
UCCC: Unsupervised Community-consensus Contrastive Clustering

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 As one of the most important tasks in computer vision, unsupervised image classi-
2 fication aims to group images into semantically meaningful clusters without using
3 any labels. In this paper, we propose a one-stage clustering method called *Unsuper-*
4 *vised Community-consensus Contrastive Clustering* (UCCC), which performs both
5 instance- and cluster-level contrastive learning. In our framework, instance-level
6 contrastive learning is capable of learning discriminative features, thereby helping
7 construct reliable communities; cluster-level contrastive learning is conducted with
8 the aid of the established communities, and further produces community-consensus
9 cluster predictions. In particular, we design a novel instance-based loss function for
10 the cluster-level contrastive learning. We demonstrate analytically that the gradient
11 of our loss function could alleviate cluster degeneracy and thus prevent from a
12 trivial solution, where the clusters are collapsed into a single entity. Extensive
13 experimental results show that UCCC consistently outperforms state-of-the-art
14 methods on six benchmark datasets.

15 1 Introduction

16 Deep neural networks have achieved human-level accuracy in image classification with the aid of
17 large-scale datasets that contain annotated images, i.e. images with their corresponding semantic
18 label. Nevertheless, annotating sufficient data is labor-intensive and time-consuming, establishing
19 significant barriers for adapting the image classification systems to new domains. As a result, the
20 focus of researchers is shifting to how to tackle image classification in an unsupervised manner.
21 Some works [1–4] utilize the architecture of neural networks as a prior to cluster images, and refine
22 the clusters iteratively by deriving the supervisory signal from the most confident sample [1, 2] or
23 through cluster re-assignments calculated offline [3, 4]. Though this kind of two-stage methods could
24 jointly learn representations and perform clustering, the errors accumulated during the alternation
25 might result in sub-optimal clustering performance. On the other hand, the newly proposed CC [5]
26 manages to learn discriminative representation and perform clustering simultaneously. The key idea
27 of CC is to consider both instance- and cluster-level contrastive learning. To be specific, for a given a
28 dataset, the positive and negative instance pairs are constructed through data augmentations, where
29 the positive one is composed of two augmented views of the same instance and the other pairs are
30 defined to be negative. By gathering the positive pairs and scattering the negatives pairs, the instance-
31 and cluster-level contrastive learning are conducted in the row and column space of the feature matrix.
32 In other words, the rows and columns of a feature matrix are regarded as the instance representations
33 and cluster representations, respectively. However, using cluster representations for cluster-level
34 contrastive learning may aggravate cluster degeneracy [3] and thus lead to a trivial solution, where
35 the clusters are collapsed into a single entity. To deal with the aforementioned issue, CC requires an
36 additional balance loss to maximize the entropy of cluster assignment probabilities, which however
37 turns out to be the competing between the contrastive loss and the balance term. In contrast, we

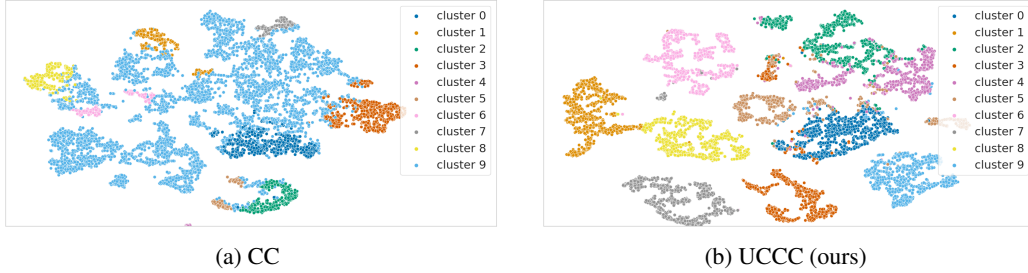


Figure 1: The t-SNE visualization of cluster-based contrastive clustering like CC [5] and our proposed instance-based contrastive clustering UCCC. Both models are trained without the balance term. It can be observed CC (Fig. 1a) suffers from severe degeneracy problem while UCCC (Fig.1b), as an instance-based contrastive clustering method, has no such issue.)

38 propose a one-stage clustering method called *Unsupervised Community-consensus Contrastive*
 39 *Clustering* (UCCC), which also adopts a dual contrastive learning framework but introduces a novel
 40 instance-based loss function for cluster-level contrastive learning.

41 Since objects in the same group are more similar to each other than to those in other groups, similar-
 42 looking objects usually belong to the same cluster while objects that can be easily distinguished tend
 43 to come from different clusters. Motivated by this observation, for a given instance, we expect its
 44 cluster prediction close to the one estimated based on its similar-looking positive community; on
 45 the other hand, the unlike-looking negative community should give to different cluster prediction.
 46 In our framework, instance-level contrastive learning is capable of learning discriminative features,
 47 thereby helping construct reliable communities; cluster-level contrastive learning is conducted
 48 with the aid of the established communities, and further produces community-consensus cluster
 49 prediction. In particular, we design the cluster-level contrastive loss function to: 1) maximize the
 50 prediction similarity between the positive instance pairs, and 2) minimize the positive-negative
 51 prediction similarity between the negative instance pairs. Doing so allows us to encourage not only
 52 self-consistent cluster prediction but also consensual cluster predictions of community.

53 The main contributions of this work are as follows: (1) We propose a one-stage clustering method
 54 called *Unsupervised Community-consensus Contrastive Clustering* (UCCC), in which a novel
 55 instance-based loss function for cluster-level contrastive learning is introduced. (2) We analytically
 56 compare the cluster- and instance-based contrastive loss function, showing the former to be a
 57 factor of data collapse while the latter could effectively alleviate this problem. Empirical results also
 58 prove that our approach does not suffer from data collapse due to the designed instance-based loss
 59 function. (3) Extensive experiments on six benchmark datasets show that our approach outperforms
 60 state-of-the-art methods in terms of three widely used clustering metrics, i.e., normalized mutual
 61 information (NMI), adjusted rand index (ARI) and accuracy (ACC).

