

Imbalanced View Contribution Evaluation and Refinement for Deep Incomplete Multi-View Clustering

Supplementary Material

Appendix A: Proof of Theorem 1

Theorem 1 (Boundedness of the metric I_ψ). *Assume that for any $S \subseteq \mathcal{V} \setminus \{v\}$, it holds that $\text{UOT}(\mathbf{p}, \mathbf{q}^{S \cup \{v\}}) \leq \text{UOT}(\mathbf{p}, \mathbf{q}^S)$. Equivalently, $E(S \cup \{v\}) \geq E(S)$. If there exists an effective improvement such that $E(\mathcal{V}) > E(\emptyset)$ (i.e., $\sum_v \psi^v > 0$), then the contribution-imbalance index satisfies*

$$I_\psi \in \left[0, 1 - \frac{1}{V}\right]. \quad (22)$$

The lower bound of 0 is achieved if and only if $\psi^1 = \dots = \psi^V = \frac{1}{V}$. In contrast, the upper bound $1 - \frac{1}{V}$ occurs in an extreme dominance scenario in which one view has $\psi^v = 1$ while all others are 0.

Proof. For any view set $\mathcal{V} = \{1, \dots, V\}$, define the UOT-based coalition utility as $E(S) = \exp(-\text{UOT}(\mathbf{p}, \mathbf{q}^{(S)}))$. For convenience, denote by \mathbf{q}^S the fused distribution constructed from the subset S .

According to the Shapley definition, the expected marginal contribution of view v is

$$\phi^v = \sum_{S \subseteq \mathcal{V} \setminus \{v\}} \frac{|S|!(V - |S| - 1)!}{V!} \times \left[\exp(-\text{UOT}(\mathbf{p}, \mathbf{q}^{S \cup \{v\}})) - \exp(-\text{UOT}(\mathbf{p}, \mathbf{q}^S)) \right]. \quad (23)$$

The sum of all Shapley values satisfies the standard efficiency property:

$$\begin{aligned} \sum_{v=1}^V \phi^v &= E(\mathcal{V}) - E(\emptyset) \\ &= \exp(-\text{UOT}(\mathbf{p}, \mathbf{q}^{\mathcal{V}})) - \exp(-\text{UOT}(\mathbf{p}, \mathbf{q}^{\emptyset})). \end{aligned} \quad (24)$$

Based on this, define the normalized contributions (relative weights) by

$$\psi^v = \frac{\phi^v}{\sum_{u=1}^V \phi^u}, \quad v = 1, \dots, V. \quad (25)$$

The I_ψ is then used to quantify the ‘‘contribution imbalance’’:

$$I_\psi = \frac{1}{2V^2 \bar{\psi}} \sum_{i=1}^V \sum_{j=1}^V |\psi^i - \psi^j|, \quad \bar{\psi} = \frac{1}{V} \sum_{v=1}^V \psi^v. \quad (26)$$

Basic assumption For any $S \subseteq \mathcal{V} \setminus \{v\}$,

$$\text{UOT}(\mathbf{p}, \mathbf{q}^{S \cup \{v\}}) \leq \text{UOT}(\mathbf{p}, \mathbf{q}^S). \quad (27)$$

Equivalently, $E(S \cup \{v\}) \geq E(S)$. *Intuition:* adding an additional view for fusion cannot worsen the OT distance to the target distribution; thus the utility never decreases. (Under the entropy-regularized/unbalanced OT fusion model used in our framework, this is a standard and empirically valid assumption.)

