# SELF-PACED ENCODING WITH ADAPTIVE GRAPH REGULARIZATION FOR MULTI-VIEW CLUSTERING

#### **Anonymous authors**

000

001

002003004

010 011

012

013

014

015

016

017

018

019

020

021

022

024

025

026

027 028 029

031

033

034

037

040

041

042

043

044

046 047

048

Paper under double-blind review

#### **ABSTRACT**

Multi-view graph clustering is a powerful technique for learning discriminative node representations by integrating complementary information from diverse views. However, existing methods often suffer from rigid fusion schemes, ignore sample difficulty during training, and struggle to capture both global semantics and local structures through graph-based regularization. To address these issues, we propose SPEAG, a novel framework for Self-Paced Encoding with Adaptive Graph Regularization. SPEAG combines view-specific graph autoencoders with a unified learning objective that incorporates self-paced training, adaptive view fusion, and structure-aware regularization. Specifically, a self-paced neighborhood expansion strategy is introduced, where the k-nearest neighbor graph is gradually densified to learn from easy instances first and hard ones later. Meanwhile, each view's embedding is adaptively weighted based on its importance, and a fusion representation is formed for global consistency. To encourage distributional alignment and enhance cluster compactness, SPEAG integrates a Maximum Mean Discrepancy (MMD) loss across views and a self-supervised clustering objective based on soft assignment refinement. Extensive experiments on real-world datasets demonstrate that SPEAG achieves superior clustering accuracy and robustness compared to existing multi-view graph clustering methods.

# 1 Introduction

Multi-view clustering (MVC) Liu et al. (2022); Fang et al. (2023) seeks to partition unlabeled data by jointly exploiting all views, and recent deep MVC advances leverage powerful neural representations. Representative methods include CoMSC Liu et al. (2021) (feature decomposition for robust representation) and DUA-Nets Geng et al. (2021) (uncertainty-aware view weighting), while CMRL Zheng et al. (2023) and SCMRL Zhou et al. (2023) further explore complementarity and semantic consensus via low-rank tensors and attention. However, many approaches emphasize view-specific features while underutilizing instance–instance relations that are crucial for clustering. Moreover, anchor-based methods reduce computation but often distort local structures, and GNN-based models (e.g., MGCN, MVGRL Kang et al. (2020); Yang & Zhu (2024); Jiang et al. (2025)) frequently fuse views heuristically (e.g., averaging) Chen et al. (2025b) and decouple representation learning from clustering; they also depend on static kNN or precomputed similarities Chen et al. (2025a) that are non-adaptive during training and sensitive to noise/outliers.

To address these challenges, we propose a novel Self-Paced and Enhanced Adaptive Graph encoding framework, dubbed SPEAG, for unsupervised multi-view graph clustering. SPEAG introduces several key innovations that are carefully integrated into a unified learning framework:

- Self-paced graph encoding with Laplacian regularization: Instead of a fixed k-NN graph, SPEAG updates each view's adjacency via an encoder-decoder; k increases during training for stable warm-up then global structure, while Laplacian terms preserve local geometry.
- Self-weighted fusion with distribution alignment: Instance-level view weights are learned
  jointly with embeddings; an MMD loss aligns fused and per-view representations, downweighting unreliable views and mitigating semantic drift across modalities for more robust
  multi-view consistency.

• Unified self-supervised clustering, end-to-end: A soft-label clustering loss tightens clusters and feeds back to the encoder; fusion, embedding, and clustering are optimized jointly, enabling mutual reinforcement, efficient cross-feedback, and extensibility within a single training pipeline.

#### 2 The Proposed Method

In this section, we propose a novel multi-view clustering via Self-Paced Encoding with Adaptive Graph regularization, whose crucial details are elaborated.

#### 2.1 NOTATIONS

Given V views  $\{X^{(v)}\}_{v=1}^V$  with  $X^{(v)} \in \mathbb{R}^{N \times d_v}$  and K clusters, where N is the number of samples and  $d_v$  the dimension of view v, we aim to learn a unified embedding  $H \in \mathbb{R}^{N \times d_h}$ . SPEAG combines view-specific graph autoencoders with a unified objective featuring self-paced training, adaptive view fusion, and structure-aware regularization. For each view we obtain a latent  $Z^{(v)} \in \mathbb{R}^{N \times d_z}$ ; pairwise distances are  $D^{(v)}$ , similarities  $W^{(v)}$ , their symmetrized form  $A^{(v)}$ , and normalized Laplacian  $\hat{L}^{(v)}$ . We fuse the view latents into a global feature R and produce the consensus embedding H, with  $w^{(v)}$  denoting the adaptive reliability weight of view v.

