# Center-Oriented Prototype Contrastive Clustering

1ˢᵗ Shihao Dong
*School of Computer Science*
*Nanjing University of Information Science and Technology*
Nanjing, China
dongshihao@nuist.edu.cn

2ⁿᵈ Xiaotong Zhou
*School of Computer Science*
*Nanjing University of Information Science and Technology*
Nanjing, China
xiaotong_zhou@nuist.edu.cn

3ʳᵈ Yuhui Zheng
*The State Key Laboratory of Tibetan Intelligence*
*Qinghai Normal University*
Xining, China
zhengyh@vip.126.com

4ᵗʰ Huiying Xu⋆
*Computer Science and Technology*
*Zhejiang Normal University*
Jinhua, China
xhy@zjnu.edu.cn

5ᵗʰ Xinzhong Zhu
*Computer Science and Technology*
*Zhejiang Normal University*
Jinhua, China
zxz@zjnu.edu.cn

*Abstract*—Contrastive learning is widely used in clustering tasks due to its discriminative representation. However, the conflict problem between classes is difficult to solve effectively. Existing methods try to solve this problem through prototype contrast, but there is a deviation between the calculation of hard prototypes and the true cluster center. To address this problem, we propose a center-oriented prototype contrastive clustering framework, which consists of a soft prototype contrastive module and a dual consistency learning module. In short, the soft prototype contrastive module uses the probability that the sample belongs to the cluster center as a weight to calculate the prototype of each category, while avoiding inter-class conflicts and reducing prototype drift. The dual consistency learning module aligns different transformations of the same sample and the neighborhoods of different samples respectively, ensuring that the features have transformation-invariant semantic information and compact intra-cluster distribution, while providing reliable guarantees for the calculation of prototypes. Extensive experiments on five datasets show that the proposed method is effective compared to the SOTA. Our code is published on https://github.com/LouisDong95/CPCC.

*Index Terms*—self-supervised learning, contrastive learning, representation learning, deep clustering.

## I. INTRODUCTION

In the data-driven world, deep clustering, as a basic tool for exploring the intrinsic structure and patterns of high-dimensional data, has been a research hotspot in the fields of deep learning and data mining. In recent years, contrastive learning [1], [2] has been widely used in clustering tasks due to its ability to efficiently learn feature representations of data without expensive labels by using the information of the data itself as training signals.

Although some works [3]–[6] have achieved good performance by combining deep clustering with contrastive learning, contrastive learning treats the same samples and their transformations as positive pairs, while the remaining samples are
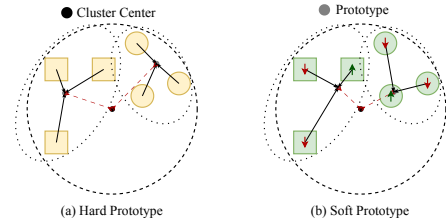
Fig. 1. **Our motivation**. The circular squares represent samples with different transformations. The prototypes are calculated by random sampling. The hard prototype calculation treats each sample equally, causing the calculated prototype to deviate from the true cluster center. The soft prototype assigns higher weights to samples close to the cluster center and reduces the weights of distant points, so that the calculated prototype is closer to the true cluster center

treated as negative pairs. Negative pairs contain samples of the same class with the same semantic information, which inevitably leads to inter-class conflicts, thereby destroying the discriminative information. Therefore, how to alleviate false pairs has always been a hot topic of concern in contrastive learning. Most existing methods alleviate the impact of false pairs by reducing the weight of false negative pairs and mining false positive pairs. GCC [7] constructs positive pairs through the adjacency graph, and NNM [6] constructs positive pairs by combining local and global neighbors. GDCL [8] deleted false negative samples that are similar to the original samples through a bias elimination strategy. TCL [9] removed the false negative samples from the contrastive loss denominator. HSAN [10] reduced the weight between easy sample pairs and increased the weight between difficult sample pairs. DIVIDE [11] reduces the weight of false positive samples through random walks. Although conflicts between classes are effectively mitigated, it does not fundamentally solve this problem. To avoid inter-class conflicts, PCL [12] calculates $K$ prototypes and considers the sample and its corresponding prototype as a positive pair, while the remaining prototypes are considered negative samples. ProPos [13] calculates $2K$

prototypes in pairs through data augmentation, and considers prototypes of the same class as positive pairs, which not only avoids class conflicts but also greatly improves the efficiency of contrastive learning. Although the prototype contrast method effectively avoids inter-class conflicts, it treats samples equally when calculating the prototype, which inevitably introduces noise, causing the prototype estimation to deviate from the true cluster center and inaccurate prototype calculation, as shown in the Fig. 1.

