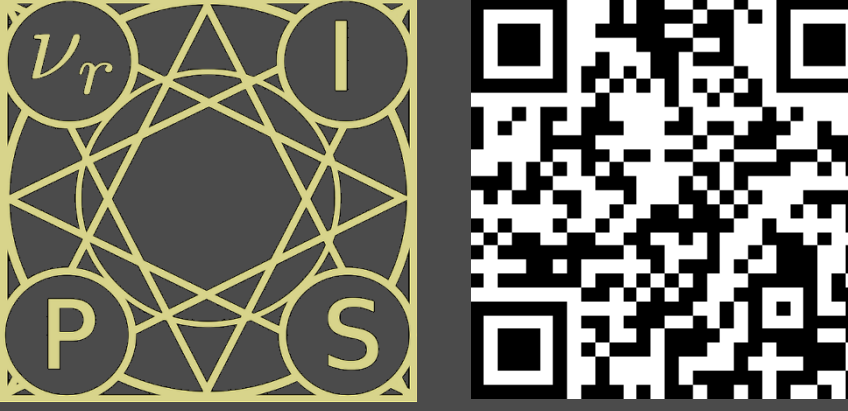# DM2C: Deep Mixed-Modal Clustering

Yangbangyan Jiang[1,2], Qianqian Xu[3], Zhiyong Yang[1,2], Xiaochun Cao[1,2,5], Qingming Huang[2,3,4,5*]

[1]Institute of Information Engineering, CAS    [2]University of Chinese Academy of Sciences

[3]Institute of Computing Technology, CAS    [4]BDKM, CAS    [5]Peng Cheng Laboratory

## Motivation

Traditional multi-modal learning requires extra pairing information among modalities for feature alignment.

- e.g., full/partial modality pairing, 'must/cannot link' constraints, co-occurrence frequency...

Table 1: Types of learning under multiple modalities

| Type | Supervision | |
|---|---|---|
| | Class Label | Modality Pairing |
| Supervised Multi-modal Learning | ✓ | ✓ |
| Unsupervised Multi-modal Learning | ✗ | ✓ |
| Unsupervised Mixed-modal Learning | ✗ | ✗ |

## Mixed-modal Clustering

**Mixed-modal**: Each instance is represented in only one modality.

Dataset $\mathcal{D}$ ⟶ $\mathcal{D}_A = \left\{ x_i^{(a)} \right\}_{i=1}^{n_a}$ and $\mathcal{D}_B = \left\{ x_i^{(b)} \right\}_{i=1}^{n_b}$

