



中南大學
CENTRAL SOUTH UNIVERSITY



NEURAL INFORMATION
PROCESSING SYSTEMS



山东师范大学
SHANDONG NORMAL UNIVERSITY

Disentangled Cross-Modal Representation Learning with Enhanced Mutual Supervision

Lu Gao^{1*}, Wenlan Chen^{1*}, Daoyuan Wang¹, Fei Guo^{1†}, Cheng Liang^{2†}

¹School of Computer Science and Engineering, Central South University

²School of Information Science and Engineering, Shandong Normal University

CONTENT

01

Research Background

02

Methods

03

Experiments

04

Conclusion

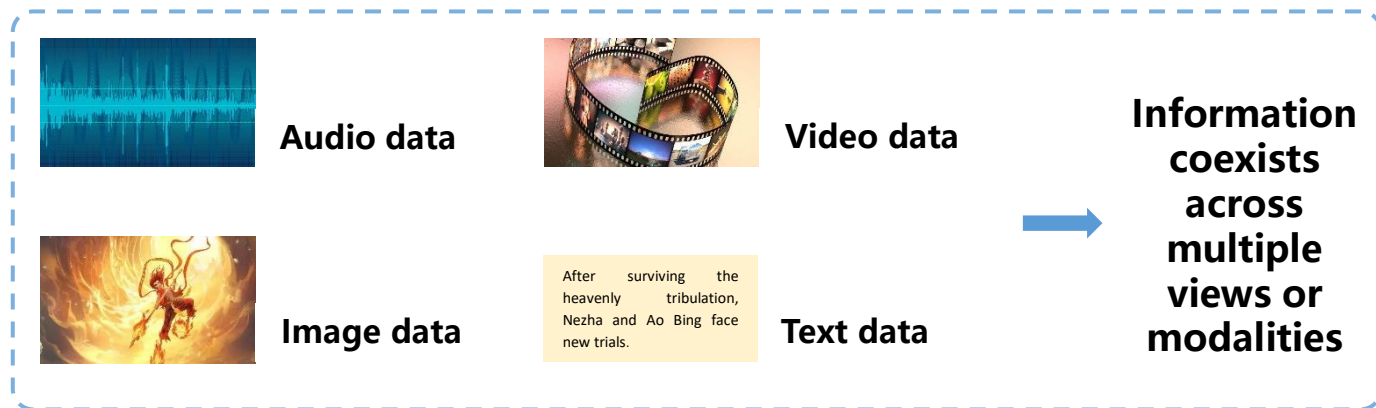
01

Research Background





Research Background



Goal

- Extract semantically aligned representations from heterogeneous modalities such as images and text.

Challenges



Existing multimodal VAE-based models often struggle to effectively align heterogeneous modalities.