To address the above problem, we propose a **C**enter-oriented **P**rototype **C**ontrastive **C**lustering method, termed CPCC. Our method consists of soft prototype contrastive (SPC) and dual consistency learning (DCL). Specifically, SPC considers prototypes from the same category as positive pairs and prototypes from different categories as negative pairs. The prototype is calculated by taking the square of the probability that each sample belongs to its cluster center as the weight, so that high-confidence samples are assigned higher weights, otherwise, they are assigned lower weights. This strategy can not only avoid category conflicts but also reduce the deviation of prototypes from cluster centers. In addition, DCL learns the consistency of the features of the same sample under different transformations and the consistency of the features between different samples and their neighbors, resulting in further compression of the intra-class feature space. The contributions are summarized as follows:

- We propose a novel prototype contrastive clustering method to avoid inter-cluster conflicts and improve intra-cluster compactness through prototype contrast and consistency learning.
- Compared with hard prototypes, we propose a soft prototype construction method that can avoid prototype drift.
- Extensive experimental results on five benchmark datasets demonstrate the effectiveness of the proposed method compared with the existing SOTA.

## II. RELATED WORK

### A. Deep Clustering

Deep clustering utilizes deep learning techniques to perform clustering tasks. It guides the learning process through the features and structure of the data itself without relying on human intervention. According to the network model, deep clustering can be divided into generative network-based and discriminative network-based deep clustering. 1) Generative network-based, such as DEC [14] realized clustering on embedded features by reconstructing the data through AE. VaDE [15] sampled from Gaussian Mixture Models distributions which not only realized clustering but also generated new samples. ClusterGAN [16] realized clustering and generated samples through the generative adversarial network. 2) Discriminative networks-based, such as SCAN [3] pre-trained a pair of siamese networks with shared weights by contrastive learning and combined them with neighbor assignment consistency to achieve clustering. NNM [6] extended neighbor assignment consistency to local neighbors based on SCAN. CC [4] im-
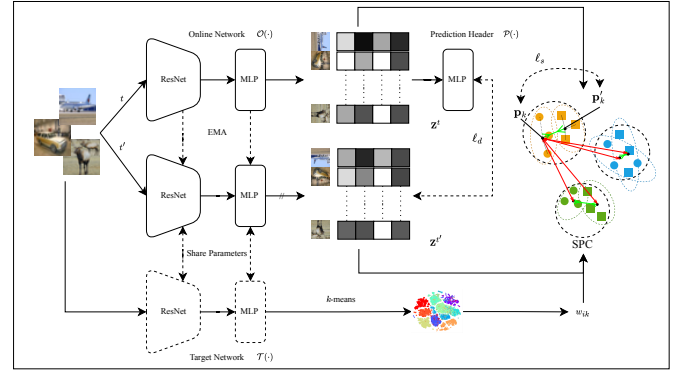


Fig. 2. **Illustration of CPCC framework.** Our framework consists of an online network $\mathcal{O}(\cdot)$ and a target network $\mathcal{T}(\cdot)$. The target network synchronizes parameters with the online network through a moving average strategy. First, the cluster center of the original data set is obtained through the $\mathcal{T}(\cdot)$, and the weight $\mathbf{w}_i$ of each sample is calculated based on the distance between the feature and its center. Secondly, the features $\mathbf{Z}^t$ and $\mathbf{Z}^{t'}$ of different transformed samples are obtained through the $\mathcal{O}(\cdot)$ and the $\mathcal{T}(\cdot)$ respectively, $2K$ prototypes of different transformations are calculated based on the features and its weights, and the network is trained through SPC loss and DCL loss

proved the discrimination of the cluster space through bi-level contrasts. SPICE [5] combined contrastive learning and pseudo-labeling for semi-supervised training of networks.

### B. Contrastive Learning

Recently, contrastive learning has received much attention for its ability to learn discriminative feature representations from unlabeled data. Such methods are usually combined with a specific task to implement clustering based on MoCo [1] or SimCLR [2], whose loss functions are usually implemented by InfoNEC or NT-Xnet. However, these two methods usually require a large number of negative samples, which leads to class conflicts. One way to avoid conflicts is to not require negative sample pairs: BYOL [17] is unique in that it does not require negative pairs for contrastive learning, but rather learns the feature representations by maximizing the consistency of the outputs of the online network and the target network, thus avoiding the selection of negative samples. Another way is through prototype contrastive: PCL [12] is a contrastive framework between samples and their prototypes, and Pro-Pos [13] is a contrastive sample between prototypes. It takes prototypes of the same type as positive pairs and prototypes of different types as negative pairs, which not only avoids inter-class conflicts but also reduces computational complexity.

## III. METHOD

### A. Soft Prototype Contrastive

Given a dataset $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^N \in \mathbb{R}^{N \times D}$ and a set of transformations $T$, contrastive learning constructs pairs of samples by data augmentation $\{\mathbf{X}^t, \mathbf{X}^{t'}\}, t, t' \in T$. Where different transformations of the same sample are considered positive pairs and the remaining samples are considered negative pairs. The discrimination of each sample representation

can be effectively increased by increasing the similarity of positive pairs and decreasing the similarity of negative pairs. The formula is as follows:

$$\ell_i = -\frac{1}{N}\sum_{i=1}^{N}\log\frac{e^{(\mathbf{z}_i^t)^\top \mathbf{z}_i^{t'}/\tau}}{\sum_{j=1,j\neq i}^{N}[e^{(\mathbf{z}_i^t)^\top \mathbf{z}_j^t/\tau} + e^{(\mathbf{z}_i^t)^\top \mathbf{z}_j^{t'}/\tau}]}. \quad (1)$$