62 2 Related Work

63 **Unsupervised Image Clustering.** As one of the most important tasks in computer vision, unsu-
 64 pervised image classification aims to group images into semantically meaningful clusters without
 65 using any labels. Recently, Van Gansbeke et al. propose a two-staged approach, SCAN, where feature
 66 learning and clustering are decoupled. In particular, SCAN first employs a self-supervised task [7]
 67 to obtain high-level feature representations then clusters those representations by nearest neighbors.
 68 Another work that incorporates self-supervised representation learning into clustering is CC [5],
 69 which takes semantic labels as a special representation and conducts the instance- and the cluster-level
 70 contrastive learning simultaneously. In this work we adopt a dual contrastive framework like CC as
 71 multi-staged approaches take much more time to deploy and hardly achieves improvements compared
 72 with one-stage methods. The main difference between CC and ours is in how we perform cluster-level
 73 contrastive learning. CC conducts cluster-level learning in the column space of representation vectors,
 74 which would aggravate cluster degeneracy and thus require additional balance term. Instead, we
 75 perform cluster-level contrastive learning in the row space of representations by adopting a novel loss
 76 function, leading to substantial improvements in the clustering performance.

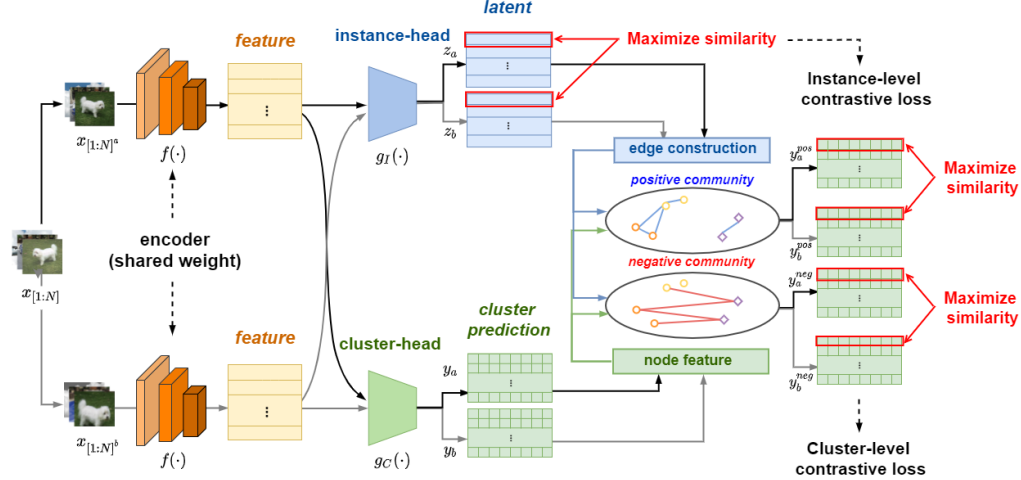


Figure 2: Our unsupervised community-consensus contrastive clustering (UCCC) framework consists of an encoder f and two MLPs that correspond to an instance head g_I and a cluster head g_C . Two random data augmentations are applied on each input image to obtain data pairs. Given data pairs, the shared encoder is used to extract features from different augmentations. These feature are projected into two subspace to conduct instance- and cluster-level contrastive learning using the corresponding projection head. Taking the features in two subspace as prior, two community graph are constructed to produce consensual cluster predictions with the guidance of the cluster-level contrastive loss.

77 3 Unsupervised Community-Consensus Contrastive Clustering

78 As shown in Fig. 2 our model consist of two sub-branch: the Instance-level Contrastive Branch (ICB)
 79 and the Cluster-level Contrastive Branch (CCB). Both branches take the output of the encoder $f(\cdot)$ as
 80 input. ICB projects the input feature into a low-dimensional space where the corresponding contrastive
 81 loss is applied. The discriminative features learned by ICB can not only attain inter-cluster difference
 82 but also preserve intra-cluster distinction, and thus help construct reliable instance communities
 83 for CCB. Since the semantic label can be regarded as a special representation, CCB projects the
 84 input instances into a subspace with a dimensionality of the cluster number, and consensual cluster
 85 predictions are further achieved with the guidance of our proposed contrastive loss. In the following,
 86 we will describe how we perform a dual contrastive learning for unsupervised clustering in detail and
 87 introduce the proposed objective function at the end.

88 3.1 Instance-level Contrastive Learning

89 To learn representations without labels, we leverage a self-supervised approach SimCLR [7], which
 90 uses “instance discrimination” as a pretext task. Given a minibatch of images $\{\mathbf{x}_i\}_{i=1}^N$, we apply
 91 random image transformations (e.g., cropping or blurring) twice on each image, thus generating
 92 two different view of them (augmentation a and b). The transformed images are projected to a
 93 subspace via $\mathbf{z} = g_I(f(\mathbf{x}))$, where g_I is a two-layer MLP projection head. The resulting $2N$ data
 94 points $\{\mathbf{z}_{1^a}, \mathbf{z}_{2^a}, \dots, \mathbf{z}_{N^a}, \mathbf{z}_{1^b}, \mathbf{z}_{2^b}, \dots, \mathbf{z}_{N^b}\}$ will be used to calculate the contrastive loss as described
 95 below.

96 **Instance-level Contrastive Loss.** The common idea of contrastive learning is the following: pull
 97 together an anchor and a “positive” sample (minimize the similarity between a positive pair), and
 98 push apart the anchor from many “negative” samples (maximize the similarity between negative
 99 pairs). A positive pair often consists of instances augmented from the same sample, and negative
 100 pairs are formed by the anchor and randomly chosen samples from the minibatch. Here, we use the
 101 dot product between $L2$ normalized features, which is cosine similarity, as the metric for instance
 102 similarity. Let $\tilde{\mathbf{z}} = \frac{\mathbf{z}}{\|\mathbf{z}\|}$ denote the normalized feature. Then instance-level contrastive loss for an

103 anchor i is defined as:

$$\ell_{i^a} = -\log \frac{\exp(\tilde{\mathbf{z}}_{i^a} \cdot \tilde{\mathbf{z}}_{i^b}/\tau)}{\sum_{j=1}^N \mathbb{1}_{[j \neq i]} [\exp(\tilde{\mathbf{z}}_{i^a} \cdot \tilde{\mathbf{z}}_{j^a}/\tau) + \exp(\tilde{\mathbf{z}}_{i^a} \cdot \tilde{\mathbf{z}}_{j^b}/\tau)]}, \quad (1)$$

104 where $\mathbb{1}_{[j \neq i]} \in \{0, 1\}$ is an indicator function evaluating to 1 iff $j = i$ and τ is the instance-level
 105 temperature parameter. Considering every augmented samples within a minibatch, the instance-level
 106 loss is computed as:

$$\mathcal{L}_{ins} = \sum_{i=1}^N \ell_{i^a} + \ell_{i^b}. \quad (2)$$

107 3.2 Cluster-level Contrastive Learning

108 In cluster-level learning, we leverage the concept of "community consensus". That is, for a given
 109 instance, its similar-looking positive peers would come from the same cluster as itself; on the other
 110 hand, the unlike-looking negative peers should belong to different clusters. Following this idea, we
 111 first construct a positive graph and a negative graph to present the community relations.