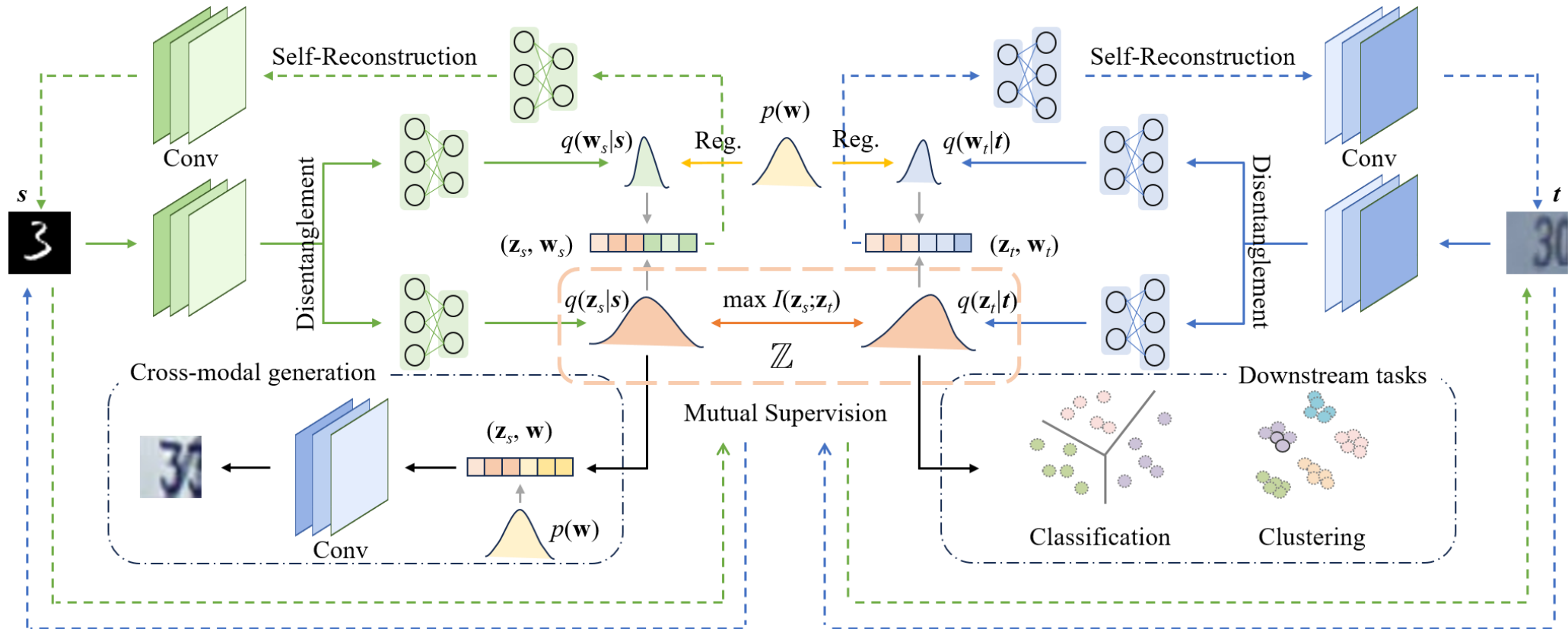
Many of these models also lack sufficient structural constraints to clearly separate the modality-specific and shared factors.

02

Methods



Methods



DCMEM

- Disentangle shared and modality-specific information, and enforce consistency through mutual supervision.
- Leverage the information bottleneck principle to promote semantic alignment and information complementarity across different modalities.

Methods

➤ Paired Setting

- Mutual Supervision
- Generative Process
- Inference Process

$$s \leftrightarrow z \leftrightarrow t$$

$$p_{\theta, \psi_z}(t, s, z_s, w_s) = p(t) p_{\psi_z}(z_s | t) p(w_s) p_{\theta}(s | z_s, w_s)$$

$$q_{\phi, \varphi}(z_s, w_s, t | s) = q_{\phi_z}(z_s | s) q_{\phi_w}(w_s | s) q_{\varphi}(t | z_s)$$



➤ ELBO Variational Framework

$$\log p_{\theta, \psi_z}(t, s) \geq \mathbb{E}_{q_{\phi}(z_s, w_s | s)} \left[\frac{q_{\varphi}(t | z_s)}{q_{\phi_z, \varphi}(t | s)} \log \frac{p_{\psi_z}(z_s | t) p(w_s) p_{\theta}(s | z_s, w_s)}{q_{\phi_z}(z_s | s) q_{\phi_w}(w_s | s) q_{\varphi}(t | z_s)} \right] + \log q_{\phi_z, \varphi}(t | s) + \log p(t).$$

Diagram labels for ELBO equation:

- Shared Representation (points to z_s)
- Prior (points to $p_{\psi_z}(z_s | t)$)
- Decoder (points to $p_{\theta}(s | z_s, w_s)$)
- Cross-modal Reconstruction (points to $q_{\phi_z, \varphi}(t | s)$)
- Encoder (points to $q_{\phi}(z_s, w_s | s)$)
- Variable Independence (points to the denominator terms $q_{\phi_z}(z_s | s)$ and $q_{\phi_w}(w_s | s)$)



➤ Structured Representation Learning

$$\max I(z_s; t; s) - I(z_s; w_s)$$

Shared Information

Redundant Information

➤ Paired Loss

$$\mathcal{L}_{Bi}(s, t)$$

➤ Shared Representations Alignment

$$\max I(z_s; z_t)$$

Methods

➤ Partially Missing Setting

- Mutual Supervision
- Generative Process
- Inference Process

$$s \leftrightarrow z \leftrightarrow t$$

$$p_{\theta, \psi_z}(s, z_s, w_s) = p_{u^t}(z_s) p(w_s) p_{\theta}(s | z_s, w_s)$$

$$q_{\phi}(z_s, w_s | s) = q_{\phi_z}(z_s | s) q_{\phi_w}(w_s | s)$$



➤ ELBO Variational Framework

$$\log p_{\theta, \psi_z}(s) \geq \mathbb{E}_{q_{\phi}(z_s, w_s | s)} \log \frac{p_{\theta}(s | z_s, w_s) p(w_s) p_{u^t}(z_s)}{q_{\phi_z}(z_s | s) q_{\phi_w}(w_s | s)},$$