Although contrastive learning methods achieve high-quality feature representation capability by this property, the false negative pairs in the denominator of Eq. (1) will lead to inter-class conflicts and affect the final result of clustering. To this end, we propose an SPC module to avoid this conflict. First, we obtain the feature $\mathbf{Z}$ of the sample through the target network $\mathcal{T}(\cdot)$, and perform $k-$means on $\mathbf{Z}$ to initialize the cluster center $\{\boldsymbol{\mu}_i\}_{i=1}^{K}$ of each class. Then, the soft assignment $\mathbf{Q}$ between each sample $\mathbf{Z}$ and the cluster centers $\boldsymbol{\mu}$ can be calculated from the student's $t$-distribution:

$$q_{ik} = \frac{(1+\|\mathbf{z}_i - \boldsymbol{\mu}_k\|^2/\alpha)^{-\frac{\alpha+1}{2}}}{\sum_{k'}(1+\|\mathbf{z}_i - \boldsymbol{\mu}_{k'}\|^2/\alpha)^{-\frac{\alpha+1}{2}}}, \quad (2)$$

where $q_{ik}$ is denoted as the probability that the $i$-th sample belongs to the $k$-th cluster, $\alpha = 1$ for all experiments. We computed the weights by raising $\mathbf{Q}$ to quadratic such that features close to the cluster center are assigned larger weights and features far from the cluster center are assigned smaller weights:

$$w_{ik} = \frac{q_{ik}^2/f_k}{\sum_{k'} q_{ik'}^2/f_{k'}}, \quad (3)$$

where $f_k = \sum_i q_{ik}$ are soft cluster frequencies. We estimate the prototypes of each class through mini-batch samples $\mathcal{B}$. Due to the presence of $\mathbf{W}$, the drift between the prototype and the cluster center caused by random sampling is avoided. The prototypes $\mathbf{P}$ and $\mathbf{P}'$ can be obtained by the following formula:

$$\mathbf{p}_k = \frac{\sum_{i=1}^{|\mathcal{B}|} w_{ik}\mathcal{O}(\mathbf{x}_i^t)}{\|\sum_{i=1}^{|\mathcal{B}|} w_{ik}\mathcal{O}(\mathbf{x}_i^t)\|_2}, \mathbf{p}_k' = \frac{\sum_{i=1}^{|\mathcal{B}|} w_{ik}\mathcal{T}(\mathbf{x}_i^{t'})}{\|\sum_{i=1}^{|\mathcal{B}|} w_{ik}\mathcal{T}(\mathbf{x}_i^{t'})\|_2}. \quad (4)$$

After we get the $2K$ prototypes $\{\mathbf{p}_1, \mathbf{p}_2, ..., \mathbf{p}_K\}$ and $\{\mathbf{p}_1', \mathbf{p}_2', ..., \mathbf{p}_K'\}$, the soft prototype contrastive loss is calculated as follows:

$$\ell_s = -\frac{1}{K}\sum_{k=1}^{K}\log\frac{e^{(\mathbf{p}_k)^\top \mathbf{p}_k'/\tau}}{\sum_{j=1,j\neq k}^{K}[e^{(\mathbf{p}_k)^\top \mathbf{p}_j/\tau} + e^{(\mathbf{p}_k)^\top \mathbf{p}_j'/\tau}]}, \quad (5)$$

where the temperature parameter $\tau$ is used to adjust the scale of similarity. With $\ell_s$, similar prototypes of the same class are close to each other in space, while prototypes of different classes are far away from each other. This avoids inter-class conflicts and prototype drift, and improves the distinction between categories and stability during training.

### B. Dual Consistency Learning

The same sample after different transformations should be consistent in the feature space, and the consistency of the features allows the network to ignore some details and learn the semantic information that is invariant to the transformations.

In addition, the neighbors of different samples should also be close in the feature space. Therefore, we not only consider different transformations of the same sample as positive pairs, but also consider the corresponding neighbors as positive pairs. The dual consistency learning loss is defined as follows:

$$\ell_d = \frac{1}{2N}\sum_i^{N}(\|\mathcal{P}(\mathbf{z}_i^t) - \mathbf{z}_i^{t'}\|_2^2 + \|\mathcal{P}(\mathbf{z}_i^{t'} + \sigma\epsilon) - \mathbf{z}_i^t\|_2^2)$$

$$= 2 - \frac{1}{N}\sum_i^{N}(<\mathcal{P}(\mathbf{z}_i^t), \mathbf{z}_i^{t'}> + <\mathcal{P}(\mathbf{z}_i^{t'} + \sigma\epsilon), \mathbf{z}_i^t>), \quad (6)$$

where $\epsilon \sim \mathcal{N}(0, \mathbf{I})$, $\sigma$ is used to control the range of neighbors. $\mathbf{z}_i + \sigma\epsilon$ denotes the neighborhood sampling for $\mathbf{z}_i$. The first part of Eq. (6) represents the consistency between the sample prediction and its transformation, and the second part represents the consistency between the sample's neighborhood prediction and its transformation. The two constrain transformation invariance and intra-class compactness respectively. Consistency learning not only aligns features of different transformations, but also provides high-confidence pseudo-labels for prototypes.

