstraction levels, infer the intrinsic hierarchical structure, and learn the level-proportion features. We conducted evaluations with image and text domains, and our quantitative analyses showed competent likelihoods and the best accuracies compared with the baselines.

1 INTRODUCTION

Clustering is one of the most traditional and frequently used machine learning tasks. Clustering models are designed to represent intrinsic data structures, such as latent Dirichlet allocation (Blei et al., 2003). The recent development of *representation learning* has contributed to generalizing model feature engineering, which also enhances data representation (Bengio et al., 2013). Therefore, representation learning has been merged into the clustering models, e.g., variational deep embedding (VaDE) (Jiang et al., 2017). Besides merging representation learning and clustering, another critical line of research is structuring the clustering result, e.g., hierarchical clustering. This paper introduces a unified model enabling nonparametric Bayesian hierarchical clustering with neural-network-based representation learning.

Autoencoder (Rumelhart et al., 1985) is a typical neural network for unsupervised representation learning and achieves a non-linear mapping from a high-dimensional input space to a low-dimensional embedding space by minimizing reconstruction errors. To turn the low-dimensional embeddings into random variables, a variational autoencoder (VAE) (Kingma & Welling, 2014) places a Gaussian prior on the embeddings. The autoencoder, whether it is probabilistic or not, has a limitation in reflecting the intrinsic hierarchical structure of data. For instance, VAE assuming a single Gaussian prior needs to be expanded to suggest an elaborate clustering structure.

Due to the limitations of modeling the cluster structure with autoencoders, prior works combine the autoencoder and the clustering algorithm. While some early cases pipeline just two models, e.g., Huang et al. (2014), a typical merging approach is to model an additional loss, such as a clustering loss, in the autoencoders (Xie et al., 2016; Guo et al., 2017; Yang et al., 2017; Nalisnick et al., 2016; Chu & Cai, 2017; Xie et al., 2017). These suggestions exhibit gains from unifying the encoding and the clustering, yet they remain at the parametric and flat-structured clustering. A more recent development releases the previous constraints by using the nonparametric Bayesian approach.

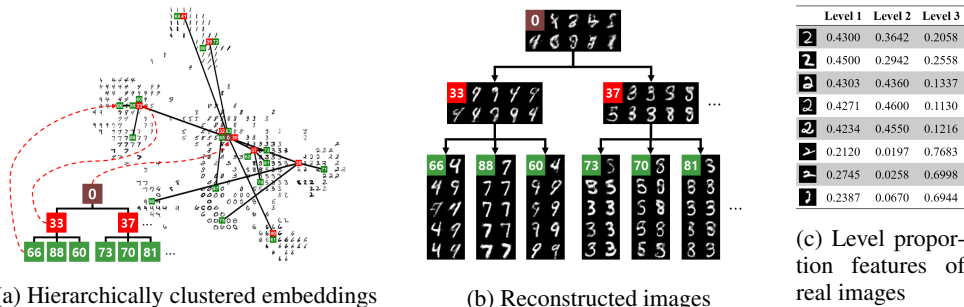


Figure 1: Example of hierarchically clustered embeddings on MNIST with three levels of hierarchy, the reconstructed digits from the hierarchical Gaussian mixture components, and the extracted level proportion features. We marked the mean of a Gaussian mixture component with the colored square, and the digit written inside the square refers to the unique index of the mixture component.

For example, the infinite mixture of VAEs (IMVAE) (Abbasnejad et al., 2017) explores the infinite space for VAE mixtures by looking for an adequate embedding space through sampling, such as the Chinese restaurant process (CRP). Whereas IMVAE remains at the flat-structured clustering, VAE-nested CRP (VAE-nCRP) (Goyal et al., 2017) captures a more complex structure, i.e., a hierarchical structure of the data, by adopting the nested Chinese restaurant process (nCRP) prior (Griffiths et al., 2004) into the cluster assignment of the Gaussian mixture model.

This paper proposes hierarchically clustered representation learning (HCRL) that is a joint model of 1) nonparametric Bayesian hierarchical clustering, and 2) representation learning with neural networks. HCRL extends a previous work on merging flat clustering and representation learning, i.e., VaDE, by incorporating inter-cluster relation modelings. Unlike a previous work of VAE-nCRP, HCRL learns the full spectrum of hierarchical clusterings, such as the level assignment and the level proportion of generating a component hierarchy. These level assignments and proportions were not modeled in VAE-nCRP, so each data instance cannot be analyzed from the perspective of generalization and specialization in a hierarchy. On the contrary, by adding level assignment and proportion modeling, a data instance can be generated from an internal component of the hierarchy, which is limited to the leaf component in VAE-nCRP. Hierarchical mixture density estimation (Vasconcelos & Lippman, 1999), where all internal and leaf components are directly modeled to generate data, is a flexible framework for hierarchical mixture modeling, such as hierarchical topic modeling (Mimno et al., 2007; Griffiths et al., 2004), with regard to the learning of the internal components.

Specifically, HCRL jointly optimizes soft-divisive hierarchical clustering in an embedding space from VAE via two mechanisms. First, HCRL includes a hierarchical-versioned Gaussian mixture model (HGMM) with a mixture of hierarchically organized Gaussian distributions. Then, HCRL sets the prior of embeddings by adopting the generative processes of HGMM. Second, to handle a dynamic hierarchy structure dealing with the clusters of unequal sizes, we explore the infinite hierarchy space by exploiting an nCRP prior. These mechanisms are fused as a unified objective function; this is done rather than concatenating the two distinct models of clustering and autoencoding. The quantitative evaluations focus on density estimation quality and hierarchical clustering accuracy, which shows that HCRL has competent likelihoods and the best accuracies compared with the baselines. When we observe our results qualitatively, we visualize 1) the hierarchical clusterings, 2) the embeddings under the hierarchy modeling, and 3) the reconstructed images from each Gaussian mixture component, as shown in Figure 1. These experiments were conducted by crossing the data domains of texts and images, so our benchmark datasets include MNIST, CIFAR-100, RCV1_v2, and 20Newsgroups.

2 PRELIMINARIES

2.1 VARIATIONAL DEEP EMBEDDING

Figure 2 presents a graphical representation and a neural architecture of VaDE (Jiang et al., 2017). The model parameters of κ , $\mu_{1:K}$, and $\sigma_{1:K}^2$, which are a proportion, means, and covariances of

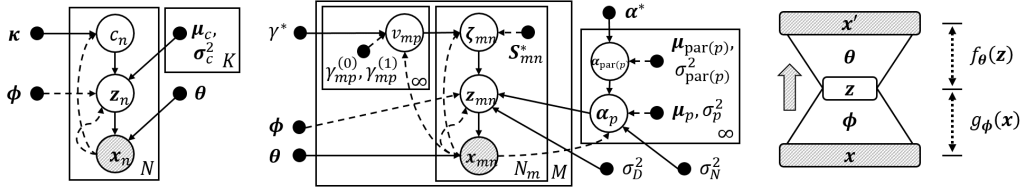


Figure 2: Graphical representation of VaDE (Jiang et al., 2017) (left), VAE-nCRP (Goyal et al., 2017) (center), and neural architecture of both models (right). In the graphical representation, the white/shaded circles represent latent/observed variables. The black dots indicate hyper or variational parameters. The solid lines represent a generative model, and dashed lines represent a variational approximation. A rectangle box means a repetition for the number of times denoted by the bottom right of the box.

mixture components, respectively, are declared outside of the neural network¹. VaDE trains model parameters to maximize the lower bound of marginal log likelihoods via the mean-field variational inference (Jordan et al., 1999). VaDE uses the Gaussian mixture model (GMM) as the prior, whereas VAE assumes a single standard Gaussian distribution on embeddings. Following the generative process of GMM, VaDE assumes that 1) the embedding draws a cluster assignment, and 2) the embedding is generated from the selected Gaussian mixture component.

VaDE uses an amortized inference as VAE, with a generative and inference networks; $\mathcal{L}(x)$ in Equation 1 denotes the evidence lower bound (ELBO), which is the lower bound on the log likelihood. It should be noted that VaDE merges the ELBO of VAE with the likelihood of GMM.

$$\log p(x) \geq \mathcal{L}(x) = \mathbb{E}_q \left[\log \frac{p(c, z, x)}{q(c, z|x)} \right] = \mathbb{E}_q \left[\log \prod_{c=1}^K \frac{\kappa_c \mathcal{N}(z|\mu_c, \sigma_c^2 I_J)}{p(c|z) \mathcal{N}(z|\tilde{\mu}, \tilde{\sigma}^2 I_J)} + \log p(x|z) \right] \quad (1)$$

2.2 VARIATIONAL AUTOENCODER NESTED CHINESE RESTAURANT PROCESS

VAE-nCRP uses the nonparametric Bayesian prior for learning tree-based hierarchies, the nCRP (Griffiths et al., 2004), so the representation could be hierarchically organized. The nCRP prior defines the distributions over children components for each parent component, recursively in a top-down way. The variational inference of the nCRP can be formalized by the nested stick-breaking construction (Wang & Blei, 2009), which is also kept in the VAE setting. The distribution over paths on the hierarchy is defined as being proportional to the product of weights corresponding to the nodes lying in each path. The weight, π_i , for the i -th node follows the Griffiths-Engen-McCloskey (GEM) distribution (Pitman et al., 2002), where π_i is constructed as $\pi_i = v_i \prod_{j=1}^{i-1} (1 - v_j)$, $v_i \sim \text{Beta}(1, \gamma)$ by a stick-breaking process. Since the nCRP provides the ELBO with the nested stick-breaking process, VAE-nCRP has a unified ELBO of VAE and the nCRP in Equation 2.

$$\mathcal{L}(x) = \mathbb{E}_q \left[\log \frac{p(v)}{q(v|x)} + \log \underbrace{\frac{p(\alpha_{par(p)}|\alpha^*) p(\alpha_p|\alpha_{par(p)}, \sigma_N^2)}{q(\alpha_p, \alpha_{par(p)}|x)}}_{(3.1)} \frac{p(\zeta|v)}{q(\zeta|x)} \underbrace{\frac{p(z|\alpha_p, \zeta, \sigma_D^2)}{q(z|x)}}_{(3.2)} + \log p(x|z) \right] \quad (2)$$

Given the ELBO of VAE-nCRP, we recognized a number of potential improvements. First, term (3.1) is for modeling the hierarchical relationship among clusters, i.e., each child is generated from its parent. VAE-nCRP trade-off is the direct dependency modeling among clusters against the mean-field variational approach. This modeling may reveal that the higher clusters in the hierarchy are more difficult to train. Second, in term (3.2), leaf mixture components generate embeddings, which implies that only leaf clusters have direct summarization ability for sub-populations. Additionally, in term (3.2), variance parameter σ_D^2 is modeled as the hyperparameter shared by all clusters. In other words, only with J -dimensional parameters, α , for the leaf mixture components, the local density modeling without variance parameters has a critical disadvantage.

For all of these weaknesses, we were able to compensate with the level proportion modeling and HGMM prior. The level assignment generated from the level proportion allows a data instance to

¹Appendix D enumerates the symbols in this paper.

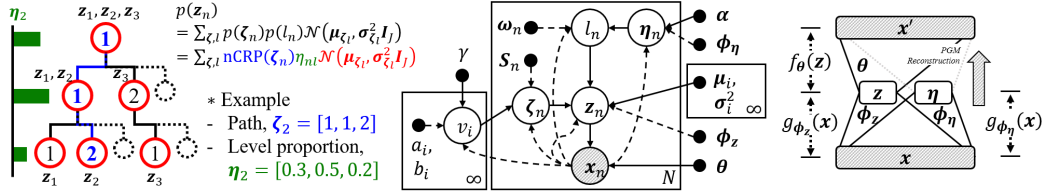


Figure 3: A simple depiction (left) of the key notations, where each numbered circle refers to the corresponding Gaussian mixture component. The graphical representation (center) and the neural architecture (right) of our proposed model, HCRL. The neural architecture of HCRL consists of two probabilistic encoder networks, g_{ϕ_η} and g_{ϕ_z} , and one probabilistic decoder network, f_θ .

select among all mixture components. We do not need direct dependency modeling between the parents and their children because all internal mixture components also generate embeddings.

3 METHODOLOGY

3.1 GENERATIVE PROCESS

The generative process of HCRL resembles the generative process of hierarchical clusterings, such as the hierarchical latent Dirichlet allocation (Griffiths et al., 2004). In detail, the generative process departs from selecting a path ζ , from the nCRP prior (phase 1). Then, we sample a level proportion (phase 2) and a level, l (phase 3), from the sampled level proportion to find the mixture component in the path, and this component of ζ_l provides the Gaussian distribution for the latent representation (phase 4). Finally, the latent representation is exploited to generate an observed datapoint (phase 5). The below formulas are the generative process with its density functions. In addition, Figure 3 illustrates a graphical representation corresponding to the described generative process. The generative process also presents our formalization of corresponding prior distributions, denoted as $p(\cdot)$, and variational distributions, denoted as $q(\cdot)$, by generation phases. The variational distributions are used in our inference methods called mean-field variational inference (MFVI) (Jordan et al., 1999) as detailed in Section 3.3.