Encoder
Decoder
Prior
Anchor Prior

Variable Independence



➤ VampPrior Anchor Prior

- Integral Form

$$p_{u^t}(z_s) = \int p(t) p_{\psi_z}(z_s | t) dt$$

- Monte Carlo Estimation

$$p_{u^t}(z_s) = \frac{1}{B} \sum_{i=1}^B p_{\psi_z}(z_s | u_i^t)$$



➤ Missing Loss

$$\mathcal{L}_s(s)$$



➤ Paired Loss

$$\mathcal{L}_{Bi}(s, t)$$



➤ Total Loss

$$\mathcal{L}(\mathcal{D}) = \sum_{s, t \in \mathcal{D}_{s, t}} \mathcal{L}_{Bi}(s, t) + \sum_{s \in \mathcal{D}_s} \mathcal{L}_s(s) + \sum_{t \in \mathcal{D}_t} \mathcal{L}_t(t).$$

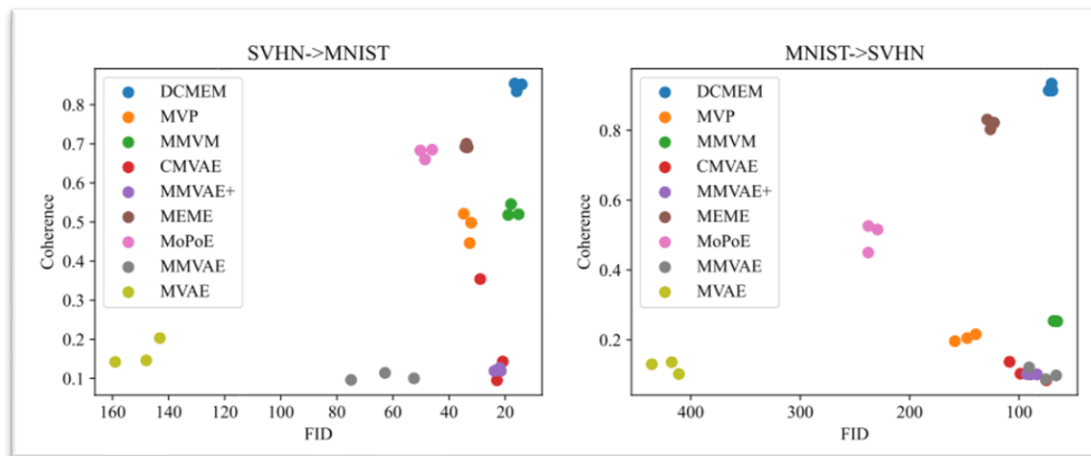
03

Experiments



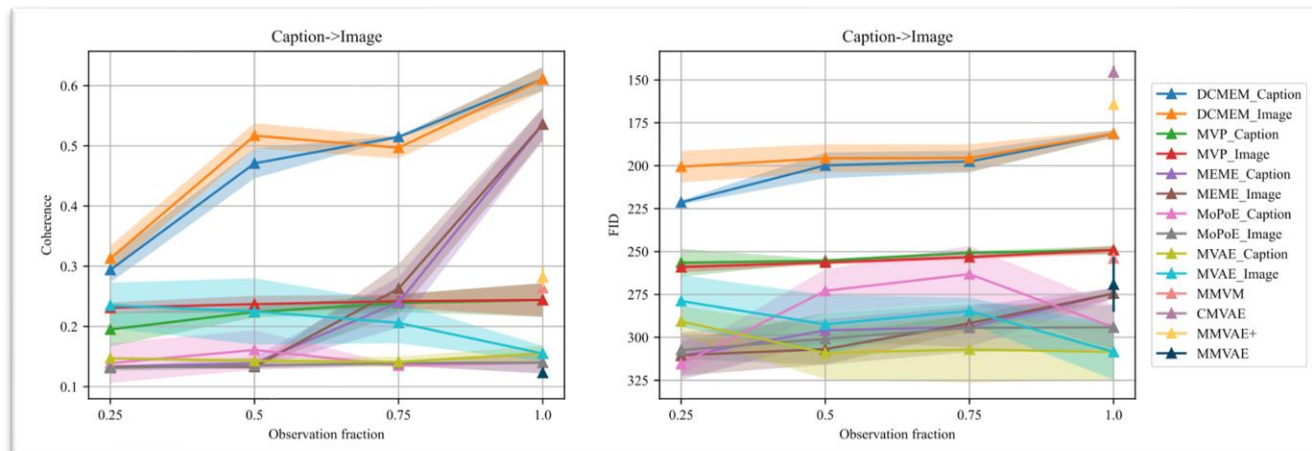


Generate performance



Paired Setting (MNIST-SVHN)

- **Task:** Cross-modal digit generation
- **Metrics:** FID (generation quality) and Coherence (generation consistency)