### C. Model Training

In summary, the total loss of CPCC is defined as:

$$\ell = \ell_d + \lambda\ell_s, \quad (7)$$

where $\lambda$ is a trade-off between $\ell_s$ and $\ell_d$. The training is divided into two stages. Since the clustering results are unreliable in the early stage of training, only $\ell_d$ is used in the pre-training stage, and the contrastive stage is trained by Eq. (7). The parameters $\boldsymbol{\theta}_{\mathcal{O}}$ of $\mathcal{O}(\cdot)$ and the parameters $\boldsymbol{\theta}_{\mathcal{T}}$ of $\mathcal{T}(\cdot)$ is updated by EMA as follows:

$$\boldsymbol{\theta}_{\mathcal{T}} = m\boldsymbol{\theta}_{\mathcal{T}} + (1-m)\boldsymbol{\theta}_{\mathcal{O}}, \quad (8)$$

Where $m \in [0, 1)$ is the coefficient that controls the rate of updating. The procedure for the optimization algorithm is summarized in Algorithm 1.

## IV. EXPERIMENTS

### A. Datasets and Evaluation Metrics

We conducted experiments on five benchmark datasets, including CIFAR-10 [18], CIFAR-20 [18], STL-10 [19], ImageNet-10 [20], and ImageNet-Dogs [20]. Datasets summary in Table I. Three widely used clustering metrics including Normalized Mutual Information (NMI), Clustering Accuracy (ACC), and Adjusted Rand Index (ARI) are utilized to evaluate our method. Higher scores indicate better clustering performance.

### B. Implementation Details

In our framework, we used ResNet-34 as the backbone to compare with other methods for fairness. To avoid uncertain cluster assignment in the early stage, pre-training is performed by $\ell_d$, and $\ell_s$ is added after 50 epochs, the learning rate keeps

**Algorithm 1** CPCC

**Input:** Dataset $\mathcal{X}$; data transformation $t, t' \sim \mathcal{T}$; epochs $E$, batch size $B$; networks $\mathcal{O}(\cdot), \mathcal{T}(\cdot)$.
**Output:** Clustering result.

1: **for** epoch = 1 to $E$ **do**
2:     Compute the features $\mathbf{Z} = \mathcal{T}(\mathbf{X})$.
3:     Perform $k$-means on $\mathbf{Z}$ and initialize cluster center $\boldsymbol{\mu}$.
4:     Calculate the weights $\mathbf{W}$ by Eq. (2) and Eq. (3).
5:     **for** batch = 1 to $\lfloor \frac{|\mathbf{X}|}{B} \rfloor$ **do**
6:         Transformation of each sample $\mathbf{X}^t, \mathbf{X}^{t'}$.
7:         Compute the features $\mathbf{Z}^t = \mathcal{O}(\mathbf{X}^t), \mathbf{Z}^{t'} = \mathcal{T}(\mathbf{X}^{t'})$.
8:         Compute the prototypes $\mathbf{P}, \mathbf{P}'$ by Eq. (4).
9:         Calculate loss $\ell_s, \ell_d$ by Eq. (5) and Eq. (6).
10:         Update parameters $\boldsymbol{\theta}_{\mathcal{O}}, \boldsymbol{\theta}_{\mathcal{T}}$ by Eq. (7) and Eq. (8).
11:     **end for**
12: **end for**
13: Perform $k$-means on $\mathbf{Z}$ to obtain the final clustering result.



(a) CPCC w/o DCL (b) CPCC w/o SPC      (c) CPCC

Fig. 3. $t$-SNE visualization for CPCC on the CIFAR-10 dataset



(a) NMI          (b) STD

Fig. 4. Ablation study on CIFAR-10

**TABLE I**
**SUMMARY OF THE DATASETS**

| Dataset | Split | Samples | Image size | Classes |
|---|---|---|---|---|
| CIFAR-10 | Train + Test | 60,000 | $32 \times 32$ | 10 |
| CIFAR-20 | Train + Test | 60,000 | $32 \times 32$ | 20 |
| STL-10 | Train + Test | 13,000 | $96 \times 96$ | 10 |
| ImageNet-10 | Train | 13,000 | $96 \times 96$ | 10 |
| ImageNet-Dogs | Train | 19,500 | $96 \times 96$ | 15 |

decaying from 0.05 to 0. The parameters $\alpha = 1$, $\tau = 0.5$, $\lambda = 0.1$, $\sigma = 0.001$ and $m \in [0, 1)$ respectively. The experiment was based on multiple averages as the final result and experimental environment contains one desktop computer with Intel Xeon(R) Silver 4310 CPU, two NVIDIA GeForce RTX 4090 GPUs, 128GB RAM, and coded with the PyTorch deep learning platform on the Ubuntu 22.04 operating system.