1. Choose a path $\zeta \sim \text{nCRP}(\zeta|\gamma)$
 - $p(\zeta) = \prod_{l=1}^L \pi_{1,\zeta_2,\dots,\zeta_l}$ where $\pi_{1,\zeta_2,\dots,\zeta_l} = \prod_{l'=1}^l \{v_{1,\zeta_2,\dots,\zeta_{l'}} (\prod_{j=1}^{\zeta_{l'}-1} (1 - v_{1,\zeta_2,\dots,j}))\}$,
 $q(\zeta|\mathbf{x}) \propto S_\zeta \triangleq \sum_{\zeta \in \text{child}(\bar{\zeta})} S_\zeta$
2. Choose a level proportion $\eta \sim \text{Dirichlet}(\eta|\alpha)$
 - $p(\eta) = \text{Dir}(\eta|\alpha)$, $q_{\phi_\eta}(\eta|\mathbf{x}) = \text{Dirichlet}(\eta|\tilde{\alpha}) \approx \text{LogisticNormal}(\eta|\tilde{\mu}_\eta, \tilde{\sigma}_\eta^2 \mathbf{I}_L)$
where $[\tilde{\mu}_\eta; \log \tilde{\sigma}_\eta^2] = g_{\phi_\eta}(\mathbf{x})$, $\tilde{\alpha}_l = \frac{1}{\tilde{\sigma}_\eta^2} (1 - \frac{2}{L} + \frac{e^{-\tilde{\mu}_\eta l}}{L^2} \sum_{l'} e^{-\tilde{\mu}_\eta l'})$
3. Choose a level $l \sim \text{Multinomial}(l|\eta)$
 - $p(l) = \text{Multinomial}(\eta)$, $q(l|\mathbf{x}) = \text{Multinomial}(l|\omega)$
where $\omega_l \propto \exp \left\{ \psi(\tilde{\alpha}_l) - \psi(\tilde{\alpha}_0) + \sum_\zeta S_\zeta \left(\sum_{j=1}^J -\frac{1}{2} \log(2\pi\sigma_{\zeta_l,j}^2) - \frac{\tilde{\sigma}_{z_j}^2}{2\sigma_{\zeta_l,j}^2} - \frac{(\tilde{\mu}_{z_j} - \mu_{\zeta_l,j})^2}{2\sigma_{\zeta_l,j}^2} \right) \right\}$
4. Choose a latent representation $z \sim \mathcal{N}(z|\mu_{\zeta_l}, \sigma_{\zeta_l}^2 \mathbf{I}_J)$
 - $p(z) = \sum_{\zeta,l} p(\zeta|\gamma) \cdot \eta_l \cdot \mathcal{N}(z|\mu_{\zeta_l}, \sigma_{\zeta_l}^2 \mathbf{I}_J)$,
 $q_{\phi_z}(z|\mathbf{x}) = \mathcal{N}(z|\tilde{\mu}_z, \tilde{\sigma}_z^2 \mathbf{I}_J)$ where $[\tilde{\mu}_z; \log \tilde{\sigma}_z^2] = g_{\phi_z}(\mathbf{x})$
5. Choose an observed datapoint $\mathbf{x} \sim \mathcal{N}(\mathbf{x}|\mu_x, \sigma_x^2 \mathbf{I}_D)$ where $[\mu_x; \log \sigma_x^2] = f_\theta(z)^2$

3.2 NEURAL ARCHITECTURE

The neural architecture of HCRL consists of two probabilistic encoders on z and η , and one probabilistic decoder on z as shown in the right part of Figure 3. This unbalanced architecture originates

²We introduce the sample distribution for the real-valued data instances, and Appendix F provides the binary case as well, which we use for MNIST.

from our modeling assumption of $p(\mathbf{x}|\mathbf{z})$, not $p(\mathbf{x}|\mathbf{z}, \boldsymbol{\eta})$. The reconstruction design of \mathbf{x} depending on the two stochastic variables of \mathbf{z} and $\boldsymbol{\eta}$ may lead to a large variance of the reconstruction on \mathbf{x} . Additionally, we cannot guarantee that both \mathbf{z} and $\boldsymbol{\eta}$ contribute to the the reconstruction on \mathbf{x} (Chen et al., 2016). Although the decoding structure of $\boldsymbol{\eta}$ is not included explicitly in the neural network architecture of HCRL, we provide the formalization of $p(\boldsymbol{\eta}|\mathbf{z})$ in Table 1 according to our generative assumptions. We call this reconstruction process, which is inherently a generative process of the traditional probabilistic graphical model (PGM), *PGM reconstruction* (see the decoding neural network part of Figure 3).

Table 1: Encoding and decoding structure on \mathbf{z} and $\boldsymbol{\eta}$ in HCRL. (s) indicates the s -th sample.

	Encoding	Decoding
\mathbf{z}	$\mathbf{z} \sim q_{\phi_z}(\mathbf{z} \mathbf{x}), \mathbf{z}^{(s)} = g_{\phi_z}(\boldsymbol{\epsilon}^{(s)}, \mathbf{x})$	$\mathbf{x} \sim p_{\theta}(\mathbf{x} \mathbf{z}), \mathbf{x}^{(s)} = f_{\theta}(\mathbf{z}^{(s)})$
$\boldsymbol{\eta}$	$\boldsymbol{\eta} \sim q_{\phi_{\eta}}(\boldsymbol{\eta} \mathbf{x}), \boldsymbol{\eta}^{(s)} = g_{\phi_{\eta}}(\boldsymbol{\epsilon}^{(s)}, \mathbf{x})$	$p(\boldsymbol{\eta} \mathbf{z}) \propto \int_{\mathbf{v}, \mathbf{z}} \sum_{\zeta, \mathbf{l}} p(\mathbf{x} \mathbf{z}) p(\mathbf{z} \zeta, \mathbf{l}) p(\mathbf{l} \boldsymbol{\eta}) p(\zeta \mathbf{v}) p(\mathbf{v})$

3.3 MEAN-FIELD VARIATIONAL INFERENCE

The formal specification can be a factorized probabilistic model as Equation 3, where $\Phi = \{\mathbf{v}, \zeta, \boldsymbol{\eta}, \mathbf{l}, \mathbf{z}\}$ denotes the set of latent variables, and \mathcal{M}_T denotes the set of all nodes in tree T .

$$p(\Phi, \mathbf{x}) = \prod_{j \notin \mathcal{M}_T} p(v_j|\gamma) \prod_{i \in \mathcal{M}_T} p(v_i|\gamma) \prod_{n=1}^N p(\zeta_n|\mathbf{v}) p(\boldsymbol{\eta}_n|\boldsymbol{\alpha}) p(l_n|\boldsymbol{\eta}_n) p(\mathbf{z}_n|\zeta_n, l_n) p_{\theta}(\mathbf{x}_n|\mathbf{z}_n) \quad (3)$$

The proportion and assignment on the mixture components for the n -th data instance are modeled by ζ_n as a path assignment; $\boldsymbol{\eta}_n$ as a level proportion; and l_n as a level assignment. v is a Beta draw used in the stick-breaking construction. The latent variables are inferred through MFVI, and therefore we assume the variational distributions are as Equation 4:

$$q(\Phi|\mathbf{x}) = \prod_{j \notin \mathcal{M}_T} p(v_j|\gamma) \prod_{i \in \mathcal{M}_T} q(v_i|a_i, b_i) \prod_{n=1}^N q(\zeta_n|\mathbf{x}_n) q_{\phi_{\eta}}(\boldsymbol{\eta}_n|\mathbf{x}_n) q(l_n|\boldsymbol{\omega}_n, \mathbf{x}_n) q_{\phi_z}(\mathbf{z}_n|\mathbf{x}_n) \quad (4)$$

where $q_{\phi_{\eta}}(\boldsymbol{\eta}_n|\mathbf{x}_n)$ and $q_{\phi_z}(\mathbf{z}_n|\mathbf{x}_n)$ should be noted because these two variational distributions follow the amortized inference of VAE. $q(\zeta|\mathbf{x}) \propto S_{\bar{\zeta}} \triangleq \sum_{\zeta \in \text{child}(\bar{\zeta})} S_{\zeta}$ is the variational distribution over path ζ , where $\text{child}(\bar{\zeta})$ means the set of all full paths that are not in T but include $\bar{\zeta}$ as a sub path. Because we specified both generative and variational distributions, we define the ELBO of HCRL, $\mathcal{L} = \mathbb{E}_q \left[\log \frac{p(\Phi, \mathbf{x})}{q(\Phi|\mathbf{x})} \right]$, in Equation 5. Appendix F enumerates the full derivation in detail.

We report that the Laplace approximation with the logistic normal distribution is applied to model the prior, $\boldsymbol{\alpha}$, of the level proportion, $\boldsymbol{\eta}$. We choose a conjugate prior of a multinomial, so $p(\boldsymbol{\eta}_n|\boldsymbol{\alpha})$ follows the Dirichlet distribution. To configure the inference network on the Dirichlet prior, the Laplace approximation is used (MacKay, 1998; Srivastava & Sutton, 2017; Hennig et al., 2012).

$$\mathcal{L}(\mathbf{x}) = \mathbb{E}_q \left[\log \frac{p(\mathbf{v})}{q(\mathbf{v}|\mathbf{x})} + \log \frac{p(\boldsymbol{\eta})}{q(\boldsymbol{\eta}|\mathbf{x})} + \log \prod_{\zeta, \mathbf{l}} \frac{p(\zeta|\mathbf{v})}{q(\zeta|\mathbf{x})} \frac{p(\mathbf{l}|\boldsymbol{\eta})}{q(\mathbf{l}|\mathbf{x})} \frac{p(\mathbf{z}|\boldsymbol{\mu}_{\zeta, \mathbf{l}}, \boldsymbol{\sigma}_{\zeta, \mathbf{l}}^2)}{q(\mathbf{z}|\mathbf{x})} + \log p(\mathbf{x}|\mathbf{z}) \right] \quad (5)$$

3.4 TRAINING ALGORITHM OF CLUSTERING HIERARCHY

This model is formalized according to the stick-breaking process scheme. Unlike the CRP, the stick-breaking process does not represent the direct sampling of the mixture component at the data instance level. Therefore, it is necessary to devise a heuristic algorithm for operations, such as *GROW*, *PRUNE*, and *MERGE*, to refine the hierarchy structure. Appendix C provides details about each operation, together with the overall training algorithm of HCRL. In the below description, an *inner* path and a *full* path refer to the path ending with an internal node and a leaf node, respectively.

- **GROW** expands the hierarchy by creating a new branch under the heavily weighted internal node. Compared with the work of Wang & Blei (2009), we modified GROW to first sample a path, $\bar{\zeta}^*$, proportional to $\sum_n q(\zeta_n = \bar{\zeta}^*)$, and then to grow the path if the sampled path is an inner path.

- **PRUNE** cuts a randomly sampled minor full path, $\bar{\zeta}^*$, satisfying $\frac{\sum_n q(\zeta_n = \bar{\zeta}^*)}{\sum_n \bar{\zeta} q(\zeta_n = \bar{\zeta})} < \delta$, where δ is the pre-defined threshold. If the removed leaf node of the full path is the last child of the parent node, we also recursively remove the parent node.
- **MERGE** combines two full paths, $\bar{\zeta}^{(i)}$ and $\bar{\zeta}^{(j)}$, with similar posterior probabilities, measured by $J(\bar{\zeta}^{(i)}, \bar{\zeta}^{(j)}) = \mathbf{q}_i \mathbf{q}_j^T / |\mathbf{q}_i| |\mathbf{q}_j|$, where $\mathbf{q}_i = [q(\zeta_1 = \bar{\zeta}^{(i)}), \dots, q(\zeta_N = \bar{\zeta}^{(i)})]$.

4 EXPERIMENTS

4.1 DATASETS AND BASELINES

Datasets: We used various hierarchically organized benchmark datasets as well as MNIST.

- **MNIST (LeCun et al., 1998):** 28x28x1 handwritten image data, with 60,000 train images and 10,000 test images. We reshaped the data to 784-d in one dimension.
- **CIFAR-100 (Krizhevsky & Hinton, 2009):** 32x32x3 colored images with 20 coarse and 100 fine classes. We used 3,072-d flattened data with 50,000 training and 10,000 testing.
- **RCV1_v2 (Lewis et al., 2004):** The preprocessed text of the Reuters Corpus Volume. We preprocessed the text by selecting the top 2,000 tf-idf words. We used the hierarchical labels up to the 4-level, and the multi-labeled documents were removed. The final preprocessed corpus consists of 11,370 training and 10,000 testing documents randomly sampled from the original test corpus.
- **20Newsgroups (Lang, 1995):** The benchmark text data extracted from 20 newsgroups, consisting 11,314 training and 7,532 testing documents. We also labeled by 4-level following the annotated hierarchical structure. We preprocessed the data through the same process as that of RCV1_v2.

Baselines: We completed our evaluation in two aspects: 1) optimizing the density estimation, and 2) clustering the hierarchical categories. First, we evaluated HCRL from the density estimation perspective by comparing it with diverse flat clustered representation learning models, and VAE-nCRP. Second, we tested HCRL from the accuracy perspective by comparing it with multiple divisive hierarchical clusterings. The below is the list of baselines. We also added the two-stage pipeline approaches, where we trained features from VaDE first and then applied the hierarchical clusterings. We reused the open source codes³ provided by the authors for several baselines, such as IDEC, DCN, VAE-nCRP, and SSC-OMP.