#### 2.2 WITHIN-VIEW RECONSTRUCTION

#### 2.2.1 Graph Embedding Autoencoder

We employ a graph convolutional autoencoder (GCAE) that ingests the feature matrix and a similarity graph per view. For view v, we first compute pairwise Euclidean distances  $D_{ij}^{(v)} = \|X_i^{(v)} - X_j^{(v)}\|_2^2$  and convert them to similarities via a Gaussian kernel  $W^{(v)} = \exp\left(-D^{(v)}/(2\sigma^2)\right)$ , where  $\sigma$  is the bandwidth controlling decay with distance. Before encoding, we symmetrize and normalize the graph:  $A^{(v)} = \frac{1}{2} (W^{(v)} + W^{(v)^{\top}})$  and  $\hat{L}^{(v)} = I - (\tilde{D}^{(v)})^{-1/2} A^{(v)} (\tilde{D}^{(v)})^{-1/2}$  with  $\tilde{D}_{ii}^{(v)} = \sum_j A_{ij}^{(v)}$ . This normalization preserves spectral properties and stabilizes message passing in the GCAE.

Feeding  $X^{(v)}$  and  $\hat{L}^{(v)}$  into the encoder yields the latent  $Z^{(v)} = \hat{L}^{(v)} \phi (\hat{L}^{(v)} X^{(v)} W_1^{(v)}) W_2^{(v)}$ , where  $W_1^{(v)}, W_2^{(v)}$  are layer parameters and  $\phi$  is the nonlinearity. We then reconstruct a row-stochastic similarity from latent distances  $\hat{D}_{ij}^{(v)} = \|Z_i^{(v)} - Z_j^{(v)}\|_2^2$  using a per-row softmax  $\bar{W}_{ij}^{(v)} = \exp(-\hat{D}_{ij}^{(v)}) / \sum_{j'=1}^N \exp(-\hat{D}_{ij'}^{(v)})$ . Finally, reconstruction fidelity is measured by the KL divergence  $L_{\rm re} = D_{\rm KL}(W^{(v)} \| \bar{W}^{(v)}) = \frac{1}{N} \sum_{i,j=1}^N W_{ij}^{(v)} \log(W_{ij}^{(v)} / \bar{W}_{ij}^{(v)})$ , which encourages  $Z^{(v)}$  to encode the view's graph structure.

#### 2.2.2 GRAPH LAPLACIAN REGULARIZATION

In our approach, we incorporate not only the graph structural information but also complementary feature information derived directly from the samples. Under the manifold assumption, if two data points are close in the original high-dimensional space Cai et al. (2008); Wen et al. (2018), their corresponding representations in the learned low-dimensional latent space should also remain close. Concretely, the consensus similarity among different views should be preserved after dimensionality reduction.

To enforce this, we introduce a graph regularization term formulated as follows:

$$L_{gls} = \sum_{i,j=1}^{N} A_{ij}^{v} ||z_{i}^{v} - z_{j}^{v}||_{2}^{2} = tr((Z^{v})^{T} L^{v} Z^{v})$$
(1)

where  $A_{ij}^v$  denotes the similarity between samples i and j in the original space of the v-th view,  $z_i^v$  is the latent representation of sample i in view v, and  $L^v = D^v - A^v$  is the graph Laplacian matrix

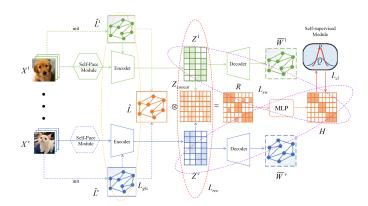


Figure 1: SPEAG adaptively selects samples for graph construction, encodes with GCN, reconstructs affinities, and fuses multi-view features to enhance clustering via multiple losses.

for view v, with  $D^v$  being the corresponding degree matrix. Here,  $tr(\cdot)$  denotes the trace operator, summing the diagonal elements of a matrix.