### C. Performance Comparison

We report the clustering methods including traditional clustering, $k$-means [21], SC [22], AC [23], NMF [24]; deep clustering methods, AE [25], DAE [26], DCGAN [27], DeCNN [28], VAE [29], JULE [30], DEC [14], DAC [20], DCCM [31], IIC [32], PICA [33] and contrastive learning based deep clustering methods, SCAN [3], GCC [7], NNM [6], CC [4], IDFD [34], SPICE [5], TCL [9], ProPos [13] and RPSC [35], As shown in Table II. the contrastive-based deep clustering methods perform better than general deep clustering methods due to the ability of contrastive learning to learn the discriminative features. In addition, prototype contrastive methods ProPos and CPCC achieve better performance than instance contrastive methods due to the reduced inter-class conflicts. Most importantly, CPCC outperforms the SOTA method in all three metrics on most datasets, further demonstrating the validity of our method.
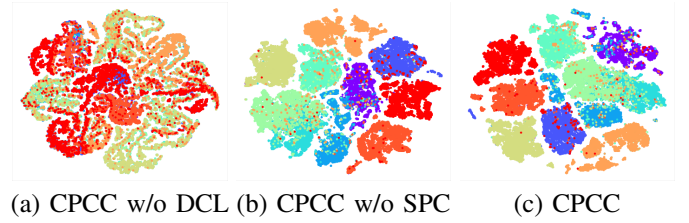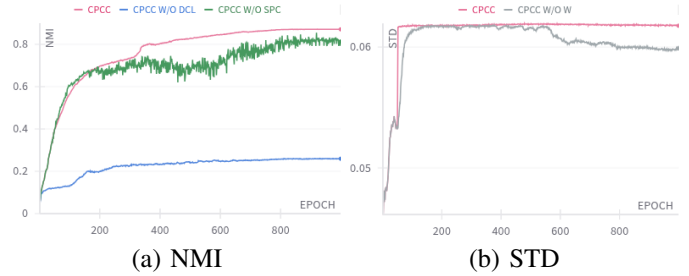
### D. Ablation Study

**Effectiveness of DCL**: To verify the effect of different modules on CPCC, we conducted the experiments shown in Table III, and used BYOL [17] as a baseline for comparison, Where $DCL_1$ and $DCL_2$ represent the first and second parts of Eq. (6) respectively. It is easy to see that both SPC and DCL are indispensable, and without them will lead to performance degradation. In particular, DCL has a great impact on the model's performance, since sample consistency allows the model to learn features that are consistent after sample transformations, which usually contain sample semantic information. We further visualize the features through $t$-SNE, as shown in Fig 3, which further demonstrates the overall distribution of features. The distribution of CPCC is accurate and compact compared to other features.

**Effectiveness of SPC**: In Fig. 4(a), CPCC w/o SPC achieves good performance but the training process is very unstable, SPC avoids inter-class conflicts, which not only improves the performance but also further improves the stability of the training process. In addition, CPCC w/o W means directly calculating the hard prototype, which has a slightly lower performance than CPCC. We combine the STD (standard deviation) of the features in training to reflect the drift of the prototype. The larger the value of STD, the more uniform the distribution of features in space, and the more accurate the estimation of the prototype. As shown in Fig. 4(b), the STD of CPCC with weights is larger and more stable, the more uniform the feature distribution, the smaller the prototype drift, therefore the closer to the real clustering center.

### E. Parameter Sensitivity Analysis

We further analyze the effect of different values of hyperparameters $\tau$ and $\lambda$ on the results as shown in Fig.5. The two