(1) Non-negativity. For any S , $E(S \cup \{v\}) - E(S) \geq 0$. Substituting into (23) gives $\phi^v \geq 0$. If $E(\mathcal{V}) > E(\emptyset)$, then from (24) we have $\sum_u \phi^u > 0$, hence by (25) the normalized weights satisfy

$$\psi^v \geq 0, \quad \sum_{v=1}^V \psi^v = 1 \implies \bar{\psi} = \frac{1}{V}. \quad (28)$$

(2) Simplification of the I_ψ index. Using (28), Eq. (26) reduces to

$$\begin{aligned} I_\psi &= \frac{1}{2V} \sum_{i=1}^V \sum_{j=1}^V |\psi^i - \psi^j| \frac{1}{2V} F(\psi), \\ \psi &\in \Delta_V \stackrel{\text{def}}{=} \{\psi \in \mathbb{R}_{\geq 0}^V : \sum_v \psi^v = 1\}. \end{aligned} \quad (29)$$

(3) Lower bound and necessary-and-sufficient condition. If $\psi_i = \frac{1}{V}$ for all i , then all pairwise differences vanish and $I_\psi = 0$. Conversely, if $I_\psi = 0$, then by (29) we have $|\psi^i - \psi^j| \equiv 0$, hence all ψ^i are equal; combined with $\sum_v \psi^v = 1$ we obtain $\psi^i = \frac{1}{V}$.

(4) Upper bound and attainability. The function $F(\psi) = \sum_{i,j} |\psi^i - \psi^j|$ is a piecewise-linear convex function of ψ , and its maximum over the convex set Δ_V is attained at an extreme point. At the extreme point $\psi^* = (1, 0, \dots, 0)$,

$$F(\psi^*) = \sum_{j=2}^V |1 - 0| + \sum_{i=2}^V |0 - 1| = 2(V - 1). \quad (30)$$

Substituting into (29) gives

$$I_\psi \leq \frac{1}{2V} \cdot 2(V - 1) = 1 - \frac{1}{V}, \quad (31)$$

with equality achieved at the extreme point. This completes the proof. \square

What’s more, the detailed procedure of the algorithm is outlined in Algorithm 1. The formula numbers in the algorithm correspond to those in the paper.

Algorithm 1 The proposed ICER

```
1: Input: Incomplete multi-view data  $\{X^v\}_{v=1}^V$ , indicator matrix  $A$ , number of clusters  $K$ , autoencoders  $\{E^v, D^v\}_{v=1}^V$ , trade-off parameters  $\lambda_1, \lambda_2$ , curriculum parameters  $(T_c, \varepsilon)$ .
2: Output: Clustering results.
3: for  $t = 1$  to  $E$  do
4:   // View distribution estimation
5:   for  $v = 1$  to  $V$  do
6:     Extract embeddings  $Z^v = E^v(\mathbf{X}^v)$ ;
7:     Compute soft assignment  $\mathbf{Q}^v$  by Eq. (3);
8:   end for
9:   Construct fused distribution  $\mathbf{q}^{(S)}$  with Eq. (4)–(5);
10:  Compute global target distribution  $\mathbf{p}$  with Eq. (6);
11:  // View contribution evaluation via Shapley value
12:  for each coalition  $\mathbf{S} \subseteq \{1, \dots, V\}$  do
13:    Form coalition-wise fused distribution  $\mathbf{q}^{(S)}$  using Eq. (4)–(5);
14:    Compute UOT distance  $\text{UOT}(\mathbf{p}, \mathbf{q}^{(S)})$  by Eq. (7);
15:    Compute coalition utility  $E(\mathbf{S})$  with Eq. (8);
16:  end for
17:  for  $v = 1$  to  $V$  do
18:    Compute Shapley value  $\phi^v$  using Eq. (12);
19:  end for
20:  Normalize contributions  $\psi^v$  (optionally compute  $I_\psi$  by Eq. (14));
21:  // Contribution-balanced enhancement (VAEL)
22:  for  $v = 1$  to  $V$  do
23:    Compute curriculum factor  $\gamma_v(t)$  with Eq. (18);
24:  end for
25:  // Loss computation and parameter update
26:  Compute reconstruction loss  $\mathcal{L}_{\text{rec}}$  with Eq. (1);
27:  Compute cross-view consistency loss  $\mathcal{L}_{\text{ccl}}$  with Eq. (9);
28:  Compute clustering loss  $\mathcal{L}_c$  with Eq. (17);
29:  Form overall objective  $\mathcal{L} = \mathcal{L}_{\text{rec}} + \lambda_1 \mathcal{L}_{\text{ccl}} + \lambda_2 \mathcal{L}_c$  as in Eq. (18);
30:  for  $v = 1$  to  $V$  do
31:    Compute view-specific weighted gradient  $g_t^{(v)} = \gamma_v(t) \nabla_{\theta_t^{(v)}} L^{(v)}$  as in Eq. (21);
32:    Update parameters  $\theta_{t+1}^{(v)} = \theta_t^{(v)} - \eta g_t^{(v)}$  by Eq. (20);
33:  end for
34: end for
35: // Clustering
36: Obtain final fused distribution  $\mathbf{P} = [p_{ij}]$ ;
37: Assign cluster labels  $\hat{y}_i = \arg \max_j p_{ij}$ .
```