1. **Variational Autoencoder (VAE) (Kingma & Welling, 2014)**
2. **Variational Deep Embedding (VaDE) (Jiang et al., 2017)**
3. **Improved Deep Embedded Clustering (IDEC) (Guo et al., 2017):** improves DEC (Xie et al., 2016) by attaching decoder structure. We use the code by the authors.
4. **Deep Clustering Network (DCN) (Yang et al., 2017):** optimizes the K-means-related cost defined on the embedding space. We used the open source code provided by the authors.
5. **Infinite Mixture of Variational Autoencoders (IMVAE) (Abbasnejad et al., 2017):** searches for the infinite embedding space by using a Bayesian nonparametric prior.
6. **Variational Autoencoder - nested Chinese Restaurant Process (VAE-nCRP) (Goyal et al., 2017):** We used the open source code provided by the authors.
7. **Hierarchical K-means (HKM) (Nister & Stewenius, 2006):** performs K-means (Lloyd, 1982) recursive in a top-down way.
8. **Mixture of Hierarchical Gaussians (MOHG) (Vasconcelos & Lippman, 1999):** infers the level-specific mixture of Gaussians.
9. **Recursive Gaussian Mixture Model (RGMM):** runs GMM recursively in a top-down manner.
10. **Recursive Scalable Sparse Subspace Clustering by Orthogonal Matching Pursuit (RSS-COMP):** performs SSC-OMP (You et al., 2016) recursively for hierarchical clustering. SSC-OMP is a well-known methods for image clustering, and we used the open source code.

4.2 QUANTITATIVE ANALYSIS

We used two measures to evaluate the learned representations in terms of the density estimations: 1) negative log likelihood (NLL), and 2) reconstruction errors (REs). Autoencoder models, such as

³<https://github.com/XifengGuo/IDEC> (IDEC); <https://github.com/boyangumn/DCN> (DCN); <https://github.com/prasoongoyal/bnp-vae> (VAE-nCRP); <http://vision.jhu.edu/code/> (SSC-OMP)

IDEC and DCN, were tested only for the REs. The NLL is estimated with 100 samples. Table 2 indicates that HCRL is best in the NLL and is competent in the REs which means that the hierarchically clustered embeddings preserve the intrinsic raw data structure.

Table 2: Test set performance of the negative log likelihood (NLL) and the reconstruction errors (REs). Replicated ten times, and the best in bold. $P^\dagger < 0.05$ (Student’s t-test). *Model-L#* means that the model trained with the #-depth hierarchy.

Model	MNIST		CIFAR-100		RCV1_v2		20Newsgroups	
	NLL	REs	NLL	REs	NLL	REs	NLL	REs
VAE	230.71	10.46	1960.06	57.54	2559.46	1434.59	2735.80	1788.22
VaDE	217.20	10.35	1921.85	53.60	2558.32	1426.38	2733.46	1782.86
IDEC	N/A	12.75	N/A	64.09	N/A	1376.26	N/A	1660.61 [†]
DCN	N/A	11.30	N/A	44.26	N/A	1361.98	N/A	1691.17
IMVAE	296.57	10.69	1992.83	40.45 [†]	2566.01	1387.02	2722.81	1718.08
VAE-nCRP-L3	718.78	32.67	2969.62	198.66	2642.88	1538.42	2712.28	1680.56
VAE-nCRP-L4	721.00	32.53	2950.73	198.97	2646.48	1542.81	2713.58	1680.71
HCRL-L3	203.24 [†]	8.70 [†]	1843.40 [†]	50.44	2554.50 [†]	1395.05	2726.75	1828.71
HCRL-L4	203.91 [†]	8.16 [†]	1849.13 [†]	50.47	2535.43 [†]	1353.34	2702.88	1711.30

VaDE generally performed better than VAE did, whereas other flat clustered representation learning models tended to be slightly different for each dataset. HCRL showed overall competent performance and better results with a deeper hierarchy of level four than of level three, which implies that capturing the deeper hierarchical structure is likely to be useful for the density estimation.

Additionally, we evaluated hierarchical clustering accuracies by following Xie et al. (2016), except for MNIST that is flat structured. Table 3 points out that HCRL has significantly better micro-averaged F-scores compared with every baseline. HCRL is able to reproduce the ground truth hierarchical structure of the data, and this trend is consistent when HCRL compared with the pipelined model, such as VaDE with a clustering model. The result of the comparisons with the clustering models, such as HKM, MOHG, RGMM, and RSSCOMP, is interesting because it experimentally proves that the joint optimization of hierarchical clustering in the embedding space improves hierarchical clustering accuracies. HCRL also presented better hierarchical accuracies than VAE-nCRP. We conjecture the reasons for the modeling aspect of VAE-nCRP: 1) the simplified prior modeling on the variance of the mixture component as just constants, and 2) the non-flexible learning of the internal components.

Table 3: Hierarchical clustering accuracies with F-scores, on CIFAR-100 with a depth of three, RCV1_v2 with a depth of four, and 20Newsgroups with a depth of four. Replicated ten times, and a confidence interval with 95%. Best in bold.

Model	CIFAR-100	RCV1_v2	20Newsgroups
HKM	0.1620 \pm 0.0077	0.2564 \pm 0.0679	0.4088 \pm 0.0426
MOHG	0.0846 \pm 0.0378	0.1026 \pm 0.0135	0.0402 \pm 0.0119
RGMM	0.1686 \pm 0.0115	0.2743 \pm 0.0521	0.4351 \pm 0.0369
RSSCOMP	0.1461 \pm 0.0228	0.2657 \pm 0.0545	0.2953 \pm 0.0474
VAE-nCRP	0.2011 \pm 0.0076	0.4128 \pm 0.0242	0.5584 \pm 0.0267
VaDE+HKM	0.1637 \pm 0.0116	0.3308 \pm 0.0664	0.4850 \pm 0.0558
VaDE+MOHG	0.1659 \pm 0.0155	0.4227 \pm 0.0927	0.4915 \pm 0.0713
VaDE+RGMM	0.1806 \pm 0.0132	0.3858 \pm 0.0615	0.4095 \pm 0.0651
VaDE+RSSCOMP	0.1923 \pm 0.0211	0.2718 \pm 0.0444	0.2905 \pm 0.0431
HCRL	0.2245 \pm 0.0137	0.4553 \pm 0.0295	0.6008 \pm 0.0973

4.3 QUALITATIVE ANALYSIS

MNIST: In Figure 1, the digits {4, 7, 9} and the digits {3, 8} are grouped together with a clear hierarchy, which was consistent between HCRL and VaDE. Also, some digits {0, 4, 2} in a round form are grouped, together, in HCRL. In addition, among the reconstructed digits from the hierarchical mixture components, the digits generated from the root have blended shapes from 0 to 9, which is natural considering the root position.

CIFAR-100: Figure 4 shows the hierarchical clustering results on CIFAR-100. Given that there were no semantic inputs from the data, the color was dominantly reflected in the clustering criteria. However, if one observes the second hierarchy, the scene images of the same sub-hierarchy are semantically consistent, although the background colors are slightly different.

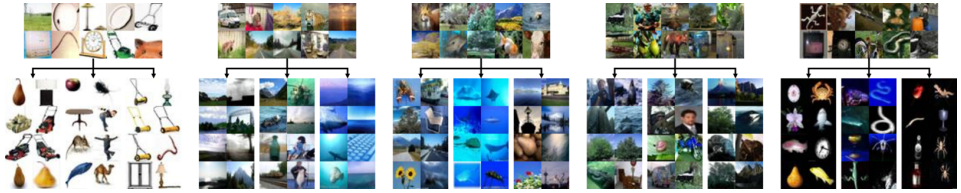


Figure 4: Example extracted sub-hierarchies on CIFAR-100

RCV1_v2: Figure 5 shows the embedding of RCV1_v2. VAE and VaDE show no hierarchy, and close sub-hierarchies are distantly embedded. VAE-nCRP guides the internal mixture components to be agglomerated at the center, and the cause of agglomeration is the generative process of VAE-nCRP, where the parameter of the internal components are inferred without direct information from data. HCRL shows a clear separation between the sub-hierarchy without the agglomeration.

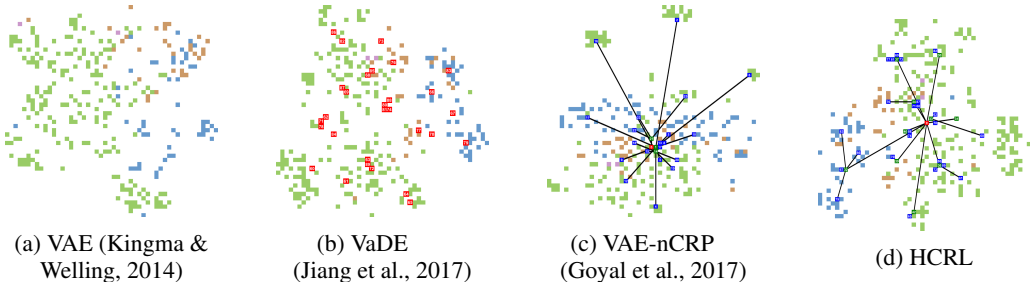


Figure 5: Comparison of embeddings on RCV1_v2, plotted using t-SNE (Maaten & Hinton, 2008). We mark the mean of a mixture component with a numbered square, colored in {red} for VaDE, {red (root), green (internal), blue (leaf)} for VAE-nCRP and HCRL. The first-level sub-hierarchies are indicated with four colors.

20Newsgroups: Figure 6 shows the example sub-hierarchies on 20Newsgroups. We enumerated topic words from documents with top-five likelihoods for each cluster, and we filtered the words by tf-idf values. We observe relatively more general contents in the internal clusters than in the leaf clusters of each internal cluster.

<p>1st subhierarchy on ‘computer’ bios pictures picture hardware screen ___ brand sesi drive connect computers ___ floppy interface transfer words hd ___ mac moving floppy screen black</p>	<p>2nd subhierarchy on ‘politics’ topic movement war noise majority ___ court criminals law crypto encryption ___ claims investigation percent adam wall ___ party conflict industry majority unit</p>	<p>3rd subhierarchy on ‘vehicles&sports’ series vehicle day finally afraid ___ hit bunch rule playing station ___ excellent hospital insurance game pick ___ pictures teams guys shot daily</p>
--	--	---

Figure 6: Example extracted sub-hierarchies on 20Newsgroups

5 CONCLUSION

In this paper, we have introduced a hierarchically clustered representation learning framework for the hierarchical mixture density estimation on deep embeddings. HCRL aims at encoding the relations among clusters as well as among instances to preserve the internal hierarchical structure of data. The main differentiated features of HCRL are 1) the crucial assumption regarding the internal mixture components for having the ability to generate data directly, and 2) the unbalanced autoencoding neural architecture for the level proportion modeling as the encoding structure, and the probabilistic model as the decoding structure. From the modeling and the evaluation, we found that HCRL enables the improvements due to the high flexibility modeling compared with the baselines.

REFERENCES

- M Ehsan Abbasnejad, Anthony Dick, and Anton van den Hengel. Infinite variational autoencoder for semi-supervised learning. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 781–790. IEEE, 2017.
- 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.
- David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.
- Xi Chen, Yan Duan, Rein Houthoofd, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *Advances in neural information processing systems*, pp. 2172–2180, 2016.
- Wenqing Chu and Deng Cai. Stacked similarity-aware autoencoders. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pp. 1561–1567. AAAI Press, 2017.
- Prasoon Goyal, Zhiting Hu, Xiaodan Liang, Chenyu Wang, and Eric Xing. Nonparametric variational auto-encoders for hierarchical representation learning. *arXiv preprint arXiv:1703.07027*, 2017.
- Thomas L Griffiths, Michael I Jordan, Joshua B Tenenbaum, and David M Blei. Hierarchical topic models and the nested chinese restaurant process. In *Advances in neural information processing systems*, pp. 17–24, 2004.
- Xifeng Guo, Long Gao, Xinwang Liu, and Jianping Yin. Improved deep embedded clustering with local structure preservation. In *International Joint Conference on Artificial Intelligence (IJCAI-17)*, pp. 1753–1759, 2017.
- Philipp Hennig, David Stern, Ralf Herbrich, and Thore Graepel. Kernel topic models. In *Artificial Intelligence and Statistics*, pp. 511–519, 2012.
- Peihao Huang, Yan Huang, Wei Wang, and Liang Wang. Deep embedding network for clustering. In *Pattern Recognition (ICPR), 2014 22nd International Conference on*, pp. 1532–1537. IEEE, 2014.
- Zhuxi Jiang, Yin Zheng, Huachun Tan, Bangsheng Tang, and Hanning Zhou. Variational deep embedding: An unsupervised and generative approach to clustering. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pp. 1965–1972. AAAI Press, 2017.
- 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.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *Proceedings of the Second International Conference on Learning Representations (ICLR 2014)*, April 2014.
- Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. 2009.
- Ken Lang. Newsweeder: Learning to filter netnews. In *Machine Learning Proceedings 1995*, pp. 331–339. Elsevier, 1995.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- David D Lewis, Yiming Yang, Tony G Rose, and Fan Li. Rcv1: A new benchmark collection for text categorization research. *Journal of machine learning research*, 5(Apr):361–397, 2004.
- Stuart Lloyd. Least squares quantization in pcm. *IEEE transactions on information theory*, 28(2): 129–137, 1982.

- Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605, 2008.
- David JC MacKay. Choice of basis for laplace approximation. *Machine learning*, 33(1):77–86, 1998.
- David Mimno, Wei Li, and Andrew McCallum. Mixtures of hierarchical topics with pachinko allocation. In *Proceedings of the 24th international conference on Machine learning*, pp. 633–640. ACM, 2007.
- Eric Nalisnick, Lars Hertel, and Padhraic Smyth. Approximate inference for deep latent gaussian mixtures. In *NIPS Workshop on Bayesian Deep Learning*, volume 2, 2016.
- David Nister and Henrik Stewenius. Scalable recognition with a vocabulary tree. In *Computer vision and pattern recognition, 2006 IEEE computer society conference on*, volume 2, pp. 2161–2168. Ieee, 2006.
- Jim Pitman et al. Combinatorial stochastic processes. Technical report, Technical Report 621, Dept. Statistics, UC Berkeley, 2002. Lecture notes for St. Flour course, 2002.
- Abel Rodriguez, David B Dunson, and Alan E Gelfand. The nested dirichlet process. *Journal of the American Statistical Association*, 103(483):1131–1154, 2008.
- David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning internal representations by error propagation. Technical report, California Univ San Diego La Jolla Inst for Cognitive Science, 1985.
- Akash Srivastava and Charles Sutton. Autoencoding variational inference for topic models. *arXiv preprint arXiv:1703.01488*, 2017.
- Naonori Ueda, Ryohei Nakano, Zoubin Ghahramani, and Geoffrey E Hinton. Smem algorithm for mixture models. In *Advances in neural information processing systems*, pp. 599–605, 1999.
- Nuno Vasconcelos and Andrew Lippman. Learning mixture hierarchies. In *Advances in Neural Information Processing Systems*, pp. 606–612, 1999.
- Chong Wang and David M Blei. Variational inference for the nested chinese restaurant process. In *Advances in Neural Information Processing Systems*, pp. 1990–1998, 2009.
- Junyuan Xie, Ross Girshick, and Ali Farhadi. Unsupervised deep embedding for clustering analysis. In *International conference on machine learning*, pp. 478–487, 2016.
- Bo Yang, Xiao Fu, Nicholas D Sidiropoulos, and Mingyi Hong. Towards k-means-friendly spaces: Simultaneous deep learning and clustering. In *International Conference on Machine Learning*, pp. 3861–3870, 2017.
- Chong You, Daniel Robinson, and René Vidal. Scalable sparse subspace clustering by orthogonal matching pursuit. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3918–3927, 2016.

A SYNTHETIC DEMO

We created a synthetic dataset that has a hierarchical structure and is sampled from the 50-dimensional Gaussian distributions, presented in Figure 7. The hierarchy, which has a branch factor of two and a depth of four, has a total of eight leaf clusters. Figure 7a shows the raw synthetic dataset in the input space of \mathbb{R}^{50} , and after running HCRL, we plot the hierarchically clustered embeddings in the latent space in Figure 7b. In addition to the embeddings, we also present a confidence ellipse with dashed lines for each learned Gaussian mixture component. Because the root component is involved in generating all of the data, it forms a large ellipse, while the leaf component summarizes the local density, so the small ellipse is learned.

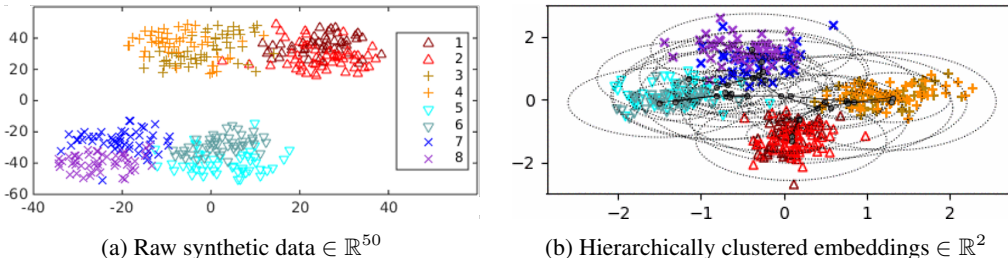


Figure 7: Synthetic data in the input space of \mathbb{R}^{50} (left), which is visualized via t-SNE (Maaten & Hinton, 2008), and hierarchically clustered embeddings in the latent space of \mathbb{R}^2 (right). We additionally show a 95% confidence ellipse with a dashed line for each Gaussian mixture component.

We show how the above embeddings learned to be hierarchically clustered in the latent space during training in Figure 8. In the learning mechanism of HCRL, we can observe the hierarchically clustered embeddings from a major deviation to a minor deviation in the data over iterations.

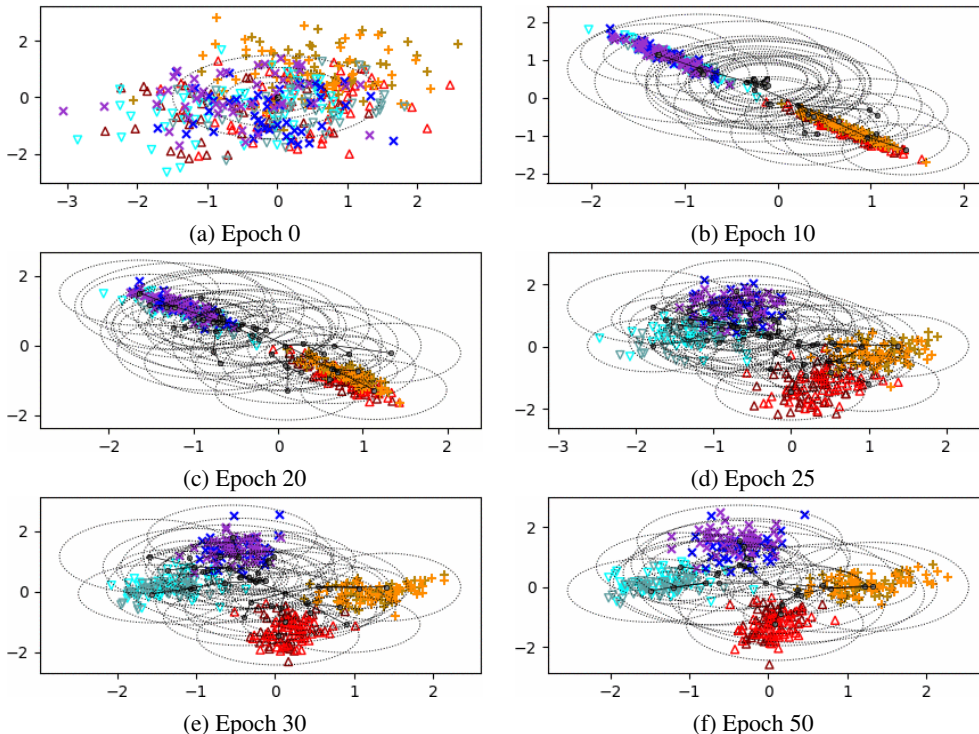


Figure 8: The process by which the embeddings of the synthetic data are learned. The dashed ellipse corresponds to the 95% contour of the learned Gaussian mixture component, whose mean is marked as the gray circle.

B EXPERIMENTAL SETTINGS

We conducted experiments for all autoencoder-based models with a neural architecture whose encoder network was set as fully connected layers with dimensions D -2000-2000-500- J for \mathbf{z} , and D -10-10- L for $\boldsymbol{\eta}$, and the decoder network is a mirror of the encoder network for \mathbf{z} . The hyperparameters of HCRL given by users, γ and α , was set to 1.0, and a vector of all entries 1 sized of L , respectively. We used the Adam optimizer (Kingma & Ba, 2014) with an init learning rate of 0.001 for MNIST dataset and 0.0001 for other datasets. Meanwhile, VAE-nCRP is targeted for grouped data. For experiments with our non-grouped datasets, we treated the group instance as a group instance having a single data instance. For parametric hierarchical clustering models, we gave the branch factor as the input parameter, [1, 20, 5], [1, 4, 7, 9], and [1, 6, 4, 3], for CIFAR-100, RCV1_v2, and 20Newsgroups, respectively. For VaDE, we set the number of clusters to the number of leaf clusters; 100 for CIFAR-100, 252 for RCV1_v2, and 72 for 20Newsgroups.

C ALGORITHMS

C.1 TRAINING ALGORITHM

Algorithm 1 summarizes the overall algorithm for HCRL. The tree-based hierarchy T is defined as (\mathbb{N}, \mathbb{P}) , where \mathbb{N} and \mathbb{P} denote a set of nodes and paths, respectively. We refer to the node at level l lying on path ζ , as $N(\zeta_{1:l}) \in \mathbb{N}$. The defined paths, \mathbb{P} , consist of full paths (ending at a leaf node), \mathbb{P}_{full} , and inner paths (ending at an internal node), $\mathbb{P}_{\text{inner}}$, as a union set.

Algorithm 1 selects an operation out of three operations: *GROW*, *PRUNE*, and *MERGE*. The *GROW* algorithm is executed for every specific iteration period, t_G . After ellapsing t_b iterations since performing the *GROW* operation, we begin to check whether the *PRUNE* or *MERGE* operation should be performed. We prioritize the *PRUNE* operation first, and if the condition of performing *PRUNE* is not satisfied, we check for the *MERGE* operation next. After performing any operation, we initialize n_b to 0, which is for locking the changed hierarchy during minimum t_b iterations to be fitted to the training data.

Algorithm 1 Training for Hierarchically Clustered Representation Learning

Input: Training examples \mathbf{x} ; the tree-based hierarchy depth, L ; period of performing *GROW*, t_{grow} ; minimum number of epochs locking the hierarchy, t_{lock} ; operation-related thresholds δ_{prune} , δ_{merge} ; a queue whose element is the set of changed paths, \mathbb{Q} ; the number of training epochs, E ; maximum length of \mathbb{Q} , Q_{max} ; grow scale, s_{grow}

Output: $T^{(E)}$, $\phi_{\mathbf{z}}$, $\phi_{\boldsymbol{\eta}}$, $\boldsymbol{\theta}$, $\boldsymbol{\omega}$, $\{a_i, b_i, \boldsymbol{\mu}_i, \boldsymbol{\sigma}_i^2\}_{i \in \mathcal{M}_{T^{(E)}}$

- 1: $\boldsymbol{\mu}_{\bar{\zeta}_{1:L}}, \boldsymbol{\sigma}_{\bar{\zeta}_{1:L}}^2 \leftarrow$ Initialize L Gaussian mixture components
- 2: $T^{(0)} \leftarrow$ Initialize the tree-based hierarchy having a single path with $\boldsymbol{\mu}_{\bar{\zeta}_{1:L}}, \boldsymbol{\sigma}_{\bar{\zeta}_{1:L}}^2$
- 3: $n_{\text{lock}} \leftarrow 0$ // for counting the number of epochs, where the hierarchy has not changed
- 4: **for** each epoch $e = 1, \dots, E$ **do**
- 5: $\phi_{\mathbf{z}}, \phi_{\boldsymbol{\eta}}, \boldsymbol{\theta} \leftarrow$ Update the network weight parameters using gradients $\nabla_{\phi_{\mathbf{z}}, \phi_{\boldsymbol{\eta}}, \boldsymbol{\theta}} \mathcal{L}(\mathbf{x})$
- 6: $\{a_i, b_i, \boldsymbol{\mu}_i, \boldsymbol{\sigma}_i^2\}_{i \in \mathcal{M}_{T^{(e-1)}}} \leftarrow$ Update node-specific params. using gradients $\nabla_{a, b, \boldsymbol{\mu}, \boldsymbol{\sigma}^2} \mathcal{L}(\mathbf{x})$
- 7: Update other variational parameters using gradients $\nabla \mathcal{L}(\mathbf{x})$
- 8: **if** $\text{mod}(e, t_{\text{grow}}) = 0$ **then**
- 9: $T^{(e)}, \mathbb{Q} \leftarrow$ *GROW*($T^{(e-1)}, \mathbb{Q}, s_{\text{grow}}, Q_{\text{max}}$) // See Algorithm 2
- 10: **end if**
- 11: **if** $T^{(e)} = T^{(e-1)}$ and $n_{\text{lock}} \geq t_{\text{lock}}$ **then**
- 12: $T^{(e)}, \mathbb{Q} \leftarrow$ *PRUNE*($T^{(e-1)}, \mathbb{Q}, \delta_{\text{prune}}$) // See Algorithm 3
- 13: **if** $T^{(e)} = T^{(e-1)}$ **then** $T^{(e)}, \mathbb{Q} \leftarrow$ *MERGE*($T^{(e-1)}, \mathbb{Q}, \delta_{\text{merge}}, Q_{\text{max}}$) // See Algorithm 4
- 14: **end if**
- 15: **if** $T^{(e)} \neq T^{(e-1)}$ **then** $n_{\text{lock}} \leftarrow 0$ **else** $n_{\text{lock}} \leftarrow n_{\text{lock}} + 1$
- 16: **end for**

C.2 ALGORITHM FOR GROW OPERATION

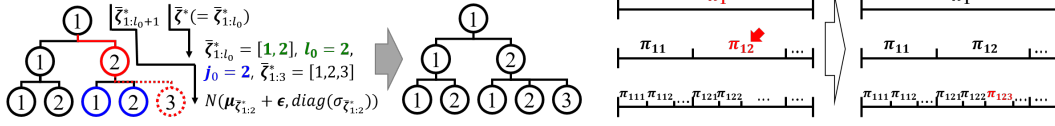


Figure 9: The illustration of GROW operation

The GROW operation expands the hierarchy by creating a new branch under the heavily weighted internal node. Compared with the work from Wang & Blei (2009), we modify GROW to firstly sample a path according to $\sum_{n=1}^N q(\zeta_n = \bar{\zeta})$, and then grow the path if the sampled path is an inner path. When we create the new Gaussian mixture component, we initialize the parameters of a corresponding Gaussian distribution depending on the mean and the variance of the parent node, as shown in line 10 of Algorithm 2.