Intuitively, when  $A^v_{ij}$  is large—implying high similarity between samples i and j—the regularization penalizes large distances  $||z^v_i - z^v_j||_2^2$  in the latent space. This encourages similar samples to stay close, preserving local structure and guiding the model to learn embeddings that reflect both feature content and intrinsic neighborhood relationships, thus maintaining the data's manifold structure.

#### 2.2.3 Self-paced Adaptive Graph Construction

Inspired by self-paced learning—progressing from easy to hard—we construct the similarity graph progressively: early training starts from a sparse k-NN backbone to stabilize optimization, and we gradually enlarge neighborhoods to enrich structure and learn more discriminative representations. Concretely, for view v we keep, for each sample i, only its k nearest neighbors; non-neighbors have zero similarity, while neighbors use edge weights  $D_{i,k+1}^{(v)} - D_{ij}^{(v)}$  (with  $D_{ij}^{(v)}$  the distance to the j-th neighbor), followed by row-wise normalization over the k neighbors.

$$W_{ij}^{v} = \frac{D_{i,k+1}^{v} - D_{i,j}^{v}}{\sum_{m=1}^{k} (D_{i,k+1}^{v} - D_{i,j}^{v})}, 1 \le j \le k,$$
(2)

where  $D_{i,m}^v$  denotes the distance between sample i and its m-th nearest neighbors.

Moreover, we calculate the distances  $D^v$  between samples based on the original data  $X^v$  only during the initialization stage. Subsequently,  $D^v$  are computed based on the learned representations  $Z^v$ :  $D^v_{ij} = ||Z^v_i - Z^v_j||_2^2$ .

By progressively increasing k and updating the similarity matrix based on the learned representation, our method can increasingly explore more reliable graph information and facilitate within-view representation learning.

# 2.3 Inter-view Self-weight Contrastive Learning

# 2.3.1 Representational Consistency Constraint

Given that various perspectives of an object inherently possess consistent characteristics, we enforce this consistency across views through a mechanism referred to as the representational consistency constraint. This constraint promotes alignment among the representations derived from different views, thereby minimizing redundancy and enhancing overall consistency.

$$L_{rcc} = \sum_{v_i, v_j} ||Z^{(v_i)} - Z^{(v_j)}||_F^2.$$
(3)

#### 2.3.2 Global Feature Generation

To integrate complementary information across views and learn compact, discriminative representations for clustering, we design a global fusion module. Given latent embeddings  $Z^1, Z^2, \ldots, Z^V \in \mathbb{R}^{n \times d_v}$  from V views, we concatenate them along the feature dimension to form the initial global representation  $Z_{\text{concat}} = [Z^1, Z^2, \ldots, Z^V] \in \mathbb{R}^{n \times (\sum_{v=1}^V d_v)}$ . For multi-view relational structure, we first compute symmetrically normalized Laplacians  $L^v$  for each view and average them to obtain the consolidated similarity matrix  $\hat{L} = \frac{1}{V} \sum_{v=1}^V L^v$ . Using this graph prior, we refine the concatenated features by propagating with the consolidated Laplacian to get  $R = \hat{L} Z_{\text{concat}}$ , which enhances intersample affinities and discriminability. Finally, we map R into a shared latent space with an MLP, yielding the global representation  $H = \hat{W}_2 \sigma(\hat{W}_1 R + b_1) + b_2 \in \mathbb{R}^{n \times d_h}$ .

#### 2.3.3 Self-weighted Contrastive Learning

Multi-view contrastive learning has demonstrated strong potential in aligning complementary information from different views. However, conventional methods typically treat all views equally, using uniform weights when computing contrastive losses. Formally, they adopt a view-invariant formulation such as:

$$L_{CL} = \sum_{m,n} L_{CL}^{m,n}(Z^m, Z^n), \tag{4}$$

where  $Z^m$  and  $R^n$  denote representations of views m and n, respectively. While this symmetric formulation facilitates consistency across views, it can undesirably amplify the influence of low-quality or noisy views by forcing them to align equally with high-quality ones. This uniform treatment may lead to representational degeneration and hinder effective feature fusion.

