

DM2C: Deep Mixed-Modal Clustering

Yangbangan Jang^{1,2}, Qianqian Xu³, Zhiyong Yang^{1,2}, Xiaochun Cao^{1,2,5}, Qingming Huang^{2,3,4,5}

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

³Institute of Computing Technology, CAS ⁴BDM, CAS ⁵Peng Cheng Laboratory



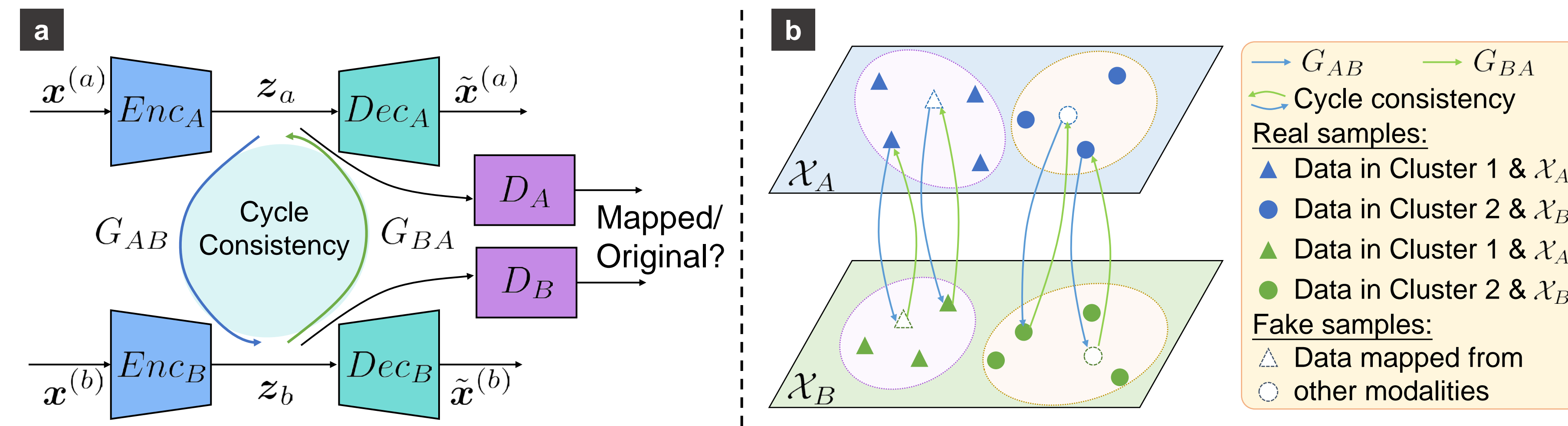
Abstract

Add your information, graphs and images to this section.