Algorithm 2 GROW Operation

```

1: function GROW( $T, \mathbb{Q}, s_{\text{grow}}, Q_{\text{max}}$ )
2:    $J(\bar{\zeta}) \leftarrow \sum_{n=1}^N q(\zeta_n = \bar{\zeta})$  for  $\bar{\zeta} \in \mathbb{P}$  // Calculate the measure
3:   Sample a path  $\bar{\zeta}^*$  with probability  $\frac{J(\bar{\zeta}^*)}{\sum_{\bar{\zeta}} J(\bar{\zeta})}$ 
4:    $Q' \leftarrow \emptyset$  // Temporary set of changed paths in this epoch
5:   if  $\bar{\zeta}^* \in \mathbb{P}_{\text{inner}}$  and  $\bar{\zeta}^* \notin Q$  s.t.  $Q \in \mathbb{Q}$  then
6:      $l_0 \leftarrow |\bar{\zeta}^*|$ 
7:     for  $l' = l_0, \dots, L - 1$  do
8:        $j_0 \leftarrow$  Maximum index for the child node whose parent path is  $\bar{\zeta}_{1:l'}^*$ 
9:        $\bar{\zeta}_{1:l'+1}^* \leftarrow [\bar{\zeta}_{1:l'}^*, j_0 + 1]$ 
10:       $N(\bar{\zeta}_{1:l'+1}^*) \leftarrow \mathcal{N}(\mu_{\bar{\zeta}_{1:l'}^*} + \epsilon, \text{diag}(\sigma_{\bar{\zeta}_{1:l'}^*}^2))$  where  $\epsilon \sim \mathcal{N}(\mathbf{0}, n_g \mathbf{I}_J)$ 
11:       $Q' \leftarrow Q' \cup \{\bar{\zeta}_{1:l'+1}^*\}$ 
12:      if  $l' < L - 1$  then
13:         $\mathbb{P}_{\text{inner}} \leftarrow \mathbb{P}_{\text{inner}} \cup \{\bar{\zeta}_{1:l'+1}^*\}$ 
14:      else
15:         $\mathbb{P}_{\text{full}} \leftarrow \mathbb{P}_{\text{full}} \cup \{\bar{\zeta}_{1:L}^*\}$ 
16:      end if
17:    end for
18:  end if
19:  enqueue  $Q'$  to  $\mathbb{Q}$ 
20:  while  $Q_{\text{max}} < |\mathbb{Q}|$  do dequeue  $\mathbb{Q}$ 
21:   $\mathbb{P} \leftarrow \mathbb{P}_{\text{full}} \cup \mathbb{P}_{\text{inner}}$ 
22:   $T \leftarrow (\mathbb{N}, \mathbb{P})$ 
23:  return  $T, \mathbb{Q}$ 
24: end function

```

C.4 ALGORITHM FOR MERGE OPERATION

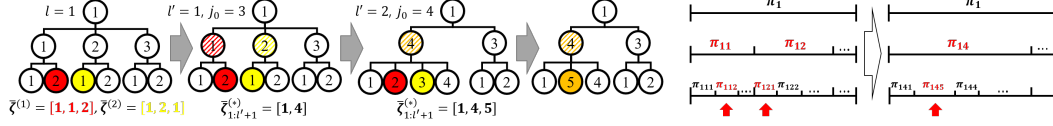


Figure 11: The illustration of MERGE operation

The MERGE operation combines two full paths with similar posterior probabilities, measured by $J(\bar{\zeta}^{(i)}, \bar{\zeta}^{(j)}) = \mathbf{q}_i \mathbf{q}_j^T / |\mathbf{q}_i| |\mathbf{q}_j|$, where $\mathbf{q}_i = [q(\zeta_1 = \bar{\zeta}^{(i)}), \dots, q(\zeta_N = \bar{\zeta}^{(i)})]$. We merged two Gaussian components by following Ueda et al. (1999). The specific meaning of combining the two paths is merging the paired two Gaussian distributions lying on the two paths by level, if the two Gaussian distributions are different. The estimation of merged Gaussian parameters, $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}$, is the weighted summation of two subject Gaussian parameters. The probability of the node at level l lying on a path ζ given \mathbf{x} , $p(\zeta_l | \mathbf{x})$, is proportional to $\sum_n \{q(l_n = l) \cdot \sum_{\zeta \in \Lambda} q(\zeta_n = \zeta)\}$, where $\Lambda = \{\zeta' | \zeta'_l = \zeta_l \text{ and } \zeta' \in \mathbb{P}_{\text{full}}\}$.

Algorithm 4 MERGE Operation

```

1: function MERGE( $T, \mathbb{Q}, \delta_{\text{merge}}, Q_{\text{max}}$ )
2:    $J(\bar{\zeta}^{(i)}, \bar{\zeta}^{(j)}) \leftarrow \frac{\mathbf{q}_i \mathbf{q}_j^T}{|\mathbf{q}_i| |\mathbf{q}_j|}$  s.t.  $\mathbf{q}_i = [q(\zeta_1 = \bar{\zeta}^{(i)}), \dots, q(\zeta_N = \bar{\zeta}^{(i)})]$  // Calculate the measure
3:    $\Omega \leftarrow \{(\bar{\zeta}^{(i)}, \bar{\zeta}^{(j)}) \mid J(\bar{\zeta}^{(i)}, \bar{\zeta}^{(j)}) \geq \delta_{\text{merge}}, \{\bar{\zeta}^{(i)}, \bar{\zeta}^{(j)}\} \subset \mathbb{P}_{\text{full}}\}$ 
4:   Randomly sample a pair of paths  $(\bar{\zeta}^{(1)}, \bar{\zeta}^{(2)}) \sim \Omega$ 
5:    $Q' \leftarrow \phi$  // Temporary set of changed paths in this epoch
6:   if  $\{(\bar{\zeta}^{(1)}, \bar{\zeta}^{(2)})\} \not\subseteq Q$  s.t.  $Q \in \mathbb{Q}$  then
7:      $l \leftarrow$  Maximum level of nodes shared by  $\bar{\zeta}^{(1)}, \bar{\zeta}^{(2)}$ 
8:      $\bar{\zeta}_{1:l}^* \leftarrow \bar{\zeta}_{1:l}^{(1)}$ 
9:     for  $l' = l, \dots, L - 1$  do
10:       $\boldsymbol{\mu}_{(1)} \leftarrow \boldsymbol{\mu}_{\bar{\zeta}_{1:l'+1}^{(1)}}$ ,  $\boldsymbol{\sigma}_{(1)}^2 \leftarrow \boldsymbol{\sigma}_{\bar{\zeta}_{1:l'+1}^{(1)}}^2$ ,  $\boldsymbol{\mu}_{(2)} \leftarrow \boldsymbol{\mu}_{\bar{\zeta}_{1:l'+1}^{(2)}}$ ,  $\boldsymbol{\sigma}_{(2)}^2 \leftarrow \boldsymbol{\sigma}_{\bar{\zeta}_{1:l'+1}^{(2)}}^2$ 
11:       $w_{(1)} \leftarrow p(\bar{\zeta}_{l'+1}^{(1)} | \mathbf{x})$ ,  $w_{(2)} \leftarrow p(\bar{\zeta}_{l'+1}^{(2)} | \mathbf{x})$ 
12:       $\boldsymbol{\mu}_* \leftarrow \frac{w_{(1)} \boldsymbol{\mu}_{(1)} + w_{(2)} \boldsymbol{\mu}_{(2)}}{w_{(1)} + w_{(2)}}$ ,  $\boldsymbol{\sigma}_*^2 \leftarrow \frac{w_{(1)} \boldsymbol{\sigma}_{(1)}^2 + w_{(2)} \boldsymbol{\sigma}_{(2)}^2}{w_{(1)} + w_{(2)}}$ 
13:       $j_0 \leftarrow$  Maximum index for the child node whose parent path is  $\bar{\zeta}_{1:l'}^*$ 
14:       $\bar{\zeta}_{1:l'+1}^* \leftarrow [\bar{\zeta}_{1:l'}^*, j_0 + 1]$ 
15:       $\mathcal{N}(\bar{\zeta}_{1:l'+1}^*) \leftarrow \mathcal{N}(\boldsymbol{\mu}_*, \text{diag}(\boldsymbol{\sigma}_*^2))$ 
16:       $\mathcal{N}(\bar{\zeta}_{1:l'+1}^{(1)}) \leftarrow \phi$ ,  $\mathcal{N}(\bar{\zeta}_{1:l'+1}^{(2)}) \leftarrow \phi$ 
17:      if  $l' < L - 1$  then
18:         $\mathbb{P}_{\text{inner}} \leftarrow \mathbb{P}_{\text{inner}} \cup \{\bar{\zeta}_{1:l'+1}^* \setminus \{\bar{\zeta}_{1:l'+1}^{(1)}, \bar{\zeta}_{1:l'+1}^{(2)}\}\}$ 
19:      else
20:         $\mathbb{P}_{\text{full}} \leftarrow \mathbb{P}_{\text{full}} \cup \{\bar{\zeta}_{1:L}^* \setminus \{\bar{\zeta}_{1:L}^{(1)}, \bar{\zeta}_{1:L}^{(2)}\}\}$ 
21:      end if
22:       $Q' \leftarrow Q' \cup \{\bar{\zeta}_{1:l'+1}^*\}$ 
23:    end for
24:  end if
25:  enqueue  $Q'$  to  $\mathbb{Q}$ 
26:  while  $Q_{\text{max}} < |\mathbb{Q}|$  do dequeue  $\mathbb{Q}$ 
27:   $\mathbb{P} \leftarrow \mathbb{P}_{\text{full}} \cup \mathbb{P}_{\text{inner}}$ 
28:   $T \leftarrow (\mathbb{N}, \mathbb{P})$ 
29:  return  $T, \mathbb{Q}$ 
30: end function

```

D NOTATIONS

The following Table 4 lists the notations used throughout this paper.

Table 4: Table of symbols

Models	Symbol	Definition	
All	\mathbf{x}/\mathbf{x}'	An observed / reconstructed datapoint	
	\mathbf{z}	A latent representation	
	D/J	The input / latent dimensionality	
	$g_*(\mathbf{x})$	A encoder network parametrized by $*$, whose input is \mathbf{x}	
	$f_\theta(\mathbf{z})$	A decoder network parametrized by θ , whose input is \mathbf{z}	
	θ	The variational parameters and weights of the decoder network f_θ	
	$\tilde{\mu}_z, \tilde{\sigma}_z^2$	The variational mean and variance for Gaussian distribution $q_{\phi_z}(\mathbf{z} \mathbf{x})$	
	μ_x, σ_x^2	The prior parameters, mean and variance, for Gaussian distribution $p_\theta(\mathbf{x} \mathbf{z})$	
	VaDE & VAE-nCRP	ϕ	The variational parameters and weights of the encoder network g_ϕ
		$\tilde{\mu}, \tilde{\sigma}^2$	The variational mean and variance for Gaussian distribution $q_\phi(\mathbf{z} \mathbf{x})$
VaDE & HCRL	N	The number of datapoints	
	$\mathbf{x}_{n=1, \dots, N}$	n -th observed datapoint	
	$\mathbf{z}_{n=1, \dots, N}$	n -th latent representation corresponding to \mathbf{x}_n	
VAE-nCRP & HCRL	L	The height of the tree-based hierarchy	
VaDE	K	The number of (finite) clusters	
	$c_{n=1, \dots, N}$	The cluster assignment of $\mathbf{z}_n, \in \{1, \dots, K\}$	
	κ	The prior parameter for multinomial distribution $p(\mathbf{c})$	
	μ_c, σ_c^2	The prior parameters, mean and variance, for Gaussian distribution of c -th cluster, $p(\mathbf{z})$	
VAE-nCRP	M	The number of sequences	
	$N_{m=1, \dots, M}$	The number datapoints in m -th sequence	
	$\mathbf{x}_{m, n=1, \dots, N}$	n -th observed datapoint in m -th sequence	
	$\mathbf{z}_{m, n=1, \dots, N}$	n -th latent representation corresponding to \mathbf{x}_{mn}	
	v_{mp}	The Beta draws of m -th sequence on node p , for the tree-based stick-breaking construction	
		γ^*	The prior parameter for Beta distribution $p(v_{mp})$
		$\gamma_{mp}^{(0)}, \gamma_{mp}^{(1)}$	The variational parameters, for Beta distribution $q(v_{mp} \mathbf{x}_m)$
		ζ_{mn}	The path assignment of \mathbf{z}_{mn}
		\mathcal{S}_{mn}^*	The variational parameter for multinomial distribution $q(\zeta_{mn} \mathbf{x}_{mn})$
		$\alpha_{\text{par}(p)}$	The J -dimensional parameter vector for the parent node of p
		α^*	The prior parameter for Gaussian distribution $p(\alpha_p)$ for the root node
		$\mu_{\text{par}(p)}, \sigma_{\text{par}(p)}^2$	The variational mean and variance for Gaussian distribution $q(\alpha_{\text{par}(p)} \mathbf{x})$
		α_p	The J -dimensional parameter vector for node p
		σ_N^2	The prior parameter, variance, for Gaussian distribution $p(\alpha_p \alpha_{\text{par}(p)})$
		μ_p, σ_p^2	The variational mean and variance for Gaussian distribution $q(\alpha_p \mathbf{x})$
		σ_D^2	The prior parameter, variance, for Gaussian distribution $p(\mathbf{z}_{mn} \zeta_{mn}, \alpha_p)$
	HCRL	ϕ_z	The variational parameters and weights of the encoder network g_{ϕ_z}
		ϕ_η	The variational parameters and weights of the encoder network g_{ϕ_η}
		$\tilde{\mu}_z, \tilde{\sigma}_z^2$	The variational mean and variance for Gaussian distribution $q_{\phi_z}(\mathbf{z} \mathbf{x})$
		$\tilde{\mu}_\eta, \tilde{\sigma}_\eta^2$	The variational mean and variance for logistic normal distribution $q_{\phi_\eta}(\boldsymbol{\eta} \mathbf{x})$
		$\tilde{\alpha}$	The variational parameter for Dirichlet distribution $q_{\phi_\eta}(\boldsymbol{\eta} \mathbf{x})$
		v_i	The Beta draws for the tree-based stick-breaking construction of node i
		γ	The prior parameter for Beta distribution $p(v_i)$
		a_i, b_i	The variational parameters, for Beta distribution $q(v_i \mathbf{x})$
		ζ_n	The path assignment of \mathbf{z}_n
		\mathcal{S}_n	The variational parameter for multinomial distribution $q(\zeta_n \mathbf{x}_n)$
	η_n	The level proportion of \mathbf{z}_n	
	α	The prior parameter for Dirichlet distribution $p(\boldsymbol{\eta}_n)$	
	l_n	The level assignment of $\mathbf{z}_n, \in \{1, \dots, L\}$	
	ω_n	The variational parameter for multinomial distribution $q(l_n \mathbf{x}_n)$	
	μ_i, σ_i^2	The prior parameters, mean and variance, for Gaussian distribution of node $i, p(\mathbf{z}_n \zeta_n, \boldsymbol{\eta}_n)$	