Multi-modal data

pairing information

Mixed-modal data

hard for feature alignment

**Goal**: learning unified representations for the modalities, then grouping the samples into $k$ categories.

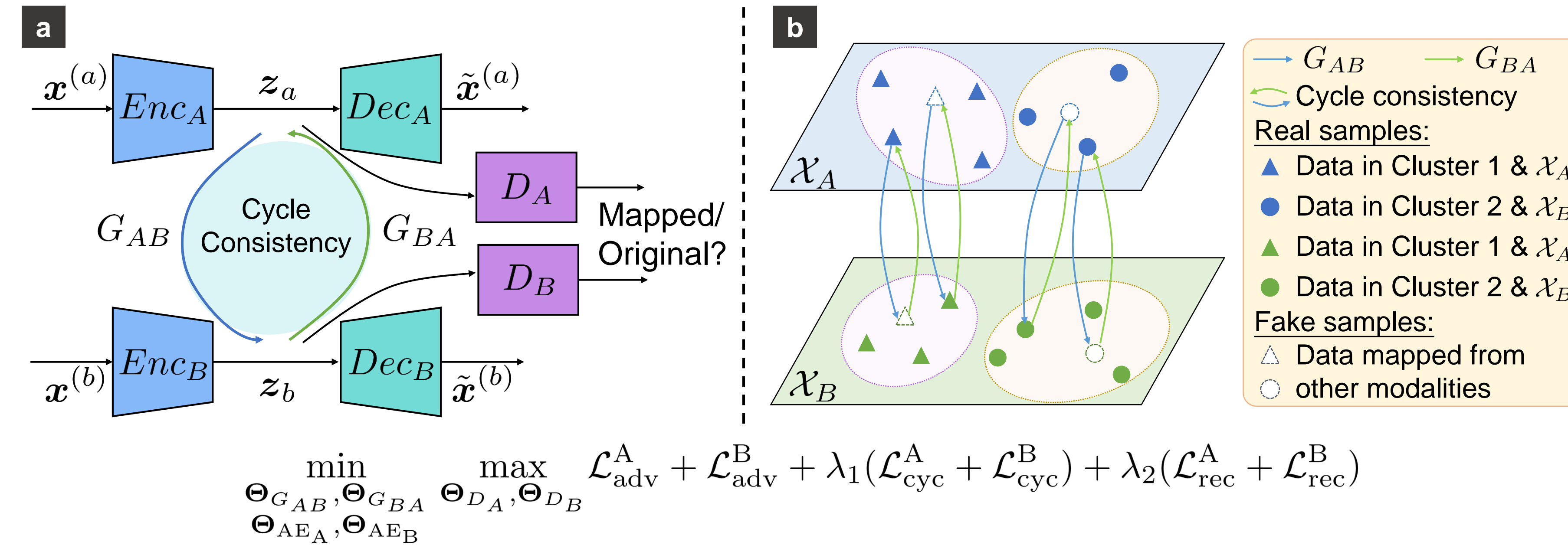## Modality unifying

Modal. 1  Modal. 2 → Joint semantic space — hard to find the correlation

Modal. 1 ↕ Modal. 2

**learn the cross-modal translation**
- easy to obtain via *cycle-consistency*
- **unifying**: transforming all the samples into a modality specific space

## Framework

$$\min_{\substack{\Theta_{G_{AB}}, \Theta_{G_{BA}} \\ \Theta_{AE_A}, \Theta_{AE_B}}} \max_{\Theta_{D_A}, \Theta_{D_B}} \mathcal{L}_{adv}^A + \mathcal{L}_{adv}^B + \lambda_1(\mathcal{L}_{cyc}^A + \mathcal{L}_{cyc}^B) + \lambda_2(\mathcal{L}_{rec}^A + \mathcal{L}_{rec}^B)$$

- Modality-specific Auto-Encoders: latent representations for each modality
$$\mathcal{L}_{rec}^A(\Theta_{AE_A}) = \|x_i^{(a)} - Dec_A(Enc_A(x_i^{(a)}))\|_2^2,$$
$$\mathcal{L}_{rec}^B(\Theta_{AE_B}) = \|x_i^{(b)} - Dec_B(Enc_B(x_i^{(b)}))\|_2^2,$$

- Unpaired Cross-Modal Mappings via *cycle-consistency*
$$\mathcal{L}_{cyc}^A(\Theta_{G_{AB}}, \Theta_{G_{BA}}) = \mathbb{E}_{z_a \sim \mathcal{X}_A} \left[\|z_a - G_{BA}(G_{AB}(z_a))\|_1\right],$$
$$\mathcal{L}_{cyc}^B(\Theta_{G_{AB}}, \Theta_{G_{BA}}) = \mathbb{E}_{z_b \sim \mathcal{X}_B} \left[\|z_b - G_{AB}(G_{BA}(z_b))\|_1\right].$$