112 **Community Peers.** Image distinction is effectively preserved in the feature space learned by ICB
 113 since the instance-level contrastive loss encourages high similarity only between the instances aug-
 114 mented from the same image. Motivated by this observation, we use the similarities of representations
 115 learned by the instance head to construct positive/negative communities.

116 Specifically, we regard the community relations
 117 as a graph structure and define the positive
 118 and negative adjacency matrices, $\mathbf{A}_{pos}, \mathbf{A}_{neg} \in$
 119 $\mathbb{R}^{2N \times 2N}$. As there might be ambiguous pairs
 120 that are not similar enough, a confidence thresh-
 121 old is required to filter out noisy information.
 122 We set the positive confidence threshold by the
 123 cosine value of an angle θ , where θ is calculated based on the number of clusters, C . Larger C usually
 124 implies more hardly distinguishable noisy pairs exist, and thus requires higher positive confidence.
 125 Table 1 summarizes the considered positive confidence in this work. See supplementary Sec. A for
 126 more details about the angle θ . The elements of positive adjacency matrix are thereby defined as:

Table 1: Positive confidence threshold with respect to the number of clusters.

# of clusters	10	15	20	200
pos_conf	0.258	0.484	0.615	0.961

$$A_{ij,pos} = \begin{cases} \text{sim}(\tilde{\mathbf{z}}_i, \tilde{\mathbf{z}}_j), & \text{if } \text{sim}(\tilde{\mathbf{z}}_i, \tilde{\mathbf{z}}_j) > \text{pos_conf} \\ 0, & \text{otherwise} \end{cases}, \quad (3)$$

127 where the augmentation notations a, b of features $\tilde{\mathbf{z}}_i, \tilde{\mathbf{z}}_j$ are ignored for simplification. Similarly, the
 128 elements of negative adjacency matrix are calculated as shown in Eq. 4. Since diverse negative peers
 129 could help the model learn how to distinguish the belonging cluster from other clusters, the negative
 130 confidence is set as 0, which is a relatively low value in comparison with the positive one.

$$A_{ij,neg} = \begin{cases} -\text{sim}(\tilde{\mathbf{z}}_i, \tilde{\mathbf{z}}_j), & \text{if } \text{sim}(\tilde{\mathbf{z}}_i, \tilde{\mathbf{z}}_j) < \text{neg_conf} \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

131 **Assignment Aggregation.** To produce cluster assignments for given samples, the input features
 132 are mapped into a subspace with a dimensionality of the cluster number via $\mathbf{y} = g_C(f(\mathbf{x}))$, where
 133 $\mathbf{y} \in \mathbb{R}^C$ denotes a cluster assignment (i -th element can be interpreted as the probability of sample \mathbf{x}
 134 being assigned to the cluster i) and g_C is a two-layer projection head. Now that the community graphs
 135 are established, we would like aggregate the information from community relations, and estimate a
 136 positive and a negative cluster assignment for each instance according to its corresponding peers. For
 137 a mini-batch of augmented samples, its cluster assignments $\{\mathbf{y}_{i^m}\}_{i=1, m \in \{a, b\}}^N$ can be concatenated
 138 into a matrix $\mathbf{Y} \in \mathbb{R}^{2N \times C}$. We consider k -hop peer relations and compute the positive/negative
 139 assignment matrix inspired from [8]:

$$\begin{cases} \mathbf{Y}_{pos} = \left(\prod_{i=1}^k \tilde{\mathbf{A}}_{pos} \right) \mathbf{Y} \\ \mathbf{Y}_{neg} = \tilde{\mathbf{A}}_{neg} \left(\prod_{i=1}^{k-1} \tilde{\mathbf{A}}_{pos} \right) \mathbf{Y} \end{cases}, \quad (5)$$

Algorithm 1: Unsupervised Community-consensus Contrastive Clustering

Input: Dataset \mathcal{X} , cluster number C , Training steps S , Batch size N , Temperature τ , confidence threshold δ , structure of encoder f , instance head g_I , and cluster head g_C

Output: cluster assignments of dataset \mathcal{X}

```
/* training */
for step  $s = 1$  to  $S$  do
    sample a minibatch  $\{\mathbf{x}_i\}_{i=1}^N$  from  $\mathcal{X}$ 
    obtain  $2N$  augmented samples  $\{\mathbf{x}_{im}\}_{i=1, m \in \{a, b\}}^N$  through data augmentation
    compute instance representations and cluster assignments by
         $\mathbf{z}_{im} = g_I(f(\mathbf{x}_{im}))$  and  $\mathbf{y}_{im} = g_C(f(\mathbf{x}_{im}))$ 
    construct the positive/negative adjacency matrix through Eq. 3-4
    compute positive and negative cluster assignments through Eq. // community fusion
    compute the instance-level loss  $\mathcal{L}_{ins}$  and cluster-level loss  $\mathcal{L}_{clu}$  through Eq. 1-2, 6
    compute overall loss  $\mathcal{L}$  by Eq. 7
    update  $f$ ,  $g_I$  and  $g_C$  to minimize  $\mathcal{L}$ 
/* testing */
for  $\mathbf{x} \in \mathcal{X}$  do
    assign  $\mathbf{x}$  to cluster  $c = \arg \max g_C(f(\mathbf{x}))$ 
```

140 with the normalized adjacency matrices $\tilde{\mathbf{A}}_* = \mathbf{D}_*^{-\frac{1}{2}} \mathbf{A}_* \mathbf{D}_*^{-\frac{1}{2}}$, $* \in \{pos, neg\}$ and \mathbf{D} denoting a
141 degree matrix where $D_{ii} = \sum_j A_{ij}$.