E GENERATIVE AND INFERENCE MODEL FOR HCRL

HCRL assumes the generative process as described in Section 3.1. Section E.1 describes the joint probability distribution, and Section E.2 presents the corresponding variational distributions. We adopt the much notation-related conventions from Wang & Blei (2009), especially on paths.

E.1 GENERATIVE MODEL

$$\begin{aligned} p_{\theta}(\mathbf{v}, \zeta, \boldsymbol{\eta}, \mathbf{l}, \mathbf{z}, \mathbf{x}) &= p(\mathbf{v}|\gamma) \prod_{n=1}^N p(\zeta_n|\mathbf{v})p(\boldsymbol{\eta}_n|\boldsymbol{\alpha})p(l_n|\boldsymbol{\eta}_n)p(\mathbf{z}_n|\zeta_n, l_n, \boldsymbol{\mu}_{1:\infty}, \boldsymbol{\sigma}_{1:\infty}^2)p_{\theta}(\mathbf{x}_n|\mathbf{z}_n) \\ &= \prod_{j \notin \mathcal{M}_T} p(v_j|\gamma) \prod_{i \in \mathcal{M}_T} p(v_i|\gamma) \prod_{n=1}^N p(\zeta_n|\mathbf{v})p(\boldsymbol{\eta}_n|\boldsymbol{\alpha})p(l_n|\boldsymbol{\eta}_n)p(\mathbf{z}_n|\zeta_n, l_n)p_{\theta}(\mathbf{x}_n|\mathbf{z}_n) \end{aligned}$$

- \mathcal{M}_T : Set of all nodes in truncated tree T
- For $j \notin \mathcal{M}_T$, $p(v_j|\gamma) = \text{Beta}(v_j|1, \gamma)$
- For $i \in \mathcal{M}_T$, $p(v_i|\gamma) = \text{Beta}(v_i|1, \gamma)$
- $p(\zeta_n = [1, \zeta_2, \dots, \zeta_L]|\mathbf{v})$

$$\begin{aligned} p(\zeta_n = [1, \zeta_2, \dots, \zeta_L]|\mathbf{v}) &= \prod_{l=1}^L \pi_{1, \zeta_2, \dots, \zeta_l} \\ &= \prod_{l=1}^L \pi_{1, \zeta_2, \dots, \zeta_{l-1}} v_{1, \zeta_2, \dots, \zeta_l} \prod_{j=1}^{\zeta_l-1} (1 - v_{1, \zeta_2, \dots, j}) \\ &= \prod_{l=1}^L \prod_{l'=1}^l \left\{ v_{1, \zeta_2, \dots, \zeta_{l'}} \left(\prod_{j=1}^{\zeta_{l'}-1} (1 - v_{1, \zeta_2, \dots, j}) \right) \right\} \end{aligned}$$

- $\pi_{1, \zeta_2, \dots, \zeta_l} = \prod_{l'=1}^l \{ v_{1, \zeta_2, \dots, \zeta_{l'}} (\prod_{j=1}^{\zeta_{l'}-1} (1 - v_{1, \zeta_2, \dots, j})) \}$
- $p(\boldsymbol{\eta}_n|\boldsymbol{\alpha}) = \text{Dirichlet}(\boldsymbol{\eta}_n|\boldsymbol{\alpha})$
- $p(l_n|\boldsymbol{\eta}_n) = \text{Multinomial}(\boldsymbol{\eta}_n)$
- $p(\mathbf{z}_n|\zeta_n = \zeta, l_n = l) = \mathcal{N}(\mathbf{z}_n|\boldsymbol{\mu}_{\zeta_l}, \boldsymbol{\sigma}_{\zeta_l}^2 \mathbf{I}_J)$
- $p_{\theta}(\mathbf{x}_n|\mathbf{z}_n)$: Probabilistic decoding of \mathbf{x}_n parametrized by θ , whose input is \mathbf{z}_n
- Tree-based stick-breaking construction
 - We will denote all Beta draws as \mathbf{v} , each of which is an independent draw from $\text{Beta}(\mathbf{v}|1, \gamma)$ (except for root $v_1 = 1$)
 - * $v_i \sim \text{Beta}(v_i|1, \gamma)$
 - The root nodes stick length: $\pi_1 = v_1 \equiv 1$
 - Stick length at second level: $\pi_{1i} = \pi_1 v_{1i} \prod_{j=1}^{i-1} (1 - v_{1j})$, $\sum_{i=1}^{\infty} \pi_{1i} = \pi_1 = 1$
 - For the segment π_{1k} , the stick lengths of its children are $\pi_{1ki} = \pi_{1k} v_{1ki} \prod_{j=1}^{i-1} (1 - v_{1kj})$, for $i = 1, 2, \dots, \infty$, $\sum_{i=1}^{\infty} \pi_{1ki} = \pi_{1k}$

E.2 INFERENCE MODEL

As VAE, we infer the random variables via the mean-field approximation, where the variational distribution, $q_{\phi_{\boldsymbol{\eta}}, \phi_{\mathbf{z}}}(\mathbf{v}, \zeta, \boldsymbol{\eta}, \mathbf{l}, \mathbf{z}|\mathbf{x})$, approximates the intractable posterior. We model the variational distributions as follows:

$$\begin{aligned} q_{\phi_{\boldsymbol{\eta}}, \phi_{\mathbf{z}}}(\mathbf{v}, \zeta, \boldsymbol{\eta}, \mathbf{l}, \mathbf{z}|\mathbf{x}) &= q(\mathbf{v}|\mathbf{a}, \mathbf{b}, \mathbf{x}) \prod_{n=1}^N q(\zeta_n|\mathbf{x}_n)q_{\phi_{\boldsymbol{\eta}}}(\boldsymbol{\eta}_n|\mathbf{x}_n)q(l_n|\boldsymbol{\omega}_n, \mathbf{x}_n)q_{\phi_{\mathbf{z}}}(\mathbf{z}_n|\mathbf{x}_n) \\ &= \prod_{j \notin \mathcal{M}_T} p(v_j|\gamma) \prod_{i \in \mathcal{M}_T} q(v_i|a_i, b_i) \prod_{n=1}^N q(\zeta_n|\mathbf{x}_n)q_{\phi_{\boldsymbol{\eta}}}(\boldsymbol{\eta}_n|\mathbf{x}_n)q(l_n|\boldsymbol{\omega}_n, \mathbf{x}_n)q_{\phi_{\mathbf{z}}}(\mathbf{z}_n|\mathbf{x}_n) \end{aligned}$$

- For $j \notin \mathcal{M}_T$, $p(v_j|\gamma) = \text{Beta}(v_j|1, \gamma)$
- For $i \in \mathcal{M}_T$, $q(v_i|a_i, b_i) \propto v_i^{a_i-1}(1-v_i)^{b_i-1} = \text{Beta}(v_i|a_i, b_i)$
 - $a_i = 1 + (L - l_i + 1) \sum_{n=1}^N \sum_{\zeta_{i_0+1}, \dots, \zeta_L} q(\zeta_n = [1, \zeta_2, \dots, \zeta_{l_0}, \zeta_{l_0+1}, \dots, \zeta_L])$
 - $b_i = \gamma + (L - l_i + 1) \sum_{n=1}^N \sum_{j, \zeta_{i_0+1}, \dots, \zeta_L: j > \zeta_{i_0}} q(\zeta_n = [1, \zeta_2, \dots, \zeta_{l_0-1}, j, \zeta_{l_0+1}, \dots, \zeta_L])$
 - * l_i : The level of the mixture component i
- $q(\zeta_n|\mathbf{x}_n) \propto S_{n\bar{\zeta}} \triangleq \sum_{\zeta \in \text{child}(\bar{\zeta})} S_{n\zeta}$
 - $\bar{\zeta}$: a path in the truncated tree T , either an *inner path* (a path ending at an internal node) or a *full path* (a path ending at a leaf node)
 - $\text{child}(\bar{\zeta})$: the set of all full paths that are not in T but include $\bar{\zeta}$ as a sub path
 - * As a special case, if $\bar{\zeta}$ is a full path, $\text{child}(\bar{\zeta})$ just contains itself
 - In the case of a full path,

$$\begin{aligned}
S_{n\bar{\zeta}} &= S_{n\zeta} \\
&= \exp \left\{ \mathbb{E}_q \left[\sum_{l=1}^L \sum_{l'=1}^l \left(\log v_{1, \zeta_2, \dots, \zeta_{l'}} + \sum_{j=1}^{\zeta_{l'}-1} \log(1 - v_{1, \zeta_2, \dots, j}) \right) \right] + Z_0 \right\} \\
&= \exp \left\{ \mathbb{E}_q \left[\sum_{l=1}^L \left(\sum_{l'=1}^l \left(\log v_{1, \zeta_2, \dots, \zeta_{l'}} + \sum_{j=1}^{\zeta_{l'}-1} \log(1 - v_{1, \zeta_2, \dots, j}) \right) + \log \mathcal{N}(\mathbf{z}_n | \boldsymbol{\mu}_{\zeta_{nl}}, \boldsymbol{\sigma}_{\zeta_{nl}}^2 \mathbf{I}_J) \right) \right] \right\} \\
&= \exp \left\{ \sum_{l=1}^L (L - l + 1) \left(\mathbb{E}_q[\log v_{1, \zeta_2, \dots, \zeta_l}] + \sum_{j=1}^{\zeta_l-1} \mathbb{E}_q[\log(1 - v_{1, \zeta_2, \dots, j})] + \log \mathcal{N}(\mathbf{z}_n | \boldsymbol{\mu}_{\zeta_{nl}}, \boldsymbol{\sigma}_{\zeta_{nl}}^2 \mathbf{I}_J) \right) \right\}
\end{aligned}$$

- In the case of an inner path, $\bar{\zeta} \triangleq [1, \bar{\zeta}_2, \dots, \bar{\zeta}_{l_0}] \subset \mathcal{M}_T$
 - * $\text{child}(\bar{\zeta}) \triangleq \{[\bar{\zeta}, \bar{\zeta}_{l_0+1}, \dots, \bar{\zeta}_L] : \bar{\zeta}_{l_0+1} > j_0\}$
 - * j_0 : maximum index for the child node whose parent path is $\bar{\zeta}$

$$\begin{aligned}
S_{n\bar{\zeta}} &= \sum_{\zeta \in \text{child}(\bar{\zeta})} S_{n\zeta} \\
&= \sum_{\zeta \in \text{child}(\bar{\zeta})} \exp \left\{ \mathbb{E}_q \left[\sum_{l=1}^L \left(\sum_{l'=1}^l \left(\log v_{1, \zeta_2, \dots, \zeta_{l'}} + \sum_{j=1}^{\zeta_{l'}-1} \log(1 - v_{1, \zeta_2, \dots, j}) \right) + \log \mathcal{N}(\mathbf{z}_n | \boldsymbol{\mu}_{\zeta_{nl}}, \boldsymbol{\sigma}_{\zeta_{nl}}^2 \mathbf{I}_J) \right) \right] \right\} \\
&= \frac{\exp \left\{ (\mathbb{E}_q[\sum_{l=l_0+1}^L \log \mathcal{N}(\mathbf{z}_n | \boldsymbol{\mu}_{\zeta_{nl}}, \boldsymbol{\sigma}_{\zeta_{nl}}^2 \mathbf{I}_J)] + \sum_{l=l_0+1}^L \log(L - l + 1) + (L - l_0)(\psi(1) - \psi(1 + \gamma))) \right\}}{(1 - \exp\{\psi(\gamma) - \psi(1 + \gamma)\})^{L-l_0}} \\
&\exp \left\{ \sum_{l=1}^{l_0} (L - l + 1) \left(\mathbb{E}_q[\log v_{1, \zeta_2, \dots, \zeta_l}] + \sum_{j=1}^{\zeta_l-1} \mathbb{E}_q[\log(1 - v_{1, \zeta_2, \dots, j})] \right) + \log \mathcal{N}(\mathbf{z}_n | \boldsymbol{\mu}_{\zeta_{nl}}, \boldsymbol{\sigma}_{\zeta_{nl}}^2 \mathbf{I}_J) \right\} \\
&\exp \left\{ \sum_{j=1}^{j_0} \mathbb{E}_q \left[(L - l_0) \log(1 - v_{1, \zeta_2, \dots, \zeta_{l_0}, j}) \right] \right\}
\end{aligned}$$

