

1 We sincerely thank all reviewers for their valuable comments. Our responses to the comments are listed below.

2 **To Reviewer #1:** Q1. About Eq.12 and reference [19] (q1 under **Correctness**).

$$\begin{aligned}
 p(A|Z) &= \xi \prod_{A_{ij}=1} \exp(f(C_i, C_j)) \prod_{A_{ij}=0} \exp(1 - f(C_i, C_j)) \\
 &\approx \prod_{A_{ij}=1} \frac{\exp(f(C_i, C_j))}{\exp(f(C_i, C_j)) + \exp(1 - f(C_i, C_j))} \prod_{A_{ij}=0} \frac{\exp(1 - f(C_i, C_j))}{\exp(f(C_i, C_j)) + \exp(1 - f(C_i, C_j))} \\
 &= \prod_{A_{ij}=1} \frac{1}{1 + \exp(1 - 2f(C_i, C_j))} \prod_{A_{ij}=0} \frac{\exp(1 - 2f(C_i, C_j))}{1 + \exp(1 - 2f(C_i, C_j))} \\
 &= \prod_{A_{ij}=1} \sigma(C_i^\top C_j) \prod_{A_{ij}=0} (1 - \sigma(C_i^\top C_j)) \quad \text{consider } f(C_i, C_j) = \frac{1}{2} C_j^\top C_j + \frac{1}{2}
 \end{aligned} \tag{1}$$

3 where $\sigma(\cdot)$ is the logistic sigmoid function and the last line in the above equation is the one used in [19]. Z does not
 4 appear on the RHS of the equation, as Z coincides with cluster memberships C (line 117 in the paper). This point can
 5 be verified since both $z_i \in \Delta_{K-1}$ and $C_i \in \Delta_{K-1}$ are the $(K - 1)$ -simplex.

6 Q2. About Claim 4.1 and ELBO (q2 under **Correctness**). If we understand it correctly, “the approximate posterior in
 7 the reconstruction term” refers to the KL term in ELBO. If so, we do consider the effects of the KL term in the proof.
 8 As in line 107 of the paper, we regard $C_i \sim \text{Dir}(\beta)$, which is a Dirichlet posterior with a Dirichlet prior $\text{Dir}(\alpha)$. In other
 9 words, our proof starts with the effects of the KL term. We set the distance metric to be $f(C_i, C_j) = 1 - \text{MSE}(C_i, C_j)$
 10 (line 140 in the paper) in the proof of Claim 4.1.

11 Q3. Does Claim 4.1 rely on a specific distribution (q1 under **Weaknesses**)? A quick answer is yes, as our proof relies
 12 on the fact $C_i \sim \text{Dir}(\beta)$. As for Gaussian variables (the case in VGAE), we are not sure if the claim still holds.

13 Q4. Laplace approximation towards Dirichlet (q3 under **Weaknesses**). In section 3, we used $\mathcal{N}(z'_i; \mu^0, \sigma^0)$ to
 14 approximate the Dirichlet distribution $\text{Dir}(z_i; \phi)$, where $z_i = \text{softmax}(z'_i)$ is based on softmax Laplace approximation.
 15 We agree sometimes it might be hard to follow and we shall add the details in the revision.

16 **To Reviewer #2:** Q1. About Heatts and $q_\phi(\cdot)$ (q1 under **Feedback**). The parameters of the variational approximation
 17 $q_\phi(\cdot)$ are determined by Eq.7 of the DGVAE framework. Putting Heatts into the framework, we derive $\mu^0, \sigma^0 =$
 18 $\text{Heatts}_\phi(A, X)$, as Heatts is a variant of GNN.

19 Q2. About DGVAE and DGAE (q1 under **Weaknesses**). We agree non-variational version performs better in most cases
 20 than its variational one, as observed by many other references [8][19]. In our opinion, the reasons are: (1) component
 21 collapsing, a good optimization of KL term results in a bad reconstruction term; (2) inappropriate priors, we use the
 22 same priors through the experiments. As for the theoretical results (Claim 4.1), we think it holds for DGAE, if we
 23 consider the Dirichlet posterior is only determined by the likelihood, i.e., without priors.

24 Q3. About the limitation of balanced cluster sizes (q2 under **Weaknesses**). We agree balanced cluster sizes could be
 25 a limitation in some cases. However, Claim 4.1 is based on asymptotic analysis. Moreover, the optimization is also
 26 affected by the Dirichlet priors, which provides our model with flexibility.

27 **To Reviewer #3:** Q1. About the interpretation of Dirichlet factors (q1 under **Weaknesses**). We agree that the Dirichlet
 28 approach does not lead to a better explanation compared to other clustering approaches, e.g., spherical k-means.
 29 Compared to the traditional graph variational auto-encoders which focus on Gaussian distributions, Dirichlet approach
 30 leads to a better explanation since its latent variables can be interpreted as cluster memberships.

31 Q2. About comparison of SOTA (q2 under **Weaknesses**). We report the experimental results of VGAE + spherical
 32 K-means (VGAE&S), for Pubmed, VGAE&S achieves 61.2%, 16.7% and 61.1% on ACC, NMI and F1 respectively;
 33 for Citeseer, VGAE&S gets 46.8%, 21.7% and 44.7% respectively; for Wiki, VGAE&S gets 27.8%, 24.3% and 19.5%
 34 respectively. The results show our method outperforms VGAE&S on all metrics. In addition, we would like to highlight
 35 that our method does not rely on any outsourcing, e.g., spherical K-means or GMM, and cluster memberships can be
 36 derived directly from the latent Z , while previous approaches such as VGAE do not have this capacity.

37 Q3. About the three contributions (q3 under **Weaknesses**). Yes, we agree Heatts is the key component that leads to the
 38 good experimental results. However, the good results are not the sole purpose of this work. More precisely, this work
 39 wants to figure out a clear way to understand how the graphs are generated and how we can improve the design of graph
 40 VAEs analytically. In this vein, we first describe the DGVAE framework (the 1st contribution); we then analyse the
 41 framework and find that the optimization of the framework favors *low pass* of GNN (the 2nd contribution); based on
 42 that, we propose Heatts. Without the 1st and 2nd contributions, the design of Heatts may reduce to trial-and-error.

43 Q4. About the corners of simplex (q1 under **Feedback**). Yes, it is imposed by the prior assumption and by optimizing
 44 the reconstruction term. At the initial stage of our model (randomly initialized parameters), the data may appear
 45 anywhere in the simplex. As the optimization goes, the Dirichlet posterior $C_i \sim \text{Dir}(\beta)$ would be “dragged” towards
 46 $\sum_k \beta_k \rightarrow 0$ (Lemma 4.2), which makes the data sampled from this distribution lie at the corners of the simplex. For
 47 the stability, the right figure of Fig. 3 provides an intuitive way to understand the process. As shown, the learning curve
 48 is smooth and drops quickly, indicating our approach is stable. The sampling strategy is illustrated in Eq.8.