- DCMEM achieves **higher sample quality and better cross-modal consistency.**

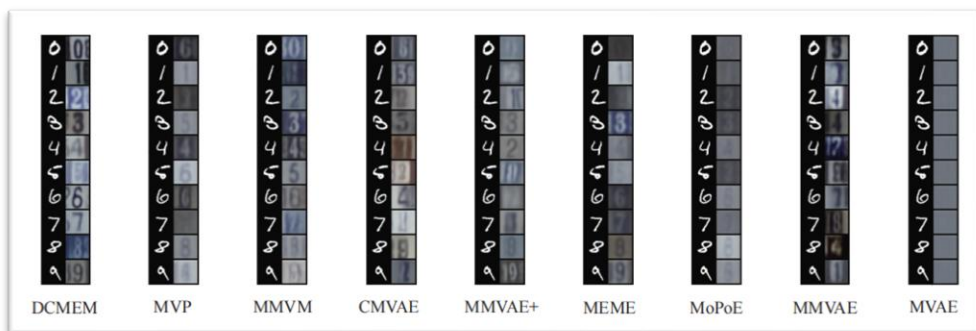


Partially Missing Setting (CUBICC)

- **Task:** Cross-modal image generation
- **Metrics:** FID and Coherence with varying proportions of observed paired data
- Achieves the **best generation quality and coherence** under all observation ratios.
- Maintains **strong cross-modal alignment** with limited paired data.

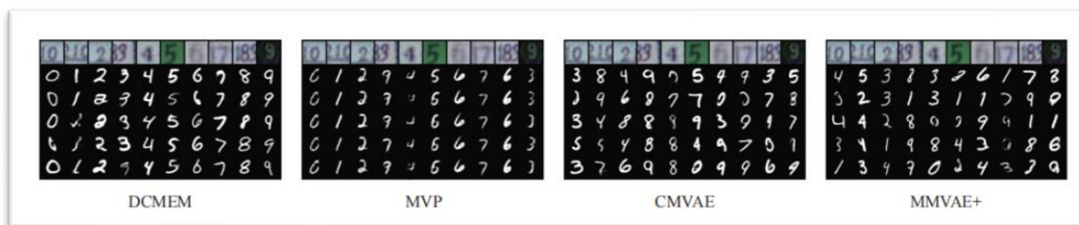


Generate performance



Qualitative Evaluation

- DCMEM produces **visually clear** and **semantically accurate** results that align well with the input modality.
- Compared to baseline methods, it better **preserves category information** and **avoids confusion or blurring** in generated samples.