- $q\phi_{\boldsymbol{\eta}}(\boldsymbol{\eta}_n|\mathbf{x}_n) = \text{Dirichlet}(\boldsymbol{\eta}_n|\tilde{\boldsymbol{\alpha}}_n)$
 - * $\tilde{\alpha}_{nl} = \frac{1}{\sigma_{nl}^2} \left(1 - \frac{2}{L} + \frac{e^{-\tilde{\mu}_{nl}}}{L^2} \sum_{l'=1}^L e^{-\tilde{\mu}_{nl'}} \right)$
- $q(l_n|\boldsymbol{\omega}_n, \mathbf{x}_n) = \text{Multinomial}(l_n|\boldsymbol{\omega}_n)$
 - * $\omega_{nl} \propto \exp \left\{ \psi(\tilde{\alpha}_{nl}) - \psi(\tilde{\alpha}_{n0}) + \sum_{\zeta} S_{n\zeta} \left(\sum_{j=1}^J -\frac{1}{2} \log(2\pi\sigma_{\zeta_{nl}, j}^2) - \frac{\tilde{\sigma}_{z_{nj}}^2}{2\sigma_{\zeta_{nl}, j}^2} - \frac{(\tilde{\mu}_{z_{nj}} - \mu_{\zeta_{nl}, j})^2}{2\sigma_{\zeta_{nl}, j}^2} \right) \right\}$
 - $\tilde{\alpha}_{n0} = \sum_{i=1}^L \tilde{\alpha}_{ni}$

* Derivation for ω_{nl}

$$\begin{aligned}
\mathcal{L}_{\omega_{nl}} &= \sum_l \omega_{nl} (\psi(\tilde{\alpha}_{nl}) - \psi(\tilde{\alpha}_{n0})) - \sum_l \omega_{nl} \log \omega_{nl} \\
&+ \sum_{\zeta'} S_{n\zeta'} \left\{ \sum_l \omega_{nl} \left(\sum_{j=1}^J -\frac{1}{2} \log(2\pi\sigma_{\zeta',j}^2) - \frac{\tilde{\sigma}_{z_{nj}}^2}{2\sigma_{\zeta',j}^2} - \frac{(\tilde{\mu}_{z_{nj}} - \mu_{\zeta',j})^2}{2\sigma_{\zeta',j}^2} \right) \right\} \\
&= \omega_{nl} (\psi(\tilde{\alpha}_{nl}) - \psi(\tilde{\alpha}_{n0})) + \sum_{\zeta'} S_{n\zeta'} \omega_{nl} \left(\sum_{j=1}^J -\frac{1}{2} \log(2\pi\sigma_{\zeta',j}^2) - \frac{\tilde{\sigma}_{z_{nj}}^2}{2\sigma_{\zeta',j}^2} - \frac{(\tilde{\mu}_{z_{nj}} - \mu_{\zeta',j})^2}{2\sigma_{\zeta',j}^2} \right) \\
&- \omega_{nl} \log \omega_{nl} + \lambda \left(\sum_l \omega_{nl} - 1 \right) \\
\frac{\partial \mathcal{L}_{\omega_{nl}}}{\partial \omega_{nl}} &= (\psi(\tilde{\alpha}_{nl}) - \psi(\tilde{\alpha}_{n0})) + \sum_{\zeta'} S_{n\zeta'} \left(\sum_{j=1}^J -\frac{1}{2} \log(2\pi\sigma_{\zeta',j}^2) - \frac{\tilde{\sigma}_{z_{nj}}^2}{2\sigma_{\zeta',j}^2} - \frac{(\tilde{\mu}_{z_{nj}} - \mu_{\zeta',j})^2}{2\sigma_{\zeta',j}^2} \right) \\
&- (\log \omega_{nl} + 1) + \lambda \\
\omega_{nl} &= \exp \left\{ \psi(\tilde{\alpha}_{nl}) - \psi(\tilde{\alpha}_{n0}) + \sum_{\zeta'} S_{n\zeta'} \left(\sum_{j=1}^J -\frac{1}{2} \log(2\pi\sigma_{\zeta',j}^2) - \frac{\tilde{\sigma}_{z_{nj}}^2}{2\sigma_{\zeta',j}^2} - \frac{(\tilde{\mu}_{z_{nj}} - \mu_{\zeta',j})^2}{2\sigma_{\zeta',j}^2} \right) - 1 + \lambda \right\} \\
\omega_{nl} &\propto \exp \left\{ \psi(\tilde{\alpha}_{nl}) - \psi(\tilde{\alpha}_{n0}) + \sum_{\zeta'} S_{n\zeta'} \left(\sum_{j=1}^J -\frac{1}{2} \log(2\pi\sigma_{\zeta',j}^2) - \frac{\tilde{\sigma}_{z_{nj}}^2}{2\sigma_{\zeta',j}^2} - \frac{(\tilde{\mu}_{z_{nj}} - \mu_{\zeta',j})^2}{2\sigma_{\zeta',j}^2} \right) \right\} \\
&- q_{\phi_z}(z_n | \mathbf{x}_n) = \mathcal{N}(z_n | \tilde{\boldsymbol{\mu}}_{z_n}, \tilde{\boldsymbol{\sigma}}_{z_n}^2 \mathbf{I}_J)
\end{aligned}$$

F EVIDENCE LOWER BOUND

In this section, we present the detailed derivation of the ELBO in Equation 6, which is the objective function for learning HCRL.

$$\begin{aligned}
\log p(\mathbf{x}) &\geq \mathcal{L}_{ELBO}(\mathbf{x}) = \mathbb{E}_q \left[\log \frac{p(\mathbf{v}, \boldsymbol{\zeta}, \boldsymbol{\eta}, \mathbf{l}, \mathbf{z}, \mathbf{x})}{q(\mathbf{v}, \boldsymbol{\zeta}, \boldsymbol{\eta}, \mathbf{l}, \mathbf{z} | \mathbf{x})} \right] \\
&= \mathbb{E}_q \left[\log \frac{\prod_{i \in \mathcal{M}_T} p(v_i | \gamma) \prod_{n=1}^N p(\boldsymbol{\zeta}_n | \mathbf{v}) p(\boldsymbol{\eta}_n | \boldsymbol{\alpha}) p(l_n | \boldsymbol{\eta}_n) p(z_n | \boldsymbol{\zeta}_n, l_n) p_{\boldsymbol{\theta}}(\mathbf{x}_n | \mathbf{z}_n)}{\prod_{i \in \mathcal{M}_T} q(v_i | a_i, b_i) \prod_{n=1}^N q(\boldsymbol{\zeta}_n | \mathbf{x}_n) q_{\phi_{\boldsymbol{\eta}}}(\boldsymbol{\eta}_n | \mathbf{x}_n) q(l_n | \boldsymbol{\omega}_n, \mathbf{x}_n) q_{\phi_z}(z_n | \mathbf{x}_n)} \right] \\
&= \sum_{i \in \mathcal{M}_T} \mathbb{E}_q[\log p(v_i | \gamma)] + \sum_{n=1}^N \mathbb{E}_q[\log p(\boldsymbol{\zeta}_n | \mathbf{v}) + \log p(\boldsymbol{\eta}_n | \boldsymbol{\alpha}) + \log p(l_n | \boldsymbol{\eta}_n) + \log p(z_n | \boldsymbol{\zeta}_n, l_n) \\
&+ \log p_{\boldsymbol{\theta}}(\mathbf{x}_n | \mathbf{z}_n)] - \sum_{i \in \mathcal{M}_T} \mathbb{E}_q[\log q(v_i | a_i, b_i)] - \sum_{n=1}^N \mathbb{E}_q[\log q(\boldsymbol{\zeta}_n | \mathbf{x}_n) + \log q_{\phi_{\boldsymbol{\eta}}}(\boldsymbol{\eta}_n | \mathbf{x}_n) \\
&+ \log q(l_n | \boldsymbol{\omega}_n, \mathbf{x}_n) + \log q_{\phi_z}(z_n | \mathbf{x}_n)] \tag{6}
\end{aligned}$$

F.1 DETAILED DERIVATION FOR ELBO

The followings are additional notations used for the detailed derivation:

- ψ : The digamma function
- $\tilde{\alpha}_{n0} = \sum_{i=1}^L \tilde{\alpha}_{ni}$, $\alpha_0 = \sum_{i=1}^L \alpha_i$