Appendix B: Proof of Remark 2

Remark 1. Let the unnormalized curriculum weight be defined as

$$\gamma_v(t) = \left(\frac{t}{T_c}\right) [(1 - \psi^v)^2 + \varepsilon], \quad v = 1, \dots, V, \quad (32)$$

where $0 \leq \psi^v \leq 1$ denotes the contribution of view v , $(1 - \psi^v)^2$ controls the amplification strength for weak views, $\varepsilon > 0$ is a small constant preventing numerical underflow, and $T_c > 0$ is the curriculum-phase length. Define the temporal factor

$$k(t) \triangleq \left(\frac{t}{T_c}\right), \quad (33)$$

then

$$\gamma_v(t) = k(t) A_v, \quad A_v \triangleq (1 - \psi^v)^2 + \varepsilon. \quad (34)$$

Under this definition, the unnormalized weight $\gamma_v(t)$ satisfies the following properties during the curriculum phase $t \in [0, T_c]$.

(1) Monotonicity. Taking the partial derivative of (32) with respect to ψ^v ,

$$\frac{\partial \gamma_v(t)}{\partial \psi^v} = k(t) \frac{\partial A_v}{\partial \psi^v} = -2k(t)(1 - \psi^v) \leq 0. \quad (35)$$

Thus, as the contribution ψ^v of view v increases, its corresponding curriculum weight $\gamma_v(t)$ is monotonically non-increasing. For any $u \neq v$, we have

$$\frac{\partial \gamma_v(t)}{\partial \psi^u} = 0. \quad (36)$$

This shows that, under the unnormalized design, the weight of each view depends solely on its own contribution: strong views (large ψ^v) receive smaller additional weight, while weak views (small ψ^v) receive larger emphasis.

(2) Boundedness. Since $0 \leq \psi^v \leq 1$, we have

$$0 \leq (1 - \psi^v)^2 \leq 1. \quad (37)$$

Hence

$$A_v = (1 - \psi^v)^2 + \varepsilon \in [\varepsilon, 1 + \varepsilon]. \quad (38)$$

For any fixed t ,

$$k(t)\varepsilon \leq \gamma_v(t) \leq k(t)(1 + \varepsilon). \quad (39)$$

If $t \in [0, T_c]$, then $k(t) \in [0, 1]$, and thus

$$0 \leq \gamma_v(t) \leq 1 + \varepsilon. \quad (40)$$

This shows that the unnormalized curriculum weight remains strictly bounded throughout the entire training process, avoiding numerical explosion or gradient blow-up.