142 **Community-consensus Contrastive Clustering Loss.** For an instance, the positive cluster assign-
143 ment can be viewed as its cluster assignment while further taking positive peers into consideration; the
144 negative cluster assignment, on the other hand, corresponds to the least possible cluster assignment
145 that it would have. In this sense, we follow the idea of "community consensus" and propose a
146 novel instance-based contrastive loss as shown in Eq. 6. In particular, the proposed loss function
147 encourages self-consistent predictions by maximize the similarity between the assignments of positive
148 instance pairs (augmented from the same image). Moreover, consensual predictions are achieved by
149 minimizing the similarity between the positive and the negative assignments of all instance pairs (any
150 two of instances within a mini-batch).

$$\mathcal{L}_{clu} = -\log \frac{\sum_{i=1}^N \exp(\mathbf{y}_{i^a, pos} \cdot \mathbf{y}_{i^b, pos} + \mathbf{y}_{i^a, neg} \cdot \mathbf{y}_{i^b, neg})}{\sum_{\forall (i, j, m, n) \in \Phi} \exp(\mathbf{y}_{i^m, pos} \cdot \mathbf{y}_{j^n, neg})}, \quad (6)$$

151 where $\Phi = \{(i, j, m, n) \mid \forall i, j \in \{1, \dots, N\} \text{ and } \forall m, n \in \{a, b\}\}$.

152 3.3 Overall Objective

153 Algorithm 1 summarizes the full training and test process of the model. Both ICB and CCB are
154 simultaneously optimized by the corresponding loss in an end-to-end manner. Hence, the overall
155 objective function is written as:

$$\mathcal{L} = \mathcal{L}_{ins} + \mathcal{L}_{clu} \quad (7)$$

156 4 Contrastive Losses for clustering: cluster-based and instance-based

157 In this section, we look at the cluster-based contrastive loss proposed in [5] and our instance-based
158 clustering loss, showing why our formulation is superior to the former one.

159 Formally, let $\mathbf{Y}_m \in \mathbb{R}^{N \times C}$ be the cluster assignments for a mini-batch under some augmentation m .
160 Since each sample belongs to only one cluster, the rows of \mathbf{Y}_m tends to be one-hot. In this sense,
161 the i -th column of \mathbf{Y}_m , denoted \mathbf{y}_m^i , can be viewed as a representation of the i -th cluster. For each
162 cluster c , a positive cluster pair is formed by combining \mathbf{y}_a^c and \mathbf{y}_b^c , namely, the representations under

163 two different augmentation, while other $2C - 2$ pairs are considered to be negative. To minimize the
 164 inter-cluster similarities to separate different clusters, the cluster-based contrastive loss is written as:

$$\mathcal{L}_{clu}^{clu} = - \sum_{\substack{c=1, \\ m \in \{a,b\}}}^C \log \frac{\exp(\mathbf{y}_a^c \cdot \mathbf{y}_b^c)}{\sum_{j=1}^C \mathbb{1}_{[i \neq c]} \left[\exp(\mathbf{y}_m^c \cdot \mathbf{y}_a^j) + \exp(\mathbf{y}_m^c \cdot \mathbf{y}_b^j) \right]} \quad (8)$$

165 It is observed that the cluster-based loss is conducted in the column space of assignment matrices,
 166 which is different from our instance-based clustering loss (Eq. 6) that measures the similarity in the
 167 row space.

168 As stated in [5] the cluster-based contrastive loss requires the following balance term to avoid the
 169 trivial solution that most instances are assigned to the same cluster.

$$\mathcal{L}_{balance} = \sum_{c=1}^C [P(\mathbf{y}_a^c) \log P(\mathbf{y}_a^c) + P(\mathbf{y}_b^c) \log P(\mathbf{y}_b^c)] \quad (9)$$

170 Here, $\mathcal{L}_{balance}$ is the negative entropy of assignment probabilities $P(\mathbf{y}_m^c) = \sum_{i=1}^N y_{im}^c / \|\mathbf{Y}\|_1$, $m \in$
 171 $\{a, b\}$ within a mini-batch under each data augmentation. Such kind of entropy maximization would
 172 lead to a sub-optimal clustering result since there is no guarantee that all clusters within a dataset
 173 should be equal-sized. Our instance-based loss in contrast with \mathcal{L}_{clu}^{clu} does not suffer from this issue
 174 and is capable of alleviating degeneracy. An analysis of the gradients with respect to the weights
 175 of the last layer of the cluster head \mathbf{w}^c (can be interpreted as a classifier for cluster c) supports this
 176 conclusion. Throughout the analysis, we assume the following assumption for simplification:

177 **Assumption 1** (Generalization). Let a, b be any two random augmentations. The statements below
 178 hold:

$$\begin{cases} \mathbf{y}_a^j \approx \mathbf{y}_b^j \approx \hat{\mathbf{y}}^j & \text{for } j \in \{1, \dots, C\} \\ \mathbf{H}_a \approx \mathbf{H}_b \approx \hat{\mathbf{H}} \end{cases},$$

179 where \mathbf{H}_m denotes the embedding matrix (prior to the last layer of cluster head) for a mini-batch
 180 under augmentation m .

181 Let $\mathbf{Y}_{m,pos}$ be the positive cluster assignments for a mini-batch under augmentation m and the
 182 corresponding embedding matrix $\mathbf{H}_{m,pos}$ is a combination of transformed \mathbf{H}_a and \mathbf{H}_b as positive
 183 cluster assignments are simply linear combinations of original cluster assignments within a mini-
 184 batch. In a similar fashion, the negative cluster assignments for a mini-batch under augmentation
 185 m and its embedding matrix are denoted by $\mathbf{Y}_{m,neg}$ and $\mathbf{H}_{m,neg}$ respectively. It is noted that the
 186 assumption of model generalization also implies:

$$\begin{cases} \mathbf{y}_{a,*}^j \approx \mathbf{y}_{b,*}^j \approx \hat{\mathbf{y}}_*^j & \text{for } j \in \{1, \dots, C\}, * \in \{pos, neg\} \\ \mathbf{H}_{a,*} \approx \mathbf{H}_{b,*} \approx \hat{\mathbf{H}}_* \end{cases}$$

187 As shown in the Supplementary, the gradient for our clustering loss \mathcal{L}_{clu} is derived as:

$$\frac{\partial \mathcal{L}_{clu}}{\partial \mathbf{w}^c} \approx -\frac{2}{N} \left[\hat{\mathbf{H}}_{pos}^\top \hat{\mathbf{y}}_{pos}^c + \hat{\mathbf{H}}_{neg}^\top \hat{\mathbf{y}}_{neg}^c \right] + \frac{1}{\Psi} \left[\left(\hat{\mathbf{E}} \hat{\mathbf{H}}_{pos} \right)^\top \hat{\mathbf{y}}_{neg}^c + \left(\hat{\mathbf{E}} \hat{\mathbf{H}}_{neg} \right)^\top \hat{\mathbf{y}}_{pos}^c \right], \quad (10)$$

188 with $\Psi = \sum_{i=1}^N \sum_{j=1}^N \exp(\hat{\mathbf{y}}_{i,pos} \cdot \hat{\mathbf{y}}_{j,neg})$ and $\hat{\mathbf{E}} \in \mathbb{R}^{N \times N}$ where $\hat{E}_{ij} = \exp(\hat{\mathbf{y}}_{i,pos} \cdot \hat{\mathbf{y}}_{j,neg})$.