TABLE II
THE CLUSTERING PERFORMANCE ON FIVE DATASETS

| Methods | CIFAR-10 | | | CIFAR-20 | | | STL-10 | | | ImageNet-10 | | | ImageNet-Dogs | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | NMI | ACC | ARI | NMI | ACC | ARI | NMI | ACC | ARI | NMI | ACC | ARI | NMI | ACC | ARI |
| $k$-means [21] | 0.087 | 0.229 | 0.049 | 0.084 | 0.130 | 0.028 | 0.125 | 0.192 | 0.061 | 0.119 | 0.241 | 0.057 | 0.055 | 0.105 | 0.020 |
| SC [22] | 0.103 | 0.247 | 0.085 | 0.090 | 0.136 | 0.022 | 0.098 | 0.159 | 0.048 | 0.151 | 0.274 | 0.076 | 0.038 | 0.111 | 0.013 |
| AC [23] | 0.105 | 0.228 | 0.065 | 0.098 | 0.138 | 0.034 | 0.239 | 0.332 | 0.140 | 0.138 | 0.242 | 0.067 | 0.037 | 0.139 | 0.021 |
| NMF [24] | 0.081 | 0.190 | 0.034 | 0.079 | 0.118 | 0.026 | 0.096 | 0.180 | 0.046 | 0.132 | 0.230 | 0.065 | 0.044 | 0.118 | 0.016 |
| AE [25] | 0.239 | 0.314 | 0.169 | 0.100 | 0.165 | 0.048 | 0.250 | 0.303 | 0.161 | 0.210 | 0.317 | 0.152 | 0.104 | 0.185 | 0.073 |
| DAE [26] | 0.251 | 0.297 | 0.163 | 0.111 | 0.151 | 0.046 | 0.224 | 0.302 | 0.152 | 0.206 | 0.304 | 0.138 | 0.104 | 0.190 | 0.078 |
| DCGAN [27] | 0.265 | 0.315 | 0.176 | 0.120 | 0.151 | 0.045 | 0.210 | 0.298 | 0.139 | 0.225 | 0.346 | 0.157 | 0.121 | 0.174 | 0.078 |
| DeCNN [28] | 0.240 | 0.282 | 0.174 | 0.092 | 0.133 | 0.038 | 0.227 | 0.299 | 0.162 | 0.186 | 0.313 | 0.142 | 0.098 | 0.175 | 0.073 |
| VAE [29] | 0.245 | 0.291 | 0.167 | 0.108 | 0.152 | 0.040 | 0.200 | 0.282 | 0.146 | 0.193 | 0.334 | 0.168 | 0.107 | 0.179 | 0.079 |
| JULE [30] | 0.192 | 0.272 | 0.138 | 0.103 | 0.137 | 0.033 | 0.182 | 0.277 | 0.164 | 0.175 | 0.300 | 0.138 | 0.054 | 0.138 | 0.028 |
| DEC [14] | 0.257 | 0.301 | 0.161 | 0.136 | 0.185 | 0.050 | 0.276 | 0.359 | 0.186 | 0.282 | 0.381 | 0.203 | 0.122 | 0.195 | 0.079 |
| DAC [20] | 0.396 | 0.522 | 0.306 | 0.185 | 0.238 | 0.088 | 0.366 | 0.470 | 0.257 | 0.394 | 0.527 | 0.302 | 0.219 | 0.275 | 0.111 |
| DCCM [31] | 0.496 | 0.623 | 0.408 | 0.285 | 0.327 | 0.173 | 0.376 | 0.482 | 0.262 | 0.608 | 0.710 | 0.555 | 0.321 | 0.383 | 0.182 |
| IIC [32] | 0.513 | 0.617 | 0.411 | 0.225 | 0.257 | 0.117 | 0.431 | 0.499 | 0.295 | - | - | - | - | - | - |
| PICA [33] | 0.591 | 0.696 | 0.512 | 0.310 | 0.337 | 0.171 | 0.611 | 0.713 | 0.531 | 0.802 | 0.870 | 0.761 | 0.352 | 0.352 | 0.201 |
| SCAN [3] | 0.797 | 0.883 | 0.772 | 0.486 | 0.507 | 0.333 | 0.698 | 0.809 | 0.646 | - | - | - | 0.612 | 0.593 | 0.457 |
| GCC [7] | 0.764 | 0.856 | 0.728 | 0.472 | 0.472 | 0.305 | 0.684 | 0.788 | 0.631 | 0.842 | 0.901 | 0.822 | 0.490 | 0.526 | 0.362 |
| NNM [6] | 0.748 | 0.843 | 0.709 | 0.484 | 0.477 | 0.316 | 0.694 | 0.808 | 0.650 | - | - | - | 0.604 | 0.586 | 0.449 |
| CC [4] | 0.705 | 0.790 | 0.637 | 0.431 | 0.429 | 0.266 | 0.764 | 0.850 | 0.726 | 0.859 | 0.893 | 0.822 | 0.445 | 0.429 | 0.274 |
| IDFD [34] | 0.711 | 0.815 | 0.663 | 0.426 | 0.425 | 0.264 | 0.643 | 0.756 | 0.575 | 0.898 | 0.954 | 0.901 | 0.546 | 0.591 | 0.413 |
| SPICE [5] | 0.734 | 0.838 | 0.705 | 0.448 | 0.468 | 0.294 | 0.817 | 0.908 | 0.812 | 0.828 | 0.921 | 0.836 | 0.572 | 0.646 | 0.479 |
| TCL [9] | 0.819 | 0.887 | 0.780 | 0.529 | 0.531 | 0.357 | 0.799 | 0.868 | 0.757 | 0.875 | 0.895 | 0.837 | 0.623 | 0.644 | 0.516 |
| ProPos [13] | <u>0.886</u> | <u>0.943</u> | <u>0.884</u> | <u>0.606</u> | <u>0.614</u> | <u>0.451</u> | 0.758 | 0.867 | 0.737 | <u>0.896</u> | <u>0.956</u> | <u>0.906</u> | <u>0.692</u> | <u>0.745</u> | <u>0.627</u> |
| RPSC [35] | 0.754 | 0.857 | 0.731 | 0.476 | 0.518 | 0.341 | **0.838** | **0.920** | **0.834** | 0.830 | 0.927 | 0.858 | 0.552 | 0.640 | 0.465 |
| Ours | **0.900** | **0.950** | **0.898** | **0.611** | **0.618** | **0.456** | <u>0.825</u> | <u>0.912</u> | <u>0.815</u> | **0.904** | **0.962** | **0.916** | **0.698** | **0.749** | **0.634** |