- Adversarial learning between Cross-modal Mappings (Generators) and Discriminators
$$\mathcal{L}_{adv}^A(\Theta_{G_{BA}}, \Theta_{D_A}) = \mathbb{E}_{z_a \sim \mathcal{X}_A}[D_A(z_a)] - \mathbb{E}_{z_b \sim \mathcal{X}_B}[D_A(G_{BA}(z_b))],$$
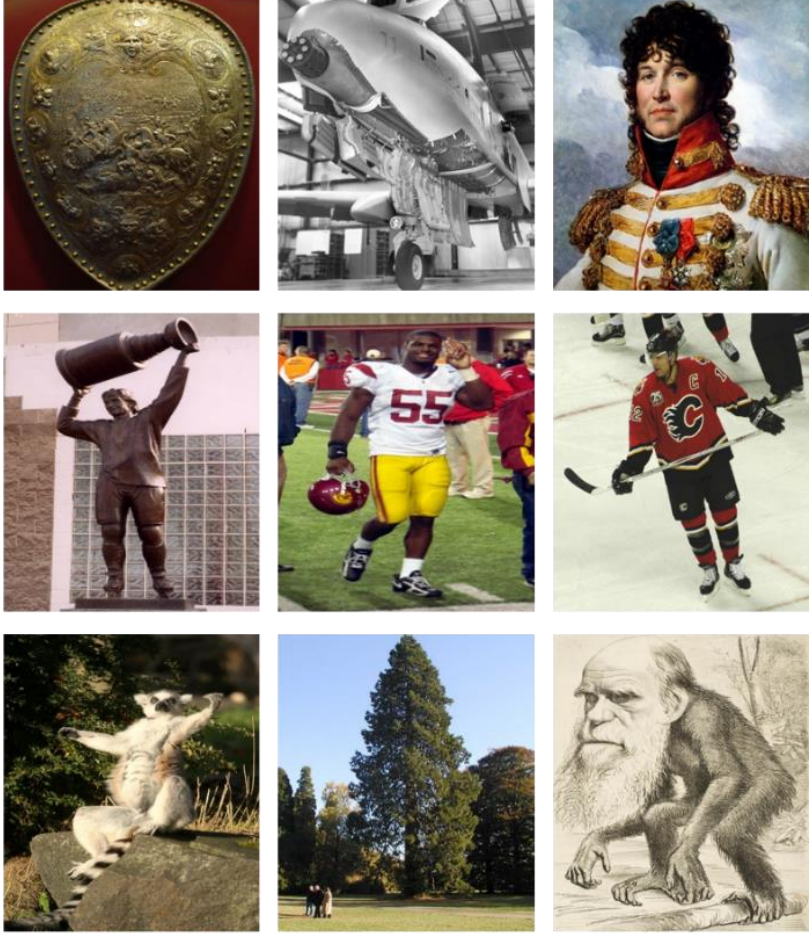$$\mathcal{L}_{adv}^B(\Theta_{G_{AB}}, \Theta_{D_B}) = \mathbb{E}_{z_b \sim \mathcal{X}_B}[D_B(z_b)] - \mathbb{E}_{z_a \sim \mathcal{X}_A}[D_B(G_{AB}(z_a))].$$

## Results

Samples on the Wikipedia dataset

Table 2: Performance comparisons on Wikipedia.

| Algorithm | Accuracy | ARI | NMI | F-score | Purity |
|---|---|---|---|---|---|
| $k$-means | 0.2291 | 0.0166 | 0.1003 | 0.1857 | 0.2301 |
| DKM | 0.2173 | 0.0108 | 0.1170 | 0.1729 | 0.2429 |
| DCN | 0.2215 | 0.0137 | 0.1172 | 0.1688 | 0.2465 |
| IDEC | 0.2153 | 0.0375 | 0.0849 | 0.1654 | 0.2606 |
| IMSAT | 0.2521 | 0.0573 | 0.1093 | 0.1738 | 0.2720 |
| Ours | 0.2720 | 0.0558 | 0.1543 | 0.1878 | 0.3075 |

Table 3: Performance comparisons on NUS-WIDE-10K.

| Algorithm | Accuracy | ARI | NMI | F-score | Purity |
|---|---|---|---|---|---|
| $k$-means | 0.2744 | 0.0044 | 0.0469 | 0.3008 | 0.5208 |
| DKM | 0.2932 | 0.0130 | 0.0116 | 0.2901 | 0.5036 |
| DCN | 0.3036 | 0.0144 | 0.0512 | 0.2959 | 0.5296 |
| IDEC | 0.3045 | 0.0006 | 0.0082 | 0.3048 | 0.5036 |
| IMSAT | 0.3080 | 0.0038 | 0.0064 | 0.3422 | 0.5036 |
| Ours | 0.3300 | 0.0710 | 0.0951 | 0.3043 | 0.5492 |