$$\begin{aligned}
\mathbb{E}_q[\log p(v_i|\gamma)] &= \int_{v'} \int_{z'} \int_{\boldsymbol{\eta}'} \sum_{l'} \sum_{\boldsymbol{\zeta}'} \prod_{j \notin \mathcal{M}_T} p(v_j = v') \prod_{i \in \mathcal{M}_T} q(v_i = v') \prod_{n=1}^N q(\boldsymbol{\zeta}_n = \boldsymbol{\zeta}') q(\boldsymbol{\eta}_n = \boldsymbol{\eta}') \\
&\quad q(l_n = l') q(\mathbf{z}_n = \mathbf{z}') \log p(v_i = v') d\boldsymbol{\eta}' d\mathbf{z}' dv' \\
&= \int_{v'} \prod_{j \notin \mathcal{M}_T} p(v_j = v') \prod_{i \in \mathcal{M}_T} q(v_i = v') \log p(v_i = v') dv' \\
&= \int_{v'} q(v_i = v') \log p(v_i = v') dv' \\
&= \int_{v'} \text{Beta}(v'|a_i, b_i) \cdot \log \text{Beta}(v'|1, \gamma) dv' \\
&= \log \Gamma(1 + \gamma) - \log \Gamma(1) - \log \Gamma(\gamma) + (1 - 1)\psi(a_i) + (\gamma - 1)\psi(b_i) \\
&\quad + (-1 - \gamma + 2)\psi(a_i + b_i) \\
&= \log \Gamma(1 + \gamma) - \log \Gamma(\gamma) + (\gamma - 1)\psi(b_i) + (1 - \gamma)\psi(a_i + b_i) \\
&= \log \Gamma(1 + \gamma) - \log \Gamma(\gamma) + (\gamma - 1)(\psi(b_i) - \psi(a_i + b_i)) \\
&= \log(\gamma \Gamma(\gamma)) - \log \Gamma(\gamma) + (\gamma - 1)(\psi(b_i) - \psi(a_i + b_i)) \\
&= \log \gamma + \log \Gamma(\gamma) - \log \Gamma(\gamma) + (\gamma - 1)(\psi(b_i) - \psi(a_i + b_i)) \\
&= \log \gamma + (\gamma - 1)(\psi(b_i) - \psi(a_i + b_i))
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}_q[\log p(\boldsymbol{\zeta}_n|\mathbf{v})] &= \frac{\exp\{\log \Gamma(L - l_0 + 1) + (L - l_0)(\psi(1) - \psi(1 + \gamma))\}}{(1 - \exp\{\psi(\gamma) - \psi(1 + \gamma)\})^{L - l_0}} \\
&\quad \times \exp\left\{\sum_{l=1}^{l_0} (L - l + 1) \left(\mathbb{E}_q[\log v_{1, \zeta_2, \dots, \zeta_l}] + \sum_{j=1}^{\zeta_l - 1} \mathbb{E}_q[\log(1 - v_{1, \zeta_2, \dots, j})] \right)\right\} \\
&\quad \times \exp\left\{\mathbb{E}_q\left[(L - l_0) \sum_{j=1}^{j_0} \log(1 - v_{1, \zeta_2, \dots, \zeta_{l_0}, j})\right]\right\}
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}_q[\log p(\boldsymbol{\eta}_n|\boldsymbol{\alpha})] &= \int_{v'} \int_{z'} \int_{\boldsymbol{\eta}'} \sum_{l'} \sum_{\boldsymbol{\zeta}'} \prod_{j \notin \mathcal{M}_T} p(v_j = v') \prod_{i \in \mathcal{M}_T} q(v_i = v') q(\boldsymbol{\zeta}_n = \boldsymbol{\zeta}') q(\boldsymbol{\eta}_n = \boldsymbol{\eta}') \\
&\quad q(l_n = l') q(\mathbf{z}_n = \mathbf{z}') \log p(\boldsymbol{\eta}_n = \boldsymbol{\eta}'|\boldsymbol{\alpha}) d\boldsymbol{\eta}' d\mathbf{z}' dv' \\
&= \int_{\boldsymbol{\eta}'} q(\boldsymbol{\eta}_n = \boldsymbol{\eta}') \log p(\boldsymbol{\eta}_n = \boldsymbol{\eta}'|\boldsymbol{\alpha}) d\boldsymbol{\eta}' = \int_{\boldsymbol{\eta}'} \text{Dir}(\boldsymbol{\eta}'|\tilde{\boldsymbol{\alpha}}_n) \cdot \log \text{Dir}(\boldsymbol{\eta}'|\boldsymbol{\alpha}) d\boldsymbol{\eta}' \\
&= \log \Gamma(\alpha_0) - \sum_{i=1}^L \log \Gamma(\alpha_i) + \sum_{i=1}^L (\alpha_i - 1)(\psi(\tilde{\alpha}_{ni}) - \psi(\tilde{\alpha}_{n0}))
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}_q[\log p(l_n|\boldsymbol{\eta}_n)] &= \int_{v'} \int_{z'} \int_{\boldsymbol{\eta}'} \sum_{l'} \sum_{\boldsymbol{\zeta}'} \prod_{j \notin \mathcal{M}_T} p(v_j = v') \prod_{i \in \mathcal{M}_T} q(v_i = v') q(\boldsymbol{\zeta}_n = \boldsymbol{\zeta}') q(\boldsymbol{\eta}_n = \boldsymbol{\eta}') \\
&\quad q(l_n = l') q(\mathbf{z}_n = \mathbf{z}') \log p(l_n = l'|\boldsymbol{\eta}_n) d\boldsymbol{\eta}' d\mathbf{z}' dv' \\
&= \int_{\boldsymbol{\eta}'} \sum_{l'} q(\boldsymbol{\eta}_n = \boldsymbol{\eta}') q(l_n = l) \log p(l_n = l|\boldsymbol{\eta}_n = \boldsymbol{\eta}') d\boldsymbol{\eta}' \\
&= \sum_{l'} q(l_n = l') \int_{\boldsymbol{\eta}'} q(\boldsymbol{\eta}_n = \boldsymbol{\eta}') \log \text{Mult}(l'|\boldsymbol{\eta}') d\boldsymbol{\eta}' \\
&= \sum_{l'} \omega_{nl'} \int_{\boldsymbol{\eta}'} \text{Dir}(\boldsymbol{\eta}'|\tilde{\boldsymbol{\alpha}}_n) \log \eta'_l d\boldsymbol{\eta}' \\
&= \sum_{l'} \omega_{nl'} (\psi(\tilde{\alpha}_{nl'}) - \psi(\tilde{\alpha}_{n0}))
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}_q[\log p(\mathbf{z}_n|\zeta_n, l_n)] &= \int_{v'} \int_{\mathbf{z}'} \int_{\boldsymbol{\eta}'} \sum_{v'} \sum_{\zeta'} \prod_{j \notin \mathcal{M}_T} p(v_j = v') \prod_{i \in \mathcal{M}_T} q(v_i = v') q(\zeta_n = \zeta) q(\boldsymbol{\eta}_n = \boldsymbol{\eta}') \\
&\quad q(l_n = l') q(\mathbf{z}_n = \mathbf{z}) \log p(\mathbf{z}_n = \mathbf{z} | \zeta_n = \zeta, l_n = l') d\boldsymbol{\eta}' d\mathbf{z}' dv' \\
&= \int_{\mathbf{z}'} \sum_{v'} \sum_{\zeta'} q(\zeta_n = \zeta) q(l_n = l') q(\mathbf{z}_n = \mathbf{z}) \log p(\mathbf{z}_n = \mathbf{z} | \zeta_n = \zeta, l_n = l') d\mathbf{z}' \\
&= \sum_{v'} \sum_{\zeta'} q(\zeta_n = \zeta) q(l_n = l') \int_{\mathbf{z}'} q(\mathbf{z}_n = \mathbf{z}) \log p(\mathbf{z}_n = \mathbf{z} | \zeta_n = \zeta, l_n = l') d\mathbf{z}' \\
&= \sum_{v'} \sum_{\zeta'} q(\zeta_n = \zeta') q(l_n = l') \int_{\mathbf{z}'} \mathcal{N}(\mathbf{z} | \tilde{\boldsymbol{\mu}}_{\mathbf{z}'}', \tilde{\boldsymbol{\sigma}}_{\mathbf{z}'}'^2, \mathbf{I}_J) \cdot \log \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}_{\zeta_{n'l'}}', \boldsymbol{\sigma}_{\zeta_{n'l'}}'^2, \mathbf{I}_J) d\mathbf{z}' \\
&= \sum_{\zeta'} S_{n\zeta'} \left\{ \sum_{v'} \omega_{n'l'} \left(\sum_{j=1}^J -\frac{1}{2} \log(2\pi\sigma_{\zeta_{n'l'},j}^2) - \frac{\tilde{\sigma}_{z_{nj}}^2}{2\sigma_{\zeta_{n'l'},j}^2} - \frac{(\tilde{\mu}_{z_{nj}} - \mu_{\zeta_{n'l'},j})^2}{2\sigma_{\zeta_{n'l'},j}^2} \right) \right\}
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}_q[\log p_{\boldsymbol{\theta}}(\mathbf{x}_n|\mathbf{z}_n)] &= \int_{v'} \int_{\mathbf{z}'} \int_{\boldsymbol{\eta}'} \sum_{v'} \sum_{\zeta'} \prod_{j \notin \mathcal{M}_T} p(v_j = v') \prod_{i \in \mathcal{M}_T} q(v_i = v') q(\zeta_n = \zeta) q(\boldsymbol{\eta}_n = \boldsymbol{\eta}') \\
&\quad q(l_n = l') q(\mathbf{z}_n = \mathbf{z}) \log p_{\boldsymbol{\theta}}(\mathbf{x}_n | \mathbf{z}_n = \mathbf{z}') d\boldsymbol{\eta}' d\mathbf{z}' dv' \\
&= \int_{\mathbf{z}'} q(\mathbf{z}_n = \mathbf{z}) \log p_{\boldsymbol{\theta}}(\mathbf{x}_n | \mathbf{z}_n = \mathbf{z}') d\mathbf{z}' \\
&\approx \frac{1}{S} \sum_{s=1}^S \log p_{\boldsymbol{\theta}}(\mathbf{x}_n | \mathbf{z}_n^{(s)}), \text{ where } \mathbf{z}^{(i,s)} = \boldsymbol{\mu}_{\mathbf{x}}^{(i)} + \boldsymbol{\sigma}_{\mathbf{x}}^{(i)} \odot \boldsymbol{\epsilon}^{(s)} \text{ and } \boldsymbol{\epsilon}^{(s)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_J) \\
&= \begin{cases} \frac{1}{S} \sum_{s=1}^S \sum_{d=1}^D x_{nd} \log \mu_{x_{nd}}^{(s)} + (1 - x_{nd}) \log(1 - \mu_{x_{nd}}^{(s)}) & \text{if } \mathbf{x}_n \text{ is binary} \\ \frac{1}{S} \sum_{s=1}^S \sum_{d=1}^D -\frac{1}{2} \log(2\pi\sigma_{x_{nd}}^{(s)2}) - \frac{(x_{nd} - \mu_{x_{nd}}^{(s)})^2}{2\sigma_{x_{nd}}^{(s)2}} & \text{if } \mathbf{x}_n \text{ is real-valued} \end{cases}
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}_q[\log q(v_i | a_i, b_i)] &= \int_{v'} \int_{\mathbf{z}'} \int_{\boldsymbol{\eta}'} \sum_{v'} \sum_{\zeta'} \prod_{j \notin \mathcal{M}_T} p(v_j = v') \prod_{i \in \mathcal{M}_T} q(v_i = v') q(\zeta_n = \zeta) q(\boldsymbol{\eta}_n = \boldsymbol{\eta}') \\
&\quad q(l_n = l') q(\mathbf{z}_n = \mathbf{z}) \log q(v_i = v') d\boldsymbol{\eta}' d\mathbf{z}' dv' \\
&= \int_{v'} \prod_{j \notin \mathcal{M}_T} p(v_j = v') \prod_{i \in \mathcal{M}_T} q(v_i = v') \log q(v_i = v') dv' \\
&= \int_{v'} q(v_i = v') \log q(v_i = v') dv' = \int_{v'} \text{Beta}(v' | a_i, b_i) \cdot \log \text{Beta}(v' | a_i, b_i) dv' \\
&= \log \Gamma(a_i + b_i) - \log \Gamma(a_i) - \log \Gamma(b_i) + (a_i - 1)\psi(a_i) + (b_i - 1)\psi(b_i) \\
&\quad + (-a_i - b_i + 2)\psi(a_i + b_i)
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}_q[\log q(\zeta_n | \mathbf{x}_n)] &= \int_{v'} \int_{\mathbf{z}'} \int_{\boldsymbol{\eta}'} \sum_{v'} \sum_{\zeta} \prod_{j \notin \mathcal{M}_T} p(v_j = v') \prod_{i \in \mathcal{M}_T} q(v_i = v') q(\zeta_n = \zeta) q(\boldsymbol{\eta}_n = \boldsymbol{\eta}') \\
&\quad q(l_n = l') q(\mathbf{z}_n = \mathbf{z}) \log q(\zeta_n = \zeta) d\boldsymbol{\eta}' d\mathbf{z}' dv' \\
&= \sum_{\zeta'} q(\zeta_n = \zeta) \log q(\zeta_n = \zeta) \\
&= \sum_{\zeta'} \frac{S_{n\zeta'}}{\sum_{\zeta''} S_{n\zeta''}} \log \frac{S_{n\zeta'}}{\sum_{\zeta''} S_{n\zeta''}}
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}_q[\log q_{\phi_{\eta}}(\boldsymbol{\eta}_n | \mathbf{x}_n)] &= \int_{v'} \int_{z'} \int_{\boldsymbol{\eta}'} \sum_{l'} \sum_{\boldsymbol{\zeta}'} \prod_{j \notin \mathcal{M}_T} p(v_j = v') \prod_{i \in \mathcal{M}_T} q(v_i = v') q(\boldsymbol{\zeta}_n = \boldsymbol{\zeta}) q(\boldsymbol{\eta}_n = \boldsymbol{\eta}') \\
&\quad q(l_n = l') q(\mathbf{z}_n = \mathbf{z}) \log q(\boldsymbol{\eta}_n = \boldsymbol{\eta}') d\boldsymbol{\eta}' dz' dv' \\
&= \int_{\boldsymbol{\eta}'} q(\boldsymbol{\eta}_n = \boldsymbol{\eta}') \log q(\boldsymbol{\eta}_n = \boldsymbol{\eta}') d\boldsymbol{\eta}' = \int_{\boldsymbol{\eta}'} \text{Dir}(\boldsymbol{\eta}' | \tilde{\boldsymbol{\alpha}}_n) \cdot \log \text{Dir}(\boldsymbol{\eta}' | \tilde{\boldsymbol{\alpha}}_n) d\boldsymbol{\eta}' \\
&= \log \Gamma(\tilde{\boldsymbol{\alpha}}_{n0}) - \sum_{i=1}^L \log \Gamma(\tilde{\boldsymbol{\alpha}}_{ni}) + \sum_{i=1}^L (\tilde{\boldsymbol{\alpha}}_{ni} - 1) (\psi(\tilde{\boldsymbol{\alpha}}_{ni}) - \psi(\tilde{\boldsymbol{\alpha}}_{n0}))
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}_q[\log q(l_n | \boldsymbol{\omega}_n, \mathbf{x}_n)] &= \int_{v'} \int_{z'} \int_{\boldsymbol{\eta}'} \sum_{l'} \sum_{\boldsymbol{\zeta}'} \prod_{j \notin \mathcal{M}_T} p(v_j = v') \prod_{i \in \mathcal{M}_T} q(v_i = v') q(\boldsymbol{\zeta}_n = \boldsymbol{\zeta}) q(\boldsymbol{\eta}_n = \boldsymbol{\eta}') \\
&\quad q(l_n = l') q(\mathbf{z}_n = \mathbf{z}) \log q(l_n = l') d\boldsymbol{\eta}' dz' dv' \\
&= \sum_{l'} q(l_n = l') \log q(l_n = l') = \sum_{l'} \text{Mult}(l' | \boldsymbol{\omega}_n) \log \text{Mult}(l' | \boldsymbol{\omega}_n) \\
&= \sum_{l'} \omega_{nl'} \cdot \log \omega_{nl'}
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}_q[\log q_{\phi_z}(\mathbf{z}_n | \mathbf{x}_n)] &= \int_{v'} \int_{z'} \int_{\boldsymbol{\eta}'} \sum_{l'} \sum_{\boldsymbol{\zeta}'} \prod_{j \notin \mathcal{M}_T} p(v_j = v') \prod_{i \in \mathcal{M}_T} q(v_i = v') q(\boldsymbol{\zeta}_n = \boldsymbol{\zeta}) q(\boldsymbol{\eta}_n = \boldsymbol{\eta}') \\
&\quad q(l_n = l') q(\mathbf{z}_n = \mathbf{z}) \log q(\mathbf{z}_n = \mathbf{z}) d\boldsymbol{\eta}' dz' dv' \\
&= \int_{z'} \int_{\boldsymbol{\eta}'} q(\mathbf{z}_n = \mathbf{z}) \log q(\mathbf{z}_n = \mathbf{z}) d\boldsymbol{\eta}' dz' \\
&= \int_{z'} q(\mathbf{z}_n = \mathbf{z}) \log q(\mathbf{z}_n = \mathbf{z}) dz' = \int_{z'} \mathcal{N}(\mathbf{z} | \tilde{\boldsymbol{\mu}}_{z'}, \tilde{\boldsymbol{\sigma}}_z^2, \mathbf{I}_J) \cdot \log \mathcal{N}(\mathbf{z} | \tilde{\boldsymbol{\mu}}_{z'}, \tilde{\boldsymbol{\sigma}}_z^2, \mathbf{I}_J) dz' \\
&= -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J (1 + \log \tilde{\sigma}_{z_{n,j}}^2)
\end{aligned}$$