(3) Lipschitz stability. Consider the vector-valued function

$$\gamma(t, \psi) = (\gamma_1(t), \dots, \gamma_V(t))^\top, \quad (41)$$

whose Jacobian matrix is

$$J(\psi) = \frac{\partial \gamma(t, \psi)}{\partial \psi} = \text{diag}\left(\frac{\partial \gamma_1}{\partial \psi^1}, \dots, \frac{\partial \gamma_V}{\partial \psi^V}\right), \quad (42)$$

where each diagonal element satisfies

$$\left| \frac{\partial \gamma_v(t)}{\partial \psi^v} \right| = 2k(t)(1 - \psi^v) \leq 2k(t). \quad (43)$$

Under any consistent matrix norm, we have

$$\|J(\psi)\| \leq 2k(t). \quad (44)$$

By the mean-value theorem, for any $\psi, \psi' \in [0, 1]^V$, there exists ξ lying between them such that

$$\gamma(t, \psi) - \gamma(t, \psi') = J(\xi)(\psi - \psi'). \quad (45)$$

Therefore,

$$\|\gamma(t, \psi) - \gamma(t, \psi')\| \leq \|J(\xi)\| \cdot \|\psi - \psi'\| \leq 2k(t) \|\psi - \psi'\|. \quad (46)$$

Since $k(t) \leq 1$, we may choose a uniform Lipschitz constant $L = 2$, yielding

$$\forall t \geq 0, \forall \psi, \psi' \in [0, 1]^V, \quad (47)$$

$$\|\gamma(t, \psi) - \gamma(t, \psi')\| \leq L \|\psi - \psi'\|.$$

Hence, during the curriculum phase, the unnormalized weights $\gamma(t, \psi)$ are globally Lipschitz continuous with respect to the contribution vector ψ , with Lipschitz constant $L = 2$.

Appendix C: Dataset Description

To realistically emulate both random and imbalanced missing conditions in multi-view datasets, we design an incomplete data generation mechanism that jointly controls the global missing level and view-wise heterogeneity. Suppose the dataset contains M views and N samples, forming an $M \times N$ observation matrix \mathbf{X} . A global missing rate r is first specified, yielding a total number of missing entries equal to MNr .

Heterogeneous Missing Distribution Across Views. Following the imbalanced missing settings commonly used in multi-view clustering, we introduce a view-specific scaling factor α_v for each view v , sampled under the constraint

$$\sum_{v=1}^M \alpha_v = M, \quad (48)$$

ensuring that the average missing rate remains r while allowing each view to exhibit distinct levels of incompleteness. The missing rate of view v is then defined as

$$r_v = \alpha_v \cdot r, \quad (49)$$

which produces heterogeneous missing patterns across views—some views become “strong” (low missingness) while others become “weak” (high missingness).

At Least One View Must Be Observed. To avoid generating completely missing samples, we enforce that each sample must remain observable in at least one view. This introduces a theoretical upper bound on the missing rate. In the worst case, each sample may lose at most $M - 1$ views; therefore, the maximal total number of missing entries is

$$N(M - 1). \quad (50)$$

Since the missing mechanism defines MNr missing entries in total, the constraint

$$MNr \leq N(M - 1) \quad (51)$$

must hold, which yields the upper bound:

$$r \leq 1 - \frac{1}{M}. \quad (52)$$

This bound is tight, and the maximum is achieved when every sample has exactly $M - 1$ missing views. For example: 2 views: $r_{\max} = 0.5$, 3 views: $r_{\max} = 0.667$, 4 views: $r_{\max} = 0.75$

Combining global control and view-specific heterogeneity, the incomplete matrix is generated by randomly masking MNr_v entries in each view according to its assigned missing rate r_v , while ensuring that each sample retains at least one observed view. This procedure allows us to simulate realistic missing patterns with both controlled randomness and imbalanced view-wise incompleteness, closely reflecting real-world multi-view scenarios.

Appendix D: Experiments

Dataset information

We provide an introduction to the five datasets used in the paper.

- **Synthetic3D[6]:** This dataset is constructed from several 3D manifold structures such as helices, spheres, and tori. Multiple views are generated via nonlinear projections or noise injection. It contains 600 samples with 3 views.
- **Reuters:** Reuters¹ is a well-known multi-view text dataset comprising 18,758 documents from 6 topic categories. Its multi-view structure is created through five

¹<https://archive.ics.uci.edu/ml/datasets.html>

language versions of each document, including English, German, French, Italian, and Spanish, naturally forming a multilingual multi-view setting.