TABLE III
IMPACT OF DIFFERENT COMBINATIONS ON PERFORMANCE

| Losses | CIFAR-10 | | | CIFAR-20 | | |
|---|---|---|---|---|---|---|
| | NMI | ACC | ARI | NMI | ACC | ARI |
| BYOL [17] | 0.794 | 0.878 | 0.766 | 0.464 | 0.45 | 0.295 |
| CPCC w/o SPC | 0.813 | 0.881 | 0.765 | 0.567 | 0.522 | 0.380 |
| CPCC w/o DCL | 0.259 | 0.332 | 0.145 | 0.168 | 0.201 | 0.065 |
| -CPCC w/o $DCL_1$ | 0.871 | 0.929 | 0.857 | 0.599 | 0.575 | 0.418 |
| -CPCC w/o $DCL_2$ | 0.813 | 0.881 | 0.765 | 0.565 | 0.522 | 0.380 |
| CPCC w/o W | 0.887 | 0.944 | 0.884 | 0.607 | 0.615 | 0.452 |
| CPCC | **0.900** | **0.950** | **0.898** | **0.611** | **0.618** | **0.456** |



Fig. 6. Convergence Analysis on CIFAR-10

the clustering effect is best when $\lambda = 0.1$. For the rest of the hyperparameters, we take empirical values.

### F. Convergence Analysis

Fig. 6 shows the convergence of $\ell_d$ and $\ell_s$. Although $\ell_s$ has been fluctuating, it is generally on a downward trend. After 50 epochs, we added $\ell_s$ to train the network together, $\ell_s$ drops rapidly and remains stable, indicating that the prototypes are well separated and the stability of training is effectively maintained.

## V. CONCLUSION

In this paper, we propose a prototype contrast clustering method to solve the inter-class conflict problem of contrastive learning-based clustering methods. By combining soft prototype contrast with feature consistency learning, we can avoid inter-class conflicts and reduce inaccurate prototype calculation caused by prototype drift, effectively improving clustering performance and stability. Experimental results prove the superiority of the proposed method.



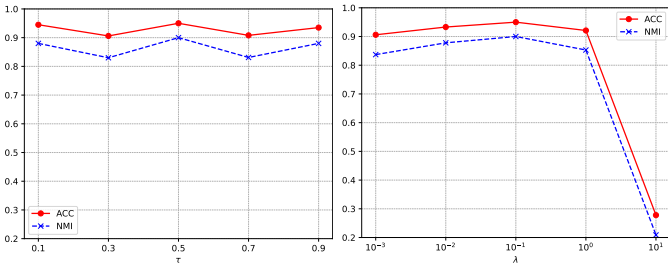(a) Analysis on value of $\tau$     (b) Analysis on value of $\lambda$

Fig. 5. Parameter sensitivity analysis on CIFAR-10

correspond to Eq. (5) and Eq. (7), respectively. The parameter $\tau$ is used to scale the similarity between sample pairs, and its value has little effect on the CPCC method. The parameter $\lambda$ is used to adjust the balance between the $\ell_s$ and the $\ell_d$, and

## REFERENCES

[1] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 9729–9738.

[2] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton, "A simple framework for contrastive learning of visual representations," in *International conference on machine learning*, 2020, pp. 1597–1607.

[3] Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, Marc Proesmans, and Luc Van Gool, "Scan: Learning to classify images without labels," in *European conference on computer vision*. Springer, 2020, pp. 268–285.

[4] Yunfan Li, Peng Hu, Zitao Liu, Dezhong Peng, Joey Tianyi Zhou, and Xi Peng, "Contrastive clustering," in *Proceedings of the AAAI conference on artificial intelligence*, 2021, vol. 35, pp. 8547–8555.

[5] Chuang Niu, Hongming Shan, and Ge Wang, "Spice: Semantic pseudo-labeling for image clustering," *IEEE Transactions on Image Processing*, vol. 31, pp. 7264–7278, 2022.

[6] Zhiyuan Dang, Cheng Deng, Xu Yang, Kun Wei, and Heng Huang, "Nearest neighbor matching for deep clustering," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 13693–13702.

[7] Huasong Zhong, Jianlong Wu, Chong Chen, Jianqiang Huang, Minghua Deng, Liqiang Nie, Zhouchen Lin, and Xian-Sheng Hua, "Graph contrastive clustering," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 9224–9233.

[8] Han Zhao, Xu Yang, Zhenru Wang, Erkun Yang, and Cheng Deng, "Graph debiased contrastive learning with joint representation clustering," in *IJCAI*, 2021, pp. 3434–3440.

[9] Yunfan Li, Mouxing Yang, Dezhong Peng, Taihao Li, Jiantao Huang, and Xi Peng, "Twin contrastive learning for online clustering," *International Journal of Computer Vision*, vol. 130, no. 9, pp. 2205–2221, 2022.