189 The discussion about the cluster-based contrastive loss is in the Supplementary.

190 **Discussion.** Since we use gradient descent [9] as our optimization algorithm, the classifiers \mathbf{w}^c is
 191 updated through $\mathbf{w}^c = \mathbf{w}^c - \eta \frac{\partial \mathcal{L}_{clu}}{\partial \mathbf{w}^c}$ every step. To understand why our formulation does not lead
 192 to cluster degeneracy, we make the following observations. First, if an instance has a high assignment
 193 score (either positive or negative) with respect to a cluster c , the negative term in Eq. 10 would pull
 194 the classifier \mathbf{w}^c close to its corresponding embedding features while the positive term would push
 195 \mathbf{w}^c away from its opposite embedding. Second, because the community graph is established based
 196 on instance-level features, the community peers are sufficiently reliable to estimate positive/negative
 197 assignment, which guarantees $\hat{\mathbf{Y}}_{neg}$ is not a zero matrix. This makes sure the growth of each cluster
 198 and prevents most instances from falling into the same entity.

Table 2: State-of-the-art comparison: The performance (%) of our model are reported in **bold font**. For fair comparison with SCAN [6], we also report the performance of our method with ResNet18 as the backbone of encoder.

Dataset	CIFAR-10			CIFAR-20			STL-10			ImageNet-10			ImageNet-Dogs			tiny-ImageNet		
	NMI	ARI	ACC	NMI	ARI	ACC	NMI	ARI	ACC	NMI	ARI	ACC	NMI	ARI	ACC	NMI	ARI	ACC
K-means [15]	8.7	4.9	22.9	8.4	2.8	13.0	12.5	6.1	19.2	11.9	5.7	24.1	5.5	2.0	10.5	6.5	0.5	2.5
SC [16]	10.3	8.5	24.7	9.0	2.2	13.6	9.8	4.8	15.9	15.1	7.6	27.4	3.8	1.3	11.1	6.3	0.4	2.2
AC [17]	10.5	6.5	22.8	9.8	3.4	13.8	23.9	14.0	33.2	13.8	6.7	24.2	3.7	2.1	13.9	6.9	0.5	2.7
NMF [18]	8.1	3.4	19.0	7.9	2.6	11.8	9.6	4.6	18.0	13.2	6.5	23.0	4.4	1.6	11.8	7.2	0.5	2.9
AE [19]	23.9	16.9	31.4	10.0	4.8	16.5	25.0	16.1	30.3	21.0	15.2	31.7	10.4	7.3	18.5	13.1	0.7	4.1
DAE [20]	25.1	16.3	29.7	11.1	4.6	15.1	22.4	15.2	30.2	20.6	13.8	30.4	10.4	7.3	18.5	12.7	0.7	3.9
DCGAN [21]	26.5	17.6	31.5	12.0	4.5	15.1	22.4	15.2	30.2	22.5	15.7	34.6	12.1	7.8	17.4	13.5	0.7	4.1
DeCNN [22]	24.0	17.4	28.2	9.2	3.8	13.3	22.7	16.2	29.9	18.6	14.2	31.3	9.8	7.3	17.5	11.1	0.6	3.5
VAE [23]	24.5	16.7	29.1	10.8	4.0	15.2	20.0	14.6	28.2	19.3	16.8	33.4	10.7	7.9	17.9	11.3	0.6	3.6
JULE [24]	19.2	13.8	27.2	10.3	3.3	13.7	18.2	16.4	27.7	17.5	13.8	30.0	5.4	2.8	13.8	10.2	0.6	3.3
DEC [1]	25.7	16.1	30.1	13.6	5.0	18.5	27.6	18.6	35.9	28.2	20.3	38.1	12.2	7.9	19.5	11.5	0.7	3.7
DAC [2]	39.6	30.6	52.2	18.5	8.8	23.8	36.6	25.7	47.0	39.4	30.2	52.7	21.9	11.1	27.5	19.0	1.7	6.6
ADC [25]	-	-	32.5	-	-	16.0	-	-	53.0	-	-	-	-	-	-	-	-	-
DDC [26]	42.4	32.9	52.4	-	-	-	37.1	26.7	48.9	43.3	34.5	57.7	-	-	-	-	-	-
DCCM [27]	49.6	40.8	62.3	28.5	17.3	32.7	37.6	26.2	48.2	60.8	55.5	71.0	32.1	18.2	38.3	22.4	3.8	10.8
IIC * [28]	-	-	61.7	-	-	25.7	-	-	59.6	-	-	-	-	-	-	-	-	-
EmbedUL [29]	-	-	81.0	-	-	35.3	-	-	66.5	-	-	-	-	-	-	-	-	-
PICA [30]	59.1	51.2	69.6	31.0	17.1	33.7	61.1	53.1	71.3	80.2	76.1	87.0	35.2	20.1	35.2	27.7	4.0	9.8
CC * [5]	70.5	63.7	79.0	43.1	26.6	42.9	76.4	72.6	85.0	85.9	82.2	89.3	44.5	27.4	42.9	34.0	7.1	14.0
SCAN [6]	71.5	66.5	81.6	44.9	28.3	44.0	67.3	61.8	79.2	-	-	-	-	-	-	-	-	-
UCCC (Ours) * ↓																		
Res18 *	74.8	71.0	84.2	47.4	31.3	46.9	73.7	70.8	84.7	87.9	88.4	94.5	48.9	36.5	52.2	36.8	10.4	20.4
Res34 *	76.8	73.2	85.5	49.3	33.1	48.4	76.8	74.3	86.7	89.3	89.9	95.3	50.9	38.5	53.5	39.1	11.8	22.4

*: one-staged clustering method.