- **Caltech101**²: Caltech101 contains images from multiple object categories and serves as a classic benchmark in visual recognition. Its multi-view version is constructed using diverse feature extraction methods, such as Gabor features, Wavelet features, HOG descriptors, and GIST global features, resulting in 6 complementary views with a total of 1,474 samples.
- **Wikipedia**: The Wikipedia ³ dataset consists of textual descriptions of Wikipedia articles and their corresponding images. Each sample provides at least two modalities: a textual view (e.g., TF-IDF or LDA features) and a visual view (e.g., SIFT or CNN features). The dataset contains 2,866 samples in total.
- **Animal**[7]: The Animal dataset contains 10,158 samples over 50 animal categories. Two deep features, DECAF[5] and VGG19[8], are extracted to form two complementary views, making it widely used for multi-view representation learning and clustering.

Compared Methods

We compare ICER with seven representative incomplete multi-view clustering methods, and their details are as follows:

- **CPSPAN**[3]: a cross-view partial sample and prototype alignment network that leverages pair-observed data and prototype alignment to address representation inconsistency and prototype bias in incomplete multi-view clustering.
- **RPCIC**[10]: a robust prototype completion framework that leverages cross-view contrastive learning and robust prototype-level alignment to mitigate noisy prototype correspondence in incomplete multi-view clustering.
- **GHICMC**[1]: a novel incomplete multi-view clustering method that integrates view-specific GCNs, hierarchical information transfer, and contrastive learning for joint representation learning and clustering.
- **PMIMC**[11]: a prototype matching learning framework that leverages relational consistency and robust prototype contrastive learning to address prototype misalignment and imputation instability in incomplete multi-view clustering.
- **FREECSL**[2]: an IMVC framework that jointly performs imputation and alignment to learn consensus semantics by leveraging cross-view consistency and consensus prototype learning.
- **NBIMVC**[9]: a neighbor-based completion framework that leverages topological relationships and nearest-neighbor information to impute missing views and en-

force cross-view alignment for incomplete multi-view clustering.

- **BURG**[4]: a distribution-based dual-consistency recovery framework that performs cross-view distribution transfer and intra-/inter-view alignment to recover missing views in incomplete multi-view clustering.

Ablation study

The table 1 ablation results on Wikipedia and Reuters further validate the roles of the three loss components. Using only \mathcal{L}_{rec} yields relatively limited performance, as it fails to handle cross-view inconsistency under imbalanced missingness. Incorporating \mathcal{L}_{ccl} consistently improves performance across all metrics, indicating its effectiveness in enhancing cross-view semantic alignment. In contrast, applying \mathcal{L}_c alone leads to noticeable performance degradation, suggesting that clustering supervision is unreliable without well-aligned representations. The full model achieves the best results on both datasets, demonstrating that the three loss components are complementary and jointly contribute to robust clustering.

Table 1. Ablation study on Wikipedia and Reuters. We report ACC, NMI, and ARI as evaluation metrics. (The missing rate is 0.3)

Datasets	\mathcal{L}_{rec}	\mathcal{L}_{ccl}	\mathcal{L}_c	ACC	NMI	ARI
Wikipedia	✓			0.4833	0.3510	0.2587
	✓	✓		0.4853	0.3713	0.3019
	✓		✓	0.2394	0.1107	0.0641
	✓	✓	✓	0.4927	0.3796	0.3073
Reuters	✓			0.3292	0.0975	0.0650
	✓	✓		0.4213	0.1571	0.1263
	✓		✓	0.2767	0.0570	0.0383
	✓	✓	✓	0.4325	0.1849	0.1444

T-SNE Visualization

Fig. 1 presents the t-SNE visualizations of feature representations on the Reuters dataset with a missing rate of 0.1. The raw features Fig. (1a) exhibit poor separability, characterized by severe overlap across clusters. After training, the learned representations Fig. (1b) demonstrate significantly improved structure, forming compact and well-separated clusters. This enhanced cluster compactness and inter-cluster separability align with the notable improvement in clustering accuracy, highlighting the effectiveness of the proposed approach in capturing discriminative and robust representations under missing data.