To address this limitation, we propose an inter-view self-weighted contrastive learning strategy that adaptively modulates the contribution of each view based on its semantic alignment with a shared global representation. The core idea is to prioritize reliable, informative views in the contrastive process while suppressing the impact of unreliable ones. Specifically, we reformulate the contrastive loss as:

$$L_{sw} = \sum_{m,n} w^v \cdot L_{sw}^{m,n}(Z^m, H), \tag{5}$$

where m is the number of views,  $Z^v$  denotes the view-specific representation, H is the fused global representation, and  $W^v$  is the adaptive weight reflecting the relative reliability of the v-th view.

Since labels are unavailable in unsupervised settings, directly evaluating the quality of a view is challenging. To estimate the semantic relevance of each view, we assess the distributional discrepancy between  $Z^v$  and H. A lower discrepancy implies a higher alignment with global semantics and thus a more trustworthy view. This discrepancy is denoted as:

$$\mathcal{D}^v = \mathcal{D}(Z^v, H),\tag{6}$$

where  $\mathcal{D}(\cdot,\cdot)$  is is a distance metric based on Maximum Mean Discrepancy (MMD) Wu et al. (2024), a non-parametric criterion that measures the distance between two distributions in a Reproducing Kernel Hilbert Space (RKHS). Given two feature sets  $X_s = \{X_i^s\}_{i=1}^{n_s}$  and  $Y_t = \{y_j^t\}_{j=1}^{n_t}$ , the squared MMD is defined as:

$$MMD^{2}(X_{s}, Y_{t}) = \frac{1}{n_{s}^{2}} \sum_{i,j=1}^{n_{s}} k(X_{i}^{s}, x_{j}^{s}) + \frac{1}{n_{t}^{2}} \sum_{i,j=1}^{n_{t}} k(y_{i}^{t}, y_{j}^{t}) - \frac{2}{n_{s}n_{t}} \sum_{i=1}^{n_{s}} \sum_{j=1}^{n_{t}} k(x_{i}^{s}, y_{j}^{t}),$$

$$(7)$$

where  $k(\cdot, \cdot)$  is a kernel function. In our case, we employ a linear kernel  $k(x, y) = x^T y$ , which avoids the need for hyperparameter tuning and suits high-dimensional representations. Given that

 $Z^{v}$  and H share the same dimensions, the discrepancy for each view is computed as:

 $MMD^{2}(Z^{v}, H) = \frac{1}{N^{2}} \sum_{i,j=1}^{N} k(Z_{i}^{v}, Z_{j}^{v}) + \frac{1}{N^{2}} \sum_{i,j=1}^{N} k(H_{i}, H_{j}) - \frac{2}{N^{2}} \sum_{i,j=1}^{N} k(Z_{i}^{v}, H_{j}),$ (8)

where N denotes the total number of samples. Based on these discrepancies, we define a normalized weight allocation function to adaptively determine the importance of each view:

$$w^{v} = \mathcal{P}(\mathcal{D}^{v}) = softmax(-\mathcal{D}^{v}). \tag{9}$$

The use of the negative discrepancy ensures that views more consistent with global semantics receive higher weights. This adaptive weighting mechanism promotes semantically aligned views and effectively suppresses noisy or misleading ones, thereby enhancing the robustness and expressiveness of the learned global representations.

#### 2.4 Self-supervised Clustering Module

In unsupervised learning, we refine the unified representation H by integrating multi-view information that captures shared and complementary patterns. Since H may not be immediately clustering-friendly, we further enhance it with a self-supervised clustering objective.

#### 2.4.1 Clustering Loss via KL Divergence

We adopt a Kullback-Leibler divergence between a target distribution P and a soft assignment Q:

$$L_{cl} = D_{KL}(P||Q) = \sum_{i} \sum_{j} p_{ij} \log \frac{p_{ij}}{q_{ij}}.$$
 (10)

Here, Q is the soft label distribution and P is the sharpened target; the KL term measures information loss when approximating P by Q.

#### 2.4.2 SOFT LABEL DISTRIBUTION Q

We compute  $q_{ij}$  via a Student-t kernel between feature  $h_i$  and centroid  $\mu_i$ :

$$q_{ij} = \frac{\left(1 + \|h_i - \mu_j\|^2 / \sigma^2\right)^{-(\alpha+1)/2}}{\sum_f \left(1 + \|h_i - \mu_f\|^2 / \sigma^2\right)^{-(\alpha+1)/2}},$$
(11)

where  $\sigma$  controls the kernel scale.