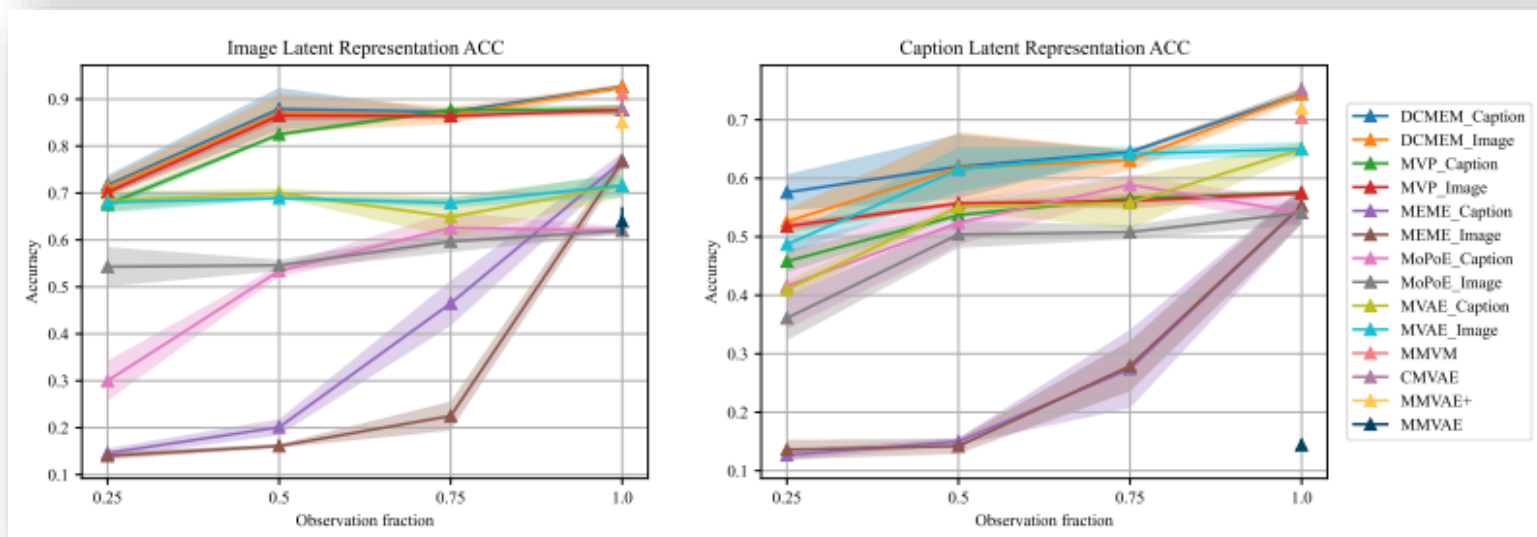
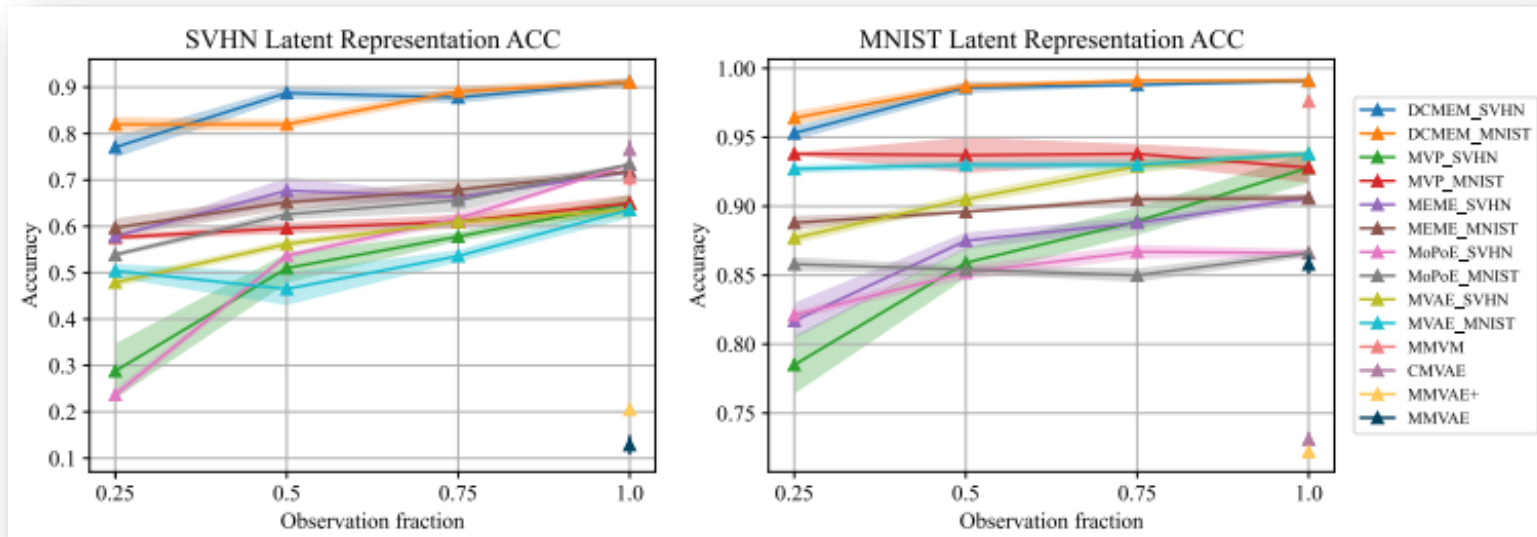


Diversity Evaluation

- DCMEM maintains **semantic consistency** while generating **diverse outputs** within each category.
- In contrast, MMVAE+ and CMVAE show **coupled category and modality-specific variables**, and MVP exhibits **limited variability**.