Analysis of VACL ablation experiment

In addition, we further evaluate the effect of applying the VACL module within the Contribution-Balanced Enhancement Strategy on the Wikipedia dataset under a missing rate

²http://www.vision.caltech.edu/Image_Datasets/Caltech101/

³<http://www.svcl.ucsd.edu/projects/crossmodal/>

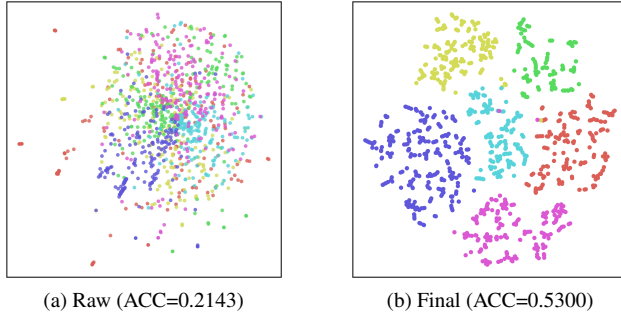


Figure 1. t-SNE visualization of clustering results for different feature spaces on the Reuters dataset (The missing rate is 0.1).

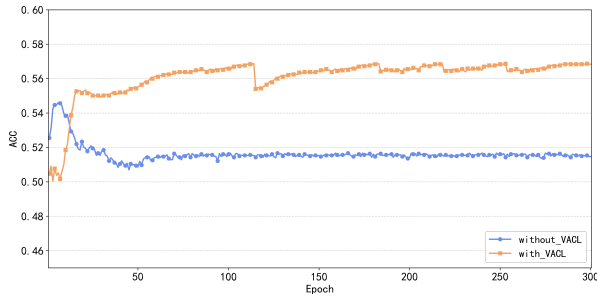


Figure 2. ACC comparison between models with and without VACL.

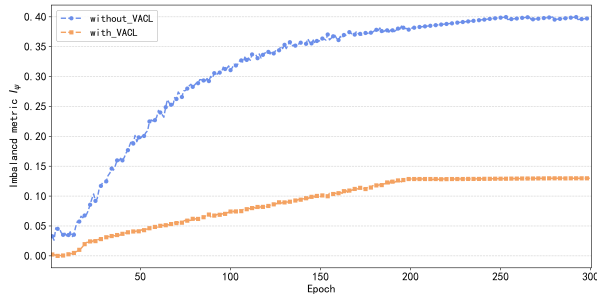


Figure 3. Imbalance metric I_{ψ} with and without VACL.

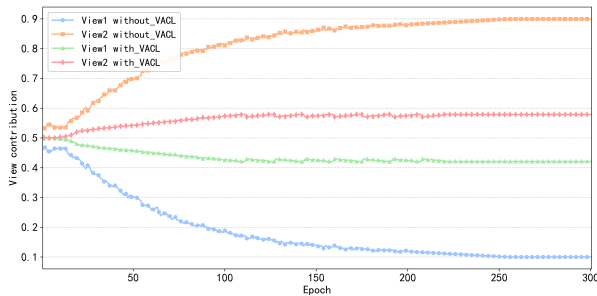


Figure 4. View contribution comparison with and without VACL.

of 0.1. The comparative results are shown in the figure below.

Fig. 2 shows that applying VACL leads to consistently higher ACC throughout the training process. At the early stage, the model equipped with VACL exhibits relatively larger fluctuations. This is because VACL initially emphasizes the view with higher contribution, causing the model to rapidly adjust its representation space, which introduces short-term instability. As training proceeds, the ACC steadily improves and becomes more stable, demonstrating that VACL facilitates more effective multi-view collaboration and helps the model escape suboptimal early-stage dynamics. In contrast, without VACL, the ACC quickly saturates at a lower level since the model is dominated by the strong view and fails to exploit complementary information from the weak view.