## 2.4.3 TARGET DISTRIBUTION P

To emphasize confident assignments and balance clusters, we set

$$p_{ij} = \frac{q_{ij}^2/f_j}{\sum_f q_{if}^2/f_f}, \qquad f_j = \sum_i q_{ij},$$
 (12)

so that larger  $q_{ij}$  contributes more while normalizing by cluster frequency.

The final label for node  $v_i$  is

$$s_i = \arg\max_j \, q_{ij}. \tag{13}$$

This self-supervised head aligns H with clustering by sharpening confident assignments, mitigating unreliable signals, and improving separability without external labels.

Table 1: Datasets Descriptions

| Dataset     | Clusters | Samples | Dimensionality            |  |  |
|-------------|----------|---------|---------------------------|--|--|
| COIL20      | 20       | 1140    | [1024, 3304, 6750]        |  |  |
| Handwritten | 10       | 2000    | [240, 76, 216, 47, 64, 6] |  |  |
| HW1256      | 10       | 2000    | [76, 216, 47, 6]          |  |  |
| Caltech     | 7        | 1400    | [40, 254, 1984, 512, 928] |  |  |
| MNIST-USPS  | 10       | 5000    | [784, 256]                |  |  |
| Fashion     | 10       | 10000   | [784, 784, 784]           |  |  |

Table 2: Clustering Results on COIL20, Handwritten, HW1256 and MNIST-USPS Datasets

| Dataset               | COIL20 |        | Handwritten |        | HW1256 |        | Caltech |        | MNIST-USPS |        | Fashion |        |
|-----------------------|--------|--------|-------------|--------|--------|--------|---------|--------|------------|--------|---------|--------|
|                       | ACC    | NMI    | ACC         | NMI    | ACC    | NMI    | ACC     | NMI    | ACC        | NMI    | ACC     | NMI    |
| $\overline{K}$ -means | 0.4142 | 0.3895 | 0.5128      | 0.4827 | 0.6564 | 0.6799 | 0.2345  | 0.0274 | 0.5191     | 0.3609 | 0.4465  | 0.1934 |
| <b>DUA-Nets</b>       | 0.7228 | 0.8272 | 0.6585      | 0.5924 | 0.7425 | 0.7933 | 0.5461  | 0.0154 | 0.9136     | 0.8359 | 0.7747  | 0.8145 |
| SGFCC                 | 0.2590 | 0.4381 | 0.3870      | 0.5501 | 0.3840 | 0.5118 | 0.4817  | 0.5262 | 0.9526     | 0.9412 | 0.9286  | 0.9180 |
| CoMSC                 | 0.5482 | 0.7382 | 0.5881      | 0.4914 | 0.7320 | 0.6793 | 0.4105  | 0.4830 | 0.7252     | 0.7025 | 0.6050  | 0.7158 |
| CMRL                  | 0.6264 | 0.7575 | 0.5439      | 0.4865 | 0.8947 | 0.8168 | 0.4082  | 0.3399 | 0.9308     | 0.8690 | 0.5483  | 0.6134 |
| ASR-ETR               | 0.6611 | 0.7940 | 0.7580      | 0.6930 | 0.7290 | 0.6487 | 0.5096  | 0.5133 | 0.7580     | 0.6930 | 0.7186  | 0.7351 |
| RCAGL                 | 0.6701 | 0.8127 | 0.8775      | 0.8061 | 0.9305 | 0.8623 | 0.6341  | 0.4871 | 0.8925     | 0.7316 | 0.7924  | 0.8097 |
| HFMVC                 | 0.4558 | 0.5956 | 0.9080      | 0.8341 | 0.8785 | 0.7927 | 0.5863  | 0.3280 | 0.9010     | 0.8431 | 0.9110  | 0.9008 |
| GCFAgg                | 0.3458 | 0.4886 | 0.8085      | 0.7752 | 0.8005 | 0.7664 | 0.3813  | 0.4321 | 0.9300     | 0.8896 | 0.8982  | 0.8714 |
| SCMVC                 | 0.5153 | 0.6451 | 0.8945      | 0.8168 | 0.7945 | 0.7047 | 0.4905  | 0.4390 | 0.9576     | 0.9525 | 0.9229  | 0.9213 |
| DCMVC                 | 0.7340 | 0.8162 | 0.8995      | 0.8718 | 0.7580 | 0.7620 | 0.3161  | 0.2460 | 0.8920     | 0.9059 | 0.7836  | 0.8745 |
| Ours                  | 0.9153 | 0.9651 | 0.9115      | 0.8467 | 0.9560 | 0.9145 | 0.6679  | 0.5345 | 0.9628     | 0.9515 | 0.9328  | 0.8935 |