199 5 EXPERIMENTS

200 The experimental evaluation is performed on six benchmark datasets, i.e. CIFAR-10, CIFAR-100
 201 [10], STL-10 [11], ImageNet-10, ImageNet-Dogs [2], and tiny-ImageNet [12]. The first two CIFAR
 202 datasets contain 60,000 images of 32x32 pixels. Following previous works [6, 5], we take 20 super-
 203 classes rather than 100 classes as the ground-truth for CIFAR-100. The next is STL-10 containing
 204 13,000 labeled images and 100,000 unlabeled images of 96x96 pixels. The additional 100,000
 205 unlabeled images are used to perform the instance-level contrastive learning. For the last three
 206 ImageNet datasets, only the training set is used. ImageNet-10 contains 13,000 images of 10 classes,
 207 ImageNet-Dogs consists of 19,500 images from 15 dog classes, and large-scaled tiny-ImageNet
 208 contains 100,000 images from 200 classes.

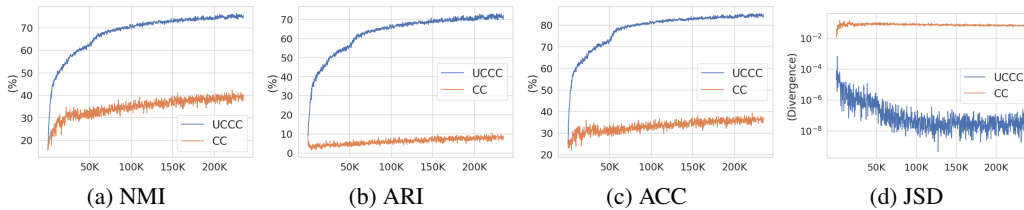
209 We adopt ResNet18/ResNet34 [13] as our backbone net and modify the stride of the first convolution
 210 layer to 1, which enables the encoder to extract more delicate features. Especially for small-sized
 211 datasets (i.e. CIFAR), we remove the first maxpooling layer and replace the activation of the first
 212 layer to Mish [14]. For ICB, the dimensionality of the final embedding space is set to 128, and the
 213 temperature parameter τ is fixed to 0.5 in all experiments. For CCB, the dimensionality of cluster
 214 assignments is naturally set to the number of clusters, and the assignment aggregation is fixed as
 215 2-hop graph fusion in all used datasets. We set the batch size to $N = 256 \times n_class/10$ except for
 216 tiny-ImageNet, where N is set to 512 due to the memory limitation, and the image size is set to 32
 217 for CIFAR datasets, 224 for other datasets. The whole model is trained from scratch for 1,000 epochs
 218 on NVIDIA Tesla V100 GPU.

219 **Evaluation criterion.** We evaluate the results based on three widely used clustering metrics in-
 220 cluding normalized mutual information (NMI), adjusted rand index (ARI) and accuracy (ACC). To
 221 further analyze the severity of cluster degeneracy, we also report Jensen–Shannon divergence (JSD)
 222 in Sec 5.2. Except for JSD, higher values of these metrics indicate better clustering performance.

223 5.1 Comparison with state-of-the-arts

224 We compare our method to the state-of-the-art on six different benchmarks. The compared state-
 225 of-the-art are mostly multi-staged methods, e.g. SC [16], NMF [18], AE [19], DAE [20], DCGAN
 226 [21], DeCNN [22], and VAE [23] obtain clustering results via k-means on the features extracted from
 227 images. According to the results shown in Table 2, UCCC outperforms these competing baselines
 228 by a large margin on all six datasets. Notably, UCCC obtains classification accuracy improvements
 229 compared with the closest competitor CC [5] by 6.5% on CIFAR-10, 5.5% on CIFAR-20, 1.7%

Figure 3: The performance evolution of UCCC and CC (w/o the balance loss) on CIFAR-10.



230 on STL-10, 6.4% on ImageNet-10, 10.6% on ImageNet-Dogs. As CC adopts a dual contrastive
 231 framework like ours, the remarkable results prove the efficacy of our design for cluster-level learning.

232 **5.2 Cluster Degeneracy**

233 Many previous works [15, 5, 6] would lead to degenerate solu-
 234 tions because such solutions are saddle points of the adopted
 235 objective function (Fig. 4). To avoid cluster degeneracy, they
 236 usually manage to assign an equal number of samples to each
 237 cluster by maximizing the entropy of clustering result. How-
 238 ever, this kind of optimization may not be a general solution
 239 especially when it deals with imbalance clusters.

240 In Fig. 3, the performance curves of our approach and CC
 241 (w/o Eq. 9) during the training phase are presented to compare
 242 ours with the degenerate one. As the training process goes, the
 243 performances of UCCC in the four considered metrics are all
 244 improved, implying cluster assignments become more reason-
 245 able. By contrast, CC hardly achieves small improvements due
 246 to the degeneracy problem, which can be reflected by the high
 247 divergence between the distribution of the ground truth and
 248 its clustering. We further perform t-SNE in the instance-level
 249 space on both UCCC and CC (Fig. 1). The result well confirms our model does not suffer from
 250 cluster degeneracy since there is no dominant cluster in our result (Fig. 1b).

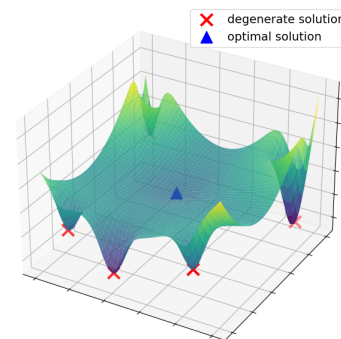


Figure 4: optimization surface

251 **5.3 Ablation Study**

252 **5.3.1 Inductive Clustering**

253 To further validate the effectiveness of UCCC, our method
 254 is evaluated on the three datasets under a more realistic
 255 inductive setting, where testing images are not available
 256 during the training phase. Specifically, we trained the
 257 model on the training set while measuring the performance
 258 using the images in testing split. Table 3 shows our relative
 259 performance drops over transductive setting. As can be
 260 seen, the performance drops are stable even though the
 261 used CNN backbones are different. In particular, the ACC
 262 drops on all three datasets are less than 3%, which well
 263 demonstrates the robustness of UCCC.