Fig. 3 further illustrates the evolution of the imbalance metric I_{ψ} . Without VACL, I_{ψ} increases sharply as training progresses, indicating that the model’s reliance on the strong view intensifies and the contribution discrepancy between views becomes increasingly severe. This growing imbalance suggests degraded cross-view cooperation and a tendency to overfit toward the dominant view. In contrast, with VACL, I_{ψ} grows very slowly and remains at a significantly lower level, confirming that VACL effectively suppresses imbalance amplification and maintains healthier multi-view collaboration during optimization.

Fig. 4 presents the normalized contribution curves of the two views. Without VACL, the contribution of View 1 rapidly decreases while View 2 overwhelmingly dominates, revealing a severe contribution-collapse phenomenon. When VACL is applied, the contributions of both views remain more balanced and evolve smoothly. The weak view is gradually reinforced while the strong view is prevented from overwhelming the training process. These results demonstrate that VACL successfully redistributes gradient focus, stabilizes the collaboration dynamics among views, and ensures that complementary information from multiple views can be effectively integrated.

References

- [1] Guoqing Chao, Kaixin Xu, Xijiong Xie, and Yongyong Chen. Global graph propagation with hierarchical information transfer for incomplete contrastive multi-view clustering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 15713–15721, 2025. 4
- [2] Yuzhuo Dai, Jiaqi Jin, Zhibin Dong, Siwei Wang, Xinwang Liu, En Zhu, Xihong Yang, Xinbiao Gan, and Yu Feng. Imputation-free and alignment-free: Incomplete multi-view clustering driven by consensus semantic learning. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 5071–5081, 2025. 4
- [3] Jiaqi Jin, Siwei Wang, Zhibin Dong, Xinwang Liu, and En Zhu. Deep incomplete multi-view clustering with cross-view partial sample and prototype alignment. In *2023 IEEE/CVF*

Conference on Computer Vision and Pattern Recognition (CVPR), pages 11600–11609, 2023. 4

- [4] Jiaqi Jin, Siwei Wang, Zhibin Dong, Xihong Yang, Xinwang Liu, En Zhu, and Kunlun He. Deep incomplete multi-view clustering with distribution dual-consistency recovery guidance. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1016–1026, 2025. 4
- [5] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017. 4
- [6] Abhishek Kumar, Piyush Rai, and Hal Daume. Co-regularized multi-view spectral clustering. *Advances in neural information processing systems*, 24, 2011. 3
- [7] Christoph H Lampert, Hannes Nickisch, and Stefan Harmeling. Attribute-based classification for zero-shot visual object categorization. *IEEE transactions on pattern analysis and machine intelligence*, 36(3):453–465, 2013. 4
- [8] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *3rd International Conference on Learning Representations*, pages 1–14. 4
- [9] Wenbiao Yan, Jihua Zhu, Yiyang Zhou, Jinqian Chen, Haozhe Cheng, Kun Yue, and Qinghai Zheng. Neighbor-based completion for addressing incomplete multiview clustering. *IEEE Transactions on Neural Networks and Learning Systems*, 36(8):15374–15384, 2025. 4
- [10] Honglin Yuan, Shiyun Lai, Xingfeng Li, Jian Dai, Yuan Sun, and Zhenwen Ren. Robust prototype completion for incomplete multi-view clustering. In *Proceedings of the 32nd ACM International Conference on Multimedia*, page 10402–10411, New York, NY, USA, 2024. Association for Computing Machinery. 4
- [11] Honglin Yuan, Yuan Sun, Fei Zhou, Jing Wen, Shihua Yuan, Xiaojian You, and Zhenwen Ren. Prototype matching learning for incomplete multi-view clustering. *IEEE Transactions on Image Processing*, 34:828–841, 2025. 4