## 2.5 Training

The training procedure is divided into two main phases: preliminary training and subsequent finetuning. During the preliminary training phase, the number of clusters k is incrementally increased to its maximum value, and the model is trained by optimizing the objective function as described in Equation (25). In the fine-tuning phase, k is held constant, and we enforce a consistency constraint on the representations. The model is then refined by minimizing the loss function presented in Equation (26). Ultimately, we apply Self-supervised Clustering Module to the consolidated representation H to derive the clustering outcomes. The entire workflow is depicted in Algorithm 1.

The preliminary training loss is given by:

$$L_{pre} = L_{rc} + \lambda_1 L_{gls}. \tag{14}$$

The fine-tuning loss is defined as:

$$L_{fine} = L_{rc} + \lambda_1 L_{als} + \lambda_2 L_{rcc} + \lambda_3 L_{sw} + \lambda_4 L_{cl}. \tag{15}$$

Here,  $\lambda_1, \lambda_2, \lambda_3$  and  $\lambda_4$  are coefficients that regulate the impact of the graph-based and consistency terms within the total loss function, respectively.

#### **EXPERIMENTS**

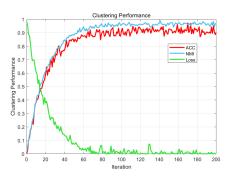
#### 3.1 Datasets

COIL20 comprises grayscale images of 20 objects across 360° poses. Handwritten and HW1256 are multi-view handwritten digits (differing in the number of views). Caltech contains multi-feature object/scene images. MNIST-USPS mixes two digit sources to form a cross-domain benchmark. Fashion consists of clothing images with multiple attributes/views. Cluster counts, sample sizes, and view dimensionalities are in Table 1.

#### 3.2 Comparative Algorithms

Baselines fall into three groups: (i) adaptive weighting/uncertainty (DUA-Nets Geng et al. (2021), RCAGL Liu et al. (2024), SCMVC Wu et al. (2024)), which modulate view contributions by reliability; (ii) subspace/anchor representations (CoMSC Liu et al. (2021), CMRL Zheng et al. (2023), AER-ETR Ji & Feng (2023)) to reduce redundancy via compact bases; and (iii) contrastive/structural constraints (HFMVC Jiang et al. (2024), DCMVC Cui et al. (2024), GCFAgg Yan et al. (2023), SGFCC Shu et al. (2024)) to enforce cross-view consistency and cluster structure. Most do not jointly leverage graph-structural guidance with contrastive consistency; SPEAG unifies both.

#### 3.3 MODEL ANALYSIS



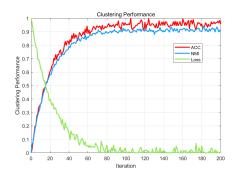
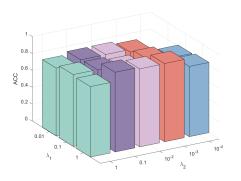


Figure 2: Clustering performance with increasing iteration on COIL20 and HW1256



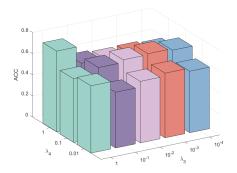


Figure 3: ACC sensitivity on Caltech: left— $\lambda_1, \lambda_2$ ; right— $\lambda_3, \lambda_4$ .

## 3.3.1 Performance Evaluation

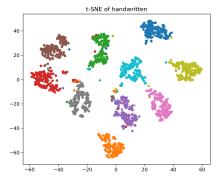
We evaluate on six benchmarks using ACC/NMI (Table 3). Findings: (1) SPEAG achieves best or second-best results on most datasets, driven by self-paced graph construction and structure-aware contrastive learning; (2) versus shallow/hybrid methods (KMeans, CoMSC, ASR-ETR, RCAGL), SPEAG better captures nonlinear cross-view relations—particularly strong on image datasets (MNIST-USPS, Fashion); (3) compared with deep baselines (DUA-Nets, CMRL, HFMVC, SCMVC, DCMVC, GCFAgg, SGFCC), SPEAG augments contrastive alignment with explicit graph supervision, yielding more clustering-friendly embeddings than methods that emphasize only consistency or only contrast.