Classification performance



Result Analysis

- DCMEM: Achieves the **best accuracy across all observation ratios**, showing strong robustness and stability.
- MMVAE+: Performs poorly on MNIST-SVHN, as its shared latent space fails to **capture discriminative features**.
- MEME / MVP: Designed for incomplete modalities but **accuracy drops sharply with higher missing rates**.

Clustering performance

- DCMEM achieves consistently **superior performance across all three types of representations**.
- Traditional multimodal VAEs (e.g., MVAE, MMVAE, MoPoE) show **clear limitations in cross-modal representation alignment**.

- DCMEM achieves the **best performance** on most metrics, showing strong **feature alignment and discriminative capability**.

Methods	SVHN Representation			MNIST Representation			Joint Representation		
	ACC	NMI	ARI	ACC	NMI	ARI	ACC	NMI	ARI
MVAE	27.9	16.0	13.1	79.2	65.5	62.6	42.7	35.3	24.5
MMVAE	22.0	10.4	10.1	21.8	10.3	10.1	22.6	10.7	10.1
MoPoE	37.9	27.2	18.5	50.5	45.6	33.0	64.1	60.5	50.7
MEME	21.9	10.3	10.0	36.5	32.1	20.4	22.4	10.6	10.1
MMVAE+	23.9	11.4	11.1	21.3	10.4	10.0	22.9	11.9	10.8
CMVAE	42.2	36.3	25.4	28.1	15.9	14.5	32.3	19.5	15.4
MMVM	42.2	27.1	20.7	<u>88.1</u>	<u>82.1</u>	<u>80.4</u>	77.5	72.2	67.5
MVP	<u>53.6</u>	<u>38.7</u>	<u>30.1</u>	81.4	79.6	73.6	<u>84.8</u>	<u>76.4</u>	<u>70.6</u>
DCMEM	91.5	80.6	82.0	99.1	97.3	98.0	99.5	98.4	98.9

Methods	Image Representation			Caption Representation			Joint Representation		
	ACC	NMI	ARI	ACC	NMI	ARI	ACC	NMI	ARI
MVAE	26.2	12.4	7.5	18.1	2.4	0.9	38.7	26.8	18.0
MMVAE	23.1	12.1	6.1	14.5	1.3	0.1	15.8	1.5	0.2
MoPoE	33.4	17.6	11.5	43.5	27.1	19.9	40.8	30.4	20.2
MEME	44.8	43.4	28.4	36.3	29.5	18.6	19.8	4.8	2.1
MMVAE+	27.7	11.9	7.1	48.7	36.4	26.8	64.4	52.6	44.1
CMVAE	<u>67.7</u>	<u>58.3</u>	<u>47.4</u>	<u>65.1</u>	53.3	<u>42.7</u>	<u>73.7</u>	<u>67.4</u>	<u>57.2</u>
MMVM	58.9	56.9	44.5	23.9	9.4	5.4	66.8	67.0	55.5
MVP	64.1	53.8	41.8	48.5	34.4	26.1	61.1	55.6	44.0
DCMEM	86.9	77.4	72.4	69.7	<u>52.2</u>	44.2	86.3	76.8	71.5