Methodology

Key idea: unify the modality spaces via cross-modal translations

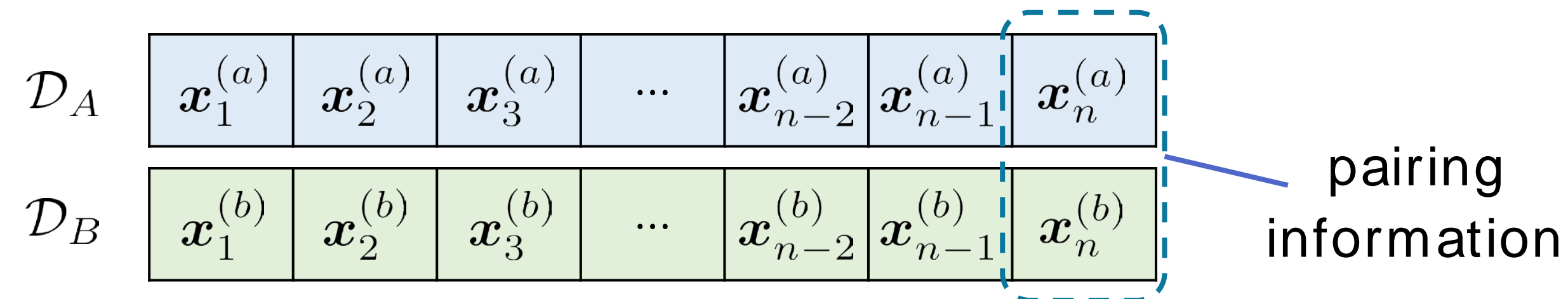


Mixed-modal data

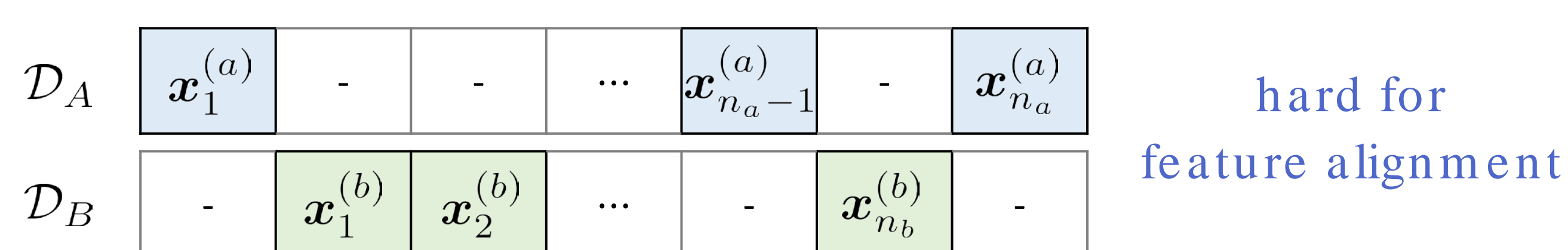
Each instance is represented in only one modality.

Dataset $\mathcal{D} \longrightarrow \mathcal{D}_A = \{\mathbf{x}_i^{(a)}\}_{i=1}^{n_a}$ and $\mathcal{D}_B = \{\mathbf{x}_i^{(b)}\}_{i=1}^{n_b}$

Multi-modal data



Mixed-modal data



Results

Add your information, graphs and images to this section.

Table 1: Dataset statistics.

Dataset	Modal.1	Modal.2	Training samples	Test samples	Categ.
Wikipedia	image	text (article)	1910	256	10
NUS-WIDE-10K	image	text (tag)	7500	2500	10

Table 2: Performance comparisons on Wikipedia.

Algorithm	Accuracy	ARI	NMI	F-score	Purity
<i>k</i> -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
<i>k</i> -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

$$\mathcal{L}(\Theta) = \mathcal{L}_{\text{adv}}^A + \mathcal{L}_{\text{adv}}^B + \lambda_1(\mathcal{L}_{\text{cyc}}^A + \mathcal{L}_{\text{cyc}}^B) + \lambda_2(\mathcal{L}_{\text{rec}}^A + \mathcal{L}_{\text{rec}}^B)$$

- Latent representations

$$\mathcal{L}_{\text{rec}}^A(\Theta_{\text{AE}_A}) = \|\mathbf{x}_i^{(a)} - \text{Dec}_A(\text{Enc}_A(\mathbf{x}_i^{(a)}))\|_2^2,$$

$$\mathcal{L}_{\text{rec}}^B(\Theta_{\text{AE}_B}) = \|\mathbf{x}_i^{(b)} - \text{Dec}_B(\text{Enc}_B(\mathbf{x}_i^{(b)}))\|_2^2,$$

- Unpaired cross-modal mappings

$$\mathcal{L}_{\text{cyc}}^A(\Theta_{G_{AB}}, \Theta_{G_{BA}}) = \mathbb{E}_{\mathbf{z}_a \sim \mathcal{X}_A} [\|\mathbf{z}_a - G_{BA}(G_{AB}(\mathbf{z}_a))\|_1],$$

$$\mathcal{L}_{\text{cyc}}^B(\Theta_{G_{AB}}, \Theta_{G_{BA}}) = \mathbb{E}_{\mathbf{z}_b \sim \mathcal{X}_B} [\|\mathbf{z}_b - G_{AB}(G_{BA}(\mathbf{z}_b))\|_1].$$

- Adversarial learning

$$\mathcal{L}_{\text{adv}}^A(\Theta_{G_{BA}}, \Theta_{D_A}) = \mathbb{E}_{\mathbf{z}_a \sim \mathcal{X}_A} [D_A(\mathbf{z}_a)] - \mathbb{E}_{\mathbf{z}_b \sim \mathcal{X}_B} [D_A(G_{BA}(\mathbf{z}_b))],$$

$$\mathcal{L}_{\text{adv}}^B(\Theta_{G_{AB}}, \Theta_{D_B}) = \mathbb{E}_{\mathbf{z}_b \sim \mathcal{X}_B} [D_B(\mathbf{z}_b)] - \mathbb{E}_{\mathbf{z}_a \sim \mathcal{X}_A} [D_B(G_{AB}(\mathbf{z}_a))].$$