[10] Yue Liu, Xihong Yang, Sihang Zhou, Xinwang Liu, Zhen Wang, Ke Liang, Wenxuan Tu, Liang Li, Jingcan Duan, and Cancan Chen, "Hard sample aware network for contrastive deep graph clustering," in *Proceedings of the AAAI conference on artificial intelligence*, 2023, vol. 37, pp. 8914–8922.

[11] Yiding Lu, Yijie Lin, Mouxing Yang, Dezhong Peng, Peng Hu, and Xi Peng, "Decoupled contrastive multi-view clustering with high-order random walks," in *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Vancouver, Canada*, 2024, pp. 14193–14201.

[12] Junnan Li, Pan Zhou, Caiming Xiong, and Steven C. H. Hoi, "Prototypical contrastive learning of unsupervised representations," in *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*, 2021.

[13] Zhizhong Huang, Jie Chen, Junping Zhang, and Hongming Shan, "Learning representation for clustering via prototype scattering and positive sampling," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 45, no. 6, pp. 7509–7524, 2023.

[14] Junyuan Xie, Ross Girshick, and Ali Farhadi, "Unsupervised deep embedding for clustering analysis," in *International conference on machine learning*. PMLR, 2016, pp. 478–487.

[15] Zhuxi Jiang, Yin Zheng, Huachun Tan, Bangsheng Tang, and Hanning Zhou, "Variational deep embedding: An unsupervised and generative approach to clustering," in *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*, 2017, pp. 1965–1972.

[16] Sudipto Mukherjee, Himanshu Asnani, Eugene Lin, and Sreeram Kannan, "Clustergan: Latent space clustering in generative adversarial networks," in *Proceedings of the AAAI conference on artificial intelligence*, 2019, vol. 33, pp. 4610–4617.

[17] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al., "Bootstrap your own latent-a new approach to self-supervised learning," *Advances in neural information processing systems*, vol. 33, pp. 21271–21284, 2020.

[18] Alex Krizhevsky, Geoffrey Hinton, et al., "Learning multiple layers of features from tiny images," 2009.

[19] Adam Coates, Andrew Ng, and Honglak Lee, "An analysis of single-layer networks in unsupervised feature learning," in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2011, pp. 215–223.

[20] Jianlong Chang, Lingfeng Wang, Gaofeng Meng, Shiming Xiang, and Chunhong Pan, "Deep adaptive image clustering," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 5879–5887.

[21] James MacQueen et al., "Some methods for classification and analysis of multivariate observations," in *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*. Oakland, CA, USA, 1967, vol. 1, pp. 281–297.

[22] Lihi Zelnik-Manor and Pietro Perona, "Self-tuning spectral clustering," *Advances in neural information processing systems*, vol. 17, 2004.

[23] K Chidananda Gowda and GJPR Krishna, "Agglomerative clustering using the concept of mutual nearest neighbourhood," *Pattern recognition*, vol. 10, no. 2, pp. 105–112, 1978.

[24] Deng Cai, Xiaofei He, Xuanhui Wang, Hujun Bao, and Jiawei Han, "Locality preserving nonnegative matrix factorization," in *Twenty-first international joint conference on artificial intelligence*, 2009.

[25] Yoshua Bengio, Pascal Lamblin, Dan Popovici, and Hugo Larochelle, "Greedy layer-wise training of deep networks," *Advances in neural information processing systems*, vol. 19, 2006.

[26] Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, Pierre-Antoine Manzagol, and Léon Bottou, "Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion.," *Journal of machine learning research*, vol. 11, no. 12, 2010.

[27] Alec Radford, Luke Metz, and Soumith Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," in *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4*, 2016.

[28] Matthew D. Zeiler, Dilip Krishnan, Graham W. Taylor, and Robert Fergus, "Deconvolutional networks," in *The Twenty-Third IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2010, San Francisco, CA, USA, 13-18 June 2010*, 2010, pp. 2528–2535.

[29] Diederik P. Kingma and Max Welling, "Auto-encoding variational bayes," in *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16*, 2014.

[30] Jianwei Yang, Devi Parikh, and Dhruv Batra, "Joint unsupervised learning of deep representations and image clusters," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 5147–5156.

[31] Jianlong Wu, Keyu Long, Fei Wang, Chen Qian, Cheng Li, Zhouchen Lin, and Hongbin Zha, "Deep comprehensive correlation mining for image clustering," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 8150–8159.

[32] Xu Ji, Joao F Henriques, and Andrea Vedaldi, "Invariant information clustering for unsupervised image classification and segmentation," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 9865–9874.

[33] Jiabo Huang, Shaogang Gong, and Xiatian Zhu, "Deep semantic clustering by partition confidence maximisation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 8849–8858.

[34] Yaling Tao, Kentaro Takagi, and Kouta Nakata, "Clustering-friendly representation learning via instance discrimination and feature decorrelation," in *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7*, 2021.

[35] Sihang Liu, Wenming Cao, Ruigang Fu, Kaixiang Yang, and Zhiwen Yu, "RPSC: robust pseudo-labeling for semantic clustering," in *Thirty-Eighth AAAI Conference on Artificial Intelligence, Vancouver, Canada*, 2024, pp. 14008–14016.