#### 3.3.2 ABLATION STUDY

We study four losses on COIL20: graph regularization  $\mathcal{L}_{gls}$ , cross-view consistency  $\mathcal{L}_{rcc}$ , self-weighted contrastive  $\mathcal{L}_{sw}$ , and self-supervised clustering  $\mathcal{L}_{cl}$ . Results show  $\mathcal{L}_{gls}$  notably improves clustering; removing any fine-tuning loss degrades performance—most severely without  $\mathcal{L}_{cl}$  (weaker instance discrimination). Dropping  $\mathcal{L}_{sw}$  harms cross-view distribution alignment, and dropping  $\mathcal{L}_{rcc}$  weakens structural consistency. The full SPEAG model is best.

Table 3: Ablation Study on COIL20 dataset in terms of ACC (%), NMI (%) and ARI(%).

|                     |                     |                    |                    |       | COIL20 |       |  |  |  |
|---------------------|---------------------|--------------------|--------------------|-------|--------|-------|--|--|--|
| $\mathcal{L}_{gls}$ | $\mathcal{L}_{rcc}$ | $\mathcal{L}_{sw}$ | $\mathcal{L}_{cl}$ | ACC   | NMI    | ARI   |  |  |  |
|                     |                     |                    |                    | 78.61 | 85.53  | 72.45 |  |  |  |
|                     | $\checkmark$        |                    |                    | 81.25 | 87.87  | 77.04 |  |  |  |
|                     |                     | $\checkmark$       |                    | 78.54 | 85.90  | 73.47 |  |  |  |
|                     |                     |                    | $\checkmark$       | 80.07 | 86.26  | 75.22 |  |  |  |
| $\checkmark$        | $\checkmark$        |                    |                    | 81.01 | 87.63  | 77.38 |  |  |  |
| $\checkmark$        |                     | $\checkmark$       |                    | 87.57 | 94.55  | 86.23 |  |  |  |
| $\checkmark$        |                     |                    | $\checkmark$       | 81.04 | 87.32  | 75.86 |  |  |  |
| $\checkmark$        | $\checkmark$        | $\checkmark$       | $\checkmark$       | 91.53 | 96.51  | 90.87 |  |  |  |



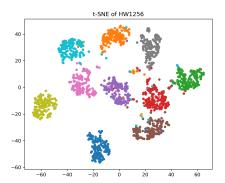


Figure 4: T-SNE visualization on the datasets handwritten and HW125

#### 3.3.3 PARAMETERS AND CONVERGENCE ANALYSIS

As iterations increase (Fig. 3), ACC/NMI rise and the loss decreases, indicating stable convergence and continuous improvement. Fig. 2 shows hyperparameter sensitivity:  $\lambda_1$  and  $\lambda_3$  have stronger effects; within reasonable ranges, larger values generally yield more robust gains.

## 4 Conclusion

In this work, we have presented SPEAG, a novel self-paced exemplar-aware graph learning framework for multi-view clustering. By integrating an exemplar-guided attention mechanism with a self-paced training strategy, SPEAG effectively balances the exploration of consistent and complementary information across views while progressively mitigating the impact of noisy or low-quality samples. Moreover, the joint learning of view-specific and consensus representations within a unified anchor graph structure allows for more robust clustering performance. Extensive experiments on multiple benchmark datasets demonstrate that our method achieves competitive or superior results compared to state-of-the-art approaches. In future work, we plan to extend SPEAG to handle streaming or dynamically evolving multi-view data, and explore its potential in semi-supervised and federated clustering scenarios.

## REFERENCES

Deng Cai, Xiaofei He, Xiaoyun Wu, and Jiawei Han. Non-negative matrix factorization on manifold. In 2008 eighth IEEE international conference on data mining, pp. 63–72. IEEE, 2008.

Zhe Chen, Xiao-Jun Wu, Tianyang Xu, Hui Li, and Josef Kittler. Deep discriminative multi-view clustering. *IEEE Transactions on Circuits and Systems for Video Technology*, 2025a.

