

1 Response for Submission 3163 “DM2C: Deep Mixed-Modal Clustering”

2 We thank all the reviewers for their careful and valuable comments.

3 To R1

4 **Q1. Ablation study** We evaluate k -means using latent modality-specific representations obtained before/after the
5 adversarial training (denoted as *modal-spec (with adv.)* and *modal-spec (w/o adv.)* respectively) in our model on the
6 Wikipedia and NUS-WIDE-10K datasets. The results are recorded in Tab. S1. We can observe that the performance
7 of our model is largely improved by the final cross-modal transformations. This indicates that the unification of
8 modality-specific representations could reduce the semantic gap between the modalities. We will add the experiments
9 and discussions if accepted.

Table S1: Ablation study on the Wikipedia and NUS-WIDE-10K datasets. The larger the better.

Dataset	Algorithm	Accuracy	ARI	NMI	F-score	Precision	Recall	Purity
Wikipedia	modal-spec (w/o adv.)	0.2301	0.0340	0.1069	0.1730	0.1289	0.2633	0.2563
	modal-spec (with adv.)	0.2395	0.0290	0.1311	0.1696	0.1256	0.2611	0.2699
	ours	0.2720	0.0558	0.1543	0.1878	0.1439	0.2700	0.3075
NUS-WIDE	modal-spec (w/o adv.)	0.2696	0.0321	0.0719	0.2323	0.3318	0.1787	0.5332
	modal-spec (with adv.)	0.2884	0.0359	0.0672	0.2542	0.3316	0.2060	0.5336
	ours	0.3300	0.0710	0.0951	0.3043	0.3579	0.2648	0.5492

10 **Q2. The value of $n_{critics}$:** This value is empirically set to 5 in the experiments.

11 To R2

12 **Q1. Cycle consistency on multiple modalities:** Perhaps due to our way of writing, it is a pity to leave you an
13 impression that only cross-domain cycle consistency is mentioned in related work. In fact, the reference [29] in our
14 paper is an application of cross-modal cycle consistency. To the best of our knowledge there is only few other relevant
15 work on cross-modal cycle consistency, be it “A Uniform Framework for Cross-Modal Visual-Audio Mutual Generation”
16 (AAAI18) and “Multi-modal Cycle-consistent Generalized Zero-Shot Learning” (ECCV18). However none of them are
17 directly available for the mixed-modal clustering task posed in our paper. We will add a detailed discussion on this
18 issue if accepted.

19 **Q2. “1-Lipschitz constraint” is not explained:** “1-Lipschitz constraint” is a requirement of the dual formulation
20 of the \mathcal{W}_1 distance. More exactly, in our case, it refers to the fact that $D_A(\cdot)$ and $D_B(\cdot)$ are 1-Lipschitz continuous,
21 which means $\|D_A(\mathbf{x}_1) - D_A(\mathbf{x}_2)\| \leq \|\mathbf{x}_1 - \mathbf{x}_2\|$ and $\|D_B(\mathbf{x}_1) - D_B(\mathbf{x}_2)\| \leq \|\mathbf{x}_1 - \mathbf{x}_2\|$ (as what is put down in
22 line 161-162).

23 **Q3. How to sample data from $\mathcal{X}_{A/B}$:** Directly sampling data from $\mathcal{X}_{A/B}$ requires transforming all the data points
24 from $\mathcal{D}_{A/B}$ into $\mathcal{X}_{A/B}$ before sampling, which is impractical for high dimensional data. Hence we adopt a much
25 simpler way that, we 1) sample a batch of data from $\mathcal{D}_{A/B}$, 2) feed them into the auto-encoder A/B to obtain the latent
26 representations lying in $\mathcal{X}_{A/B}$.

27 To R3

28 **Q1. The relationship between generators:** Ideally, as you have suggested, the generators should satisfy
29 $G_{AB} \circ G_{BA}(\cdot) = G_{BA} \circ G_{AB}(\cdot) = I(\cdot)$ considering the cycle consistency constraint. The cycle-consistency
30 regularization terms $\mathcal{L}_{cyc}^A(\Theta_{G_{AB}}, \Theta_{G_{BA}})$, $\mathcal{L}_{cyc}^B(\Theta_{G_{AB}}, \Theta_{G_{BA}})$ guarantee that $\mathbb{E}_{z_a \sim \mathcal{X}_A} [\|z_a - G_{BA}(G_{AB}(z_a))\|_1]$
31 and $\mathbb{E}_{z_b \sim \mathcal{X}_B} [\|z_b - G_{AB}(G_{BA}(z_b))\|_1]$ are small. This approximates the cycle-consistency condition. When
32 $\mathcal{L}_{cyc}^A(\Theta_{G_{AB}}, \Theta_{G_{BA}}) \rightarrow 0$, $\mathcal{L}_{cyc}^B(\Theta_{G_{AB}}, \Theta_{G_{BA}}) \rightarrow 0$, we exactly recover this condition.

33 **Q2. The choice of modality for the final clustering (only A is used in this paper):** For the clustering process, we
34 chose the modality whose data are more informative. In our setting, deep features are available for image modality (A),
35 while the text modality (B) only contains binary features. In this way, the latent representations learned for B obviously
36 have less representability than those for A. As a result, we transform all the data points into modality A.