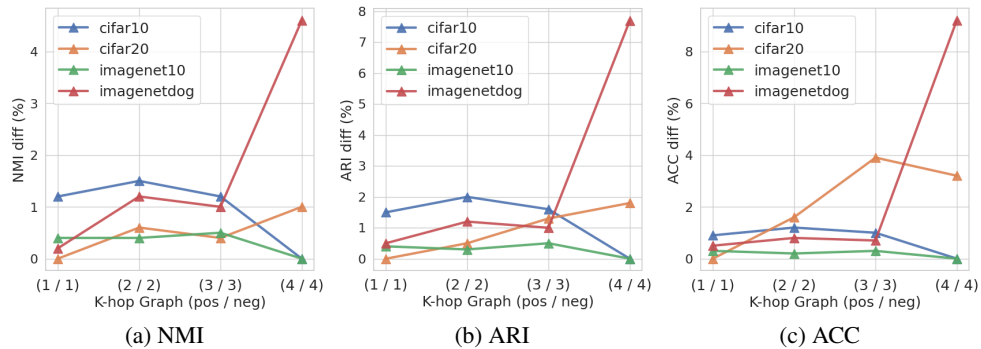
Table 3: The relative performance drops (%) under inductive setting over transductive setting.

Metrics	NMI	ARI	ACC
STL-10 (Res18)	-2.8	-3.7	-2.3
STL-10 (Res34)	-2.9	-3.4	-3.0
CIFAR-10 (Res18)	-2.4	-3.3	-2.0
CIFAR-10 (Res34)	-3.4	-2.8	-1.5
CIFAR-20 (Res18)	-3.8	-2.8	-1.1
CIFAR-20 (Res34)	-4.5	-3.8	-2.1

264 **5.3.2 K-hop Graph Fusion**

265 To study the effect of k -hop graph fusion, we take the performance with $k = 1$ as the baseline and
 266 present the trends of NMI/ARI/ACC gain with respect to varying values of k over $\{1, 2, 3, 4\}$ on four
 267 benchmark datasets. As shown in Fig. 5, we can see that performance improvements are achieved
 268 when considering more than 1-hop neighbors on the four datasets. It is also observed that using 2-hop
 269 graph fusion leads to the best performance in most cases because taking too many community peers
 270 into consideration usually comes with noisy information. Nevertheless, for ImageNet-Dogs dataset,

Figure 5: The influence of k -hop graph fusion on four datasets.



271 the results using 4-hop graph fusion outperform the one with $k = 1$ by a large margin (+3.4% in
 272 NMI, +6.5% in ARI, +8.4% in ACC). The reason might be that ImageNet-Dogs consists of relatively
 273 similar images compared to other datasets, and thus more available information is beneficial to the
 274 model learning.

275 6 Conclusion

276 Based on the observation that similar-looking objects usually belong to the same cluster while objects
 277 that can be easily distinguished tend to come from different clusters, we propose the Unsupervised
 278 Community-consensus Contrastive Clustering (UCCC) which conducts the instance- and cluster-level
 279 contrastive learning simultaneously under a unified framework. By incorporating a novel instance-
 280 based contrastive loss into cluster-level learning, our model is encouraged to produce consensual
 281 cluster assignments. We further verify that our network does not lead to degenerate solutions due to
 282 the designed cluster-level contrastive loss. Experimental evaluation shows that the proposed method
 283 outperforms prior work by large margins for a variety of datasets.

284 Broader Impact

285 For unsupervised classification, there is no need to specify in advance all the classes in the image,
 286 which reduces the dependence of deep learning on massive labeled data. The proposed method
 287 in this work adopts contrastive learning to produce reliable cluster assignments. As a one-stage
 288 deep clustering method, our work is capable of extracting discriminative features and performing
 289 clustering through one-step training. Although such training process could reduce the deployment
 290 difficulty of clustering algorithms in practical applications, providing the opportunity to learn from
 291 big unannotated datasets may cause problems concerning users' information security and personal
 292 privacy if not controlled.

293 Broadly speaking, there are two shortcuts of deep clustering. Firstly, the prediction accuracy is still
 294 lower than models trained in a supervised manner. Thus clustering is not applicable to those settings
 295 that require high accuracy and confidence, e.g., self-driving cars. Secondly, computational complexity
 296 for deep clustering is much higher than the traditional algorithms, e.g. k-means. Thus the deep model
 297 can not be applied to applications where computational resource is limited.

298 **References**

- 299 [1] Junyuan Xie, Ross Girshick, and Ali Farhadi. Unsupervised deep embedding for clustering
300 analysis. In *International conference on machine learning*, pages 478–487. PMLR, 2016.
- 301 [2] Jianlong Chang, Lingfeng Wang, Gaofeng Meng, Shiming Xiang, and Chunhong Pan. Deep
302 adaptive image clustering. In *Proceedings of the IEEE international conference on computer
303 vision*, pages 5879–5887, 2017.
- 304 [3] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering
305 for unsupervised learning of visual features. In *Proceedings of the European Conference on
306 Computer Vision (ECCV)*, pages 132–149, 2018.
- 307 [4] Mathilde Caron, Piotr Bojanowski, Julien Mairal, and Armand Joulin. Unsupervised pre-
308 training of image features on non-curated data. In *Proceedings of the IEEE/CVF International
309 Conference on Computer Vision*, pages 2959–2968, 2019.
- 310 [5] Yunfan Li, Peng Hu, Zitao Liu, Dezhong Peng, Joey Tianyi Zhou, and Xi Peng. Contrastive
311 clustering. *arXiv preprint arXiv:2009.09687*, 2020.
- 312 [6] Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, Marc Proesmans, and Luc
313 Van Gool. Scan: Learning to classify images without labels. In *European Conference on
314 Computer Vision*, pages 268–285. Springer, 2020.
- 315 [7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework
316 for contrastive learning of visual representations. In *International conference on machine
317 learning*, pages 1597–1607. PMLR, 2020.
- 318 [8] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional
319 networks. *arXiv preprint arXiv:1609.02907*, 2016.
- 320 [9] Shun-ichi Amari. Backpropagation and stochastic gradient descent method. *Neurocomputing*, 5
321 (4-5):185–196, 1993.
- 322 [10] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
323 2009.
- 324 [11] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsuper-
325 vised feature learning. In *Proceedings of the fourteenth international conference on artificial
326 intelligence and statistics*, pages 215–223. JMLR Workshop and Conference Proceedings, 2011.
- 327 [12] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7:7, 2015.
- 328 [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
329 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
330 pages 770–778, 2016.
- 331 [14] Diganta Misra. Mish: A self regularized non-monotonic activation function. *arXiv preprint
332 arXiv:1908.08681*, 2019.
- 333 [15] James MacQueen et al. Some methods for classification and analysis of multivariate observations.
334 In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*,
335 volume 1, pages 281–297. Oakland, CA, USA, 1967.
- 336 [16] Lihi Zelnik-Manor and Pietro Perona. Self-tuning spectral clustering, 2004. *Advances in Neural
337 Information Processing Systems*, 17, 2005.
- 338 [17] K Chidananda Gowda and G Krishna. Agglomerative clustering using the concept of mutual
339 nearest neighbourhood. *Pattern recognition*, 10(2):105–112, 1978.
- 340 [18] Deng Cai, Xiaofei He, Xuanhui Wang, Hujun Bao, and Jiawei Han. Locality preserving
341 nonnegative matrix factorization. In *IJCAI*, volume 9, pages 1010–1015, 2009.
- 342 [19] Yoshua Bengio, Pascal Lamblin, Dan Popovici, Hugo Larochelle, et al. Greedy layer-wise
343 training of deep networks. *Advances in neural information processing systems*, 19:153, 2007.