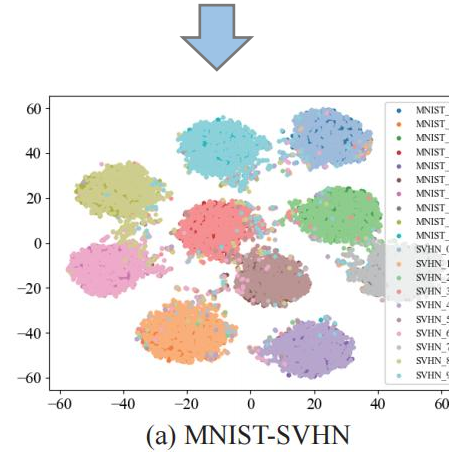
Representation Analysis

Class-level Semantic Alignment Analysis

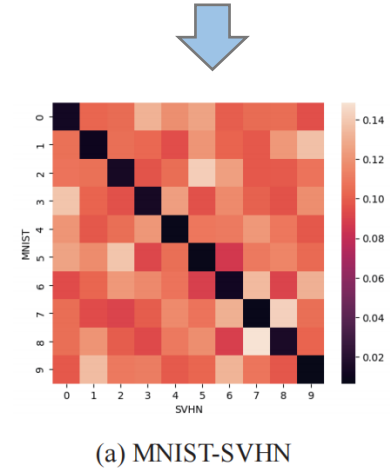
Objective: To examine whether different modalities maintain consistent semantic representations at the class level.

Method: Compute the distance matrix K_{ij} conditioned on class, where the diagonal entries represent intra-class distances and the off-diagonal entries represent inter-class distances.

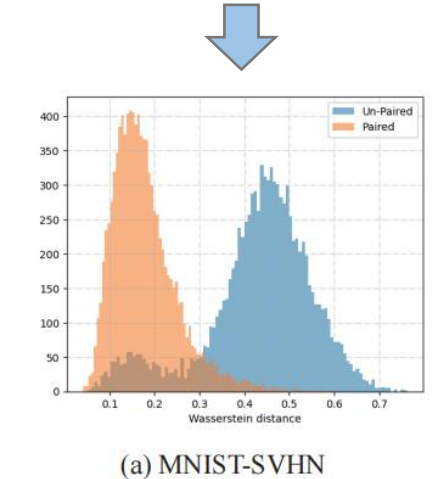
Shared Representation t-SNE Visualization



Class-level Semantic Alignment Analysis



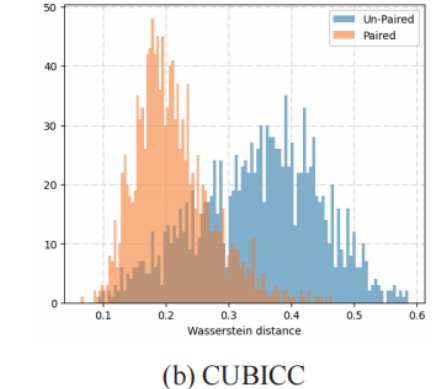
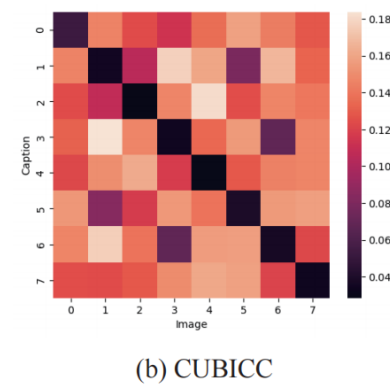
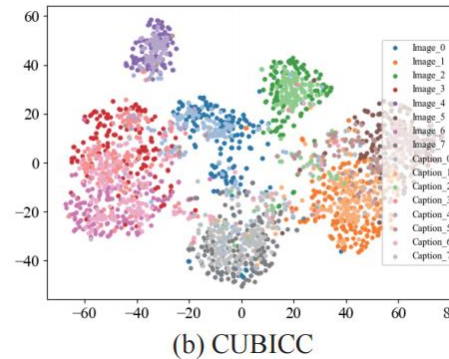
Sample-level Semantic Alignment Analysis



Sample-level Semantic Alignment Analysis

Objective: To determine whether semantically related cross-modal samples have closer latent distributions.

Method: Compute the 2-Wasserstein distance between all sample pairs and use histograms to distinguish paired from unpaired samples.



04

Conclusion



Conclusion



Contributions

01

- DCMEM separates shared and modality-specific information and enforces consistency through mutual supervision.

02

- DCMEM leverages the information bottleneck principle to promote compact and complementary feature encoding.

03

- DCMEM supports learning from both complete and partially missing multimodal data through a valid variational objective.

04

- Experimental results show superior performance on diverse tasks, including cross-modal generation, clustering, and classification.



中南大學
CENTRAL SOUTH UNIVERSITY



NEURAL INFORMATION
PROCESSING SYSTEMS



山东师范大学
SHANDONG NORMAL UNIVERSITY

Thanks!

Disentangled Cross-Modal Representation Learning with Enhanced Mutual Supervision