Zhe Chen, Xiao-Jun Wu, Tianyang Xu, Hui Li, and Josef Kittler. Multi-layer multi-level comprehensive learning for deep multi-view clustering. *Information Fusion*, 116:102785, 2025b.

- Jinrong Cui, Yuting Li, Han Huang, and Jie Wen. Dual contrast-driven deep multi-view clustering. *IEEE Transactions on Image Processing*, 2024.
- Uno Fang, Man Li, Jianxin Li, Longxiang Gao, Tao Jia, and Yanchun Zhang. A comprehensive survey on multi-view clustering. *IEEE Transactions on Knowledge and Data Engineering*, 35 (12):12350–12368, 2023.
  - Yu Geng, Zongbo Han, Changqing Zhang, and Qinghua Hu. Uncertainty-aware multi-view representation learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 7545–7553, 2021.
- Jintian Ji and Songhe Feng. Anchor structure regularization induced multi-view subspace clustering via enhanced tensor rank minimization. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 19343–19352, 2023.
  - Bingbing Jiang, Chenglong Zhang, Xinyan Liang, Peng Zhou, Jie Yang, Xingyu Wu, Junyi Guan, Weiping Ding, and Weiguo Sheng. Collaborative similarity fusion and consistency recovery for incomplete multi-view clustering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 17617–17625, 2025.
  - Xiaorui Jiang, Zhongyi Ma, Yulin Fu, Yong Liao, and Pengyuan Zhou. Heterogeneity-aware federated deep multi-view clustering towards diverse feature representations. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pp. 9184–9193, 2024.
  - Zhao Kang, Guoxin Shi, Shudong Huang, Wenyu Chen, Xiaorong Pu, Joey Tianyi Zhou, and Zenglin Xu. Multi-graph fusion for multi-view spectral clustering. *Knowledge-Based Systems*, 189:105102, 2020.
    - Chenghua Liu, Zhuolin Liao, Yixuan Ma, and Kun Zhan. Stationary diffusion state neural estimation for multiview clustering. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, pp. 7542–7549, 2022.
    - Jiyuan Liu, Xinwang Liu, Yuexiang Yang, Xifeng Guo, Marius Kloft, and Liangzhong He. Multiview subspace clustering via co-training robust data representation. *IEEE Transactions on Neural Networks and Learning Systems*, 33(10):5177–5189, 2021.
    - Suyuan Liu, Qing Liao, Siwei Wang, Xinwang Liu, and En Zhu. Robust and consistent anchor graph learning for multi-view clustering. *IEEE Transactions on Knowledge and Data Engineering*, 36 (8):4207–4219, 2024.
    - Zhenqiu Shu, Bin Li, Cunli Mao, Shengxiang Gao, and Zhengtao Yu. Structure-guided feature and cluster contrastive learning for multi-view clustering. *Neurocomputing*, 582:127555, 2024.
    - Jie Wen, Na Han, Xiaozhao Fang, Lunke Fei, Ke Yan, and Shanhua Zhan. Low-rank preserving projection via graph regularized reconstruction. *IEEE transactions on cybernetics*, 49(4):1279–1291, 2018.
    - Song Wu, Yan Zheng, Yazhou Ren, Jing He, Xiaorong Pu, Shudong Huang, Zhifeng Hao, and Lifang He. Self-weighted contrastive fusion for deep multi-view clustering. *IEEE Transactions on Multimedia*, 26:9150–9162, 2024.
  - Weiqing Yan, Yuanyang Zhang, Chenlei Lv, Chang Tang, Guanghui Yue, Liang Liao, and Weisi Lin. Gcfagg: Global and cross-view feature aggregation for multi-view clustering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 19863–19872, 2023.
  - Yang Yang and Changming Zhu. Deep multi-view clustering based on global hybrid alignment with cross-contrastive learning. *The Visual Computer*, pp. 1–13, 2024.
- Qinghai Zheng, Jihua Zhu, Zhongyu Li, Zhiqiang Tian, and Chen Li. Comprehensive multi-view representation learning. *Information Fusion*, 89:198–209, 2023.
  - Yiyang Zhou, Qinghai Zheng, Shunshun Bai, and Jihua Zhu. Semantically consistent multi-view representation learning. *Knowledge-Based Systems*, 278:110899, 2023.