- 344 [20] Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, Pierre-Antoine Manzagol,
345 and Léon Bottou. Stacked denoising autoencoders: Learning useful representations in a deep
346 network with a local denoising criterion. *Journal of machine learning research*, 11(12), 2010.
- 347 [21] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with
348 deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- 349 [22] Matthew D Zeiler, Dilip Krishnan, Graham W Taylor, and Rob Fergus. Deconvolutional net-
350 works. In *2010 IEEE Computer Society Conference on computer vision and pattern recognition*,
351 pages 2528–2535. IEEE, 2010.
- 352 [23] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*
353 *arXiv:1312.6114*, 2013.
- 354 [24] Jianwei Yang, Devi Parikh, and Dhruv Batra. Joint unsupervised learning of deep representations
355 and image clusters. In *Proceedings of the IEEE conference on computer vision and pattern*
356 *recognition*, pages 5147–5156, 2016.
- 357 [25] Philip Haeusser, Johannes Plapp, Vladimir Golkov, Elie Aljalbout, and Daniel Cremers. ASSO-
358 ciative deep clustering: Training a classification network with no labels. In *German Conference*
359 *on Pattern Recognition*, pages 18–32. Springer, 2018.
- 360 [26] Jianlong Chang, Yiwen Guo, Lingfeng Wang, Gaofeng Meng, Shiming Xiang, and Chunhong
361 Pan. Deep discriminative clustering analysis. *arXiv preprint arXiv:1905.01681*, 2019.
- 362 [27] Jianlong Wu, Keyu Long, Fei Wang, Chen Qian, Cheng Li, Zhouchen Lin, and Hongbin Zha.
363 Deep comprehensive correlation mining for image clustering. In *Proceedings of the IEEE/CVF*
364 *International Conference on Computer Vision*, pages 8150–8159, 2019.
- 365 [28] Xu Ji, João F Henriques, and Andrea Vedaldi. Invariant information clustering for unsuper-
366 vised image classification and segmentation. In *Proceedings of the IEEE/CVF International*
367 *Conference on Computer Vision*, pages 9865–9874, 2019.
- 368 [29] Sungwon Han, Sungwon Park, Sungkyu Park, Sundong Kim, and Meeyoung Cha. Mitigating
369 embedding and class assignment mismatch in unsupervised image classification. In *16th*
370 *European Conference on Computer Vision, ECCV 2020*. Springer, 2020.
- 371 [30] Jiabo Huang, Shaogang Gong, and Xiatian Zhu. Deep semantic clustering by partition confi-
372 dence maximisation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
373 *Pattern Recognition*, pages 8849–8858, 2020.

374 Checklist

375 The checklist follows the references. Please read the checklist guidelines carefully for information on
376 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
377 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
378 the appropriate section of your paper or providing a brief inline description. For example:

- 379 • Did you include the license to the code and datasets? **[Yes]** See Section ??.
- 380 • Did you include the license to the code and datasets? **[No]** The code and the data are
381 proprietary.
- 382 • Did you include the license to the code and datasets? **[N/A]**

383 Please do not modify the questions and only use the provided macros for your answers. Note that the
384 Checklist section does not count towards the page limit. In your paper, please delete this instructions
385 block and only keep the Checklist section heading above along with the questions/answers below.

- 386 1. For all authors...
 - 387 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
388 contributions and scope? **[Yes]**
 - 389 (b) Did you describe the limitations of your work? **[Yes]**
 - 390 (c) Did you discuss any potential negative societal impacts of your work? **[Yes]**
 - 391 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
392 them? **[Yes]**
- 393 2. If you are including theoretical results...
 - 394 (a) Did you state the full set of assumptions of all theoretical results? **[Yes]** see Sec. 4 and
395 the Supplementary
 - 396 (b) Did you include complete proofs of all theoretical results? **[Yes]** see the Supplementary
- 397 3. If you ran experiments...
 - 398 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
399 perimental results (either in the supplemental material or as a URL)? **[Yes]** see the
400 Supplementary
 - 401 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
402 were chosen)? **[Yes]** see Sec. 5
 - 403 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
404 ments multiple times)? **[No]**
 - 405 (d) Did you include the total amount of compute and the type of resources used (e.g., type
406 of GPUs, internal cluster, or cloud provider)? **[Yes]** see Sec. 5
- 407 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - 408 (a) If your work uses existing assets, did you cite the creators? **[Yes]** see Sec. 5
 - 409 (b) Did you mention the license of the assets? **[Yes]**
 - 410 (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]**
411 reference code is in the Supplementary
 - 412 (d) Did you discuss whether and how consent was obtained from people whose data you're
413 using/curating? **[No]**
 - 414 (e) Did you discuss whether the data you are using/curating contains personally identifiable
415 information or offensive content? **[No]**
- 416 5. If you used crowdsourcing or conducted research with human subjects...
 - 417 (a) Did you include the full text of instructions given to participants and screenshots, if
418 applicable? **[N/A]**
 - 419 (b) Did you describe any potential participant risks, with links to Institutional Review
420 Board (IRB) approvals, if applicable? **[N/A]**
 - 421 (c) Did you include the estimated hourly wage paid to participants and the total amount
422 spent on participant compensation? **[N/A]**

1 A Positive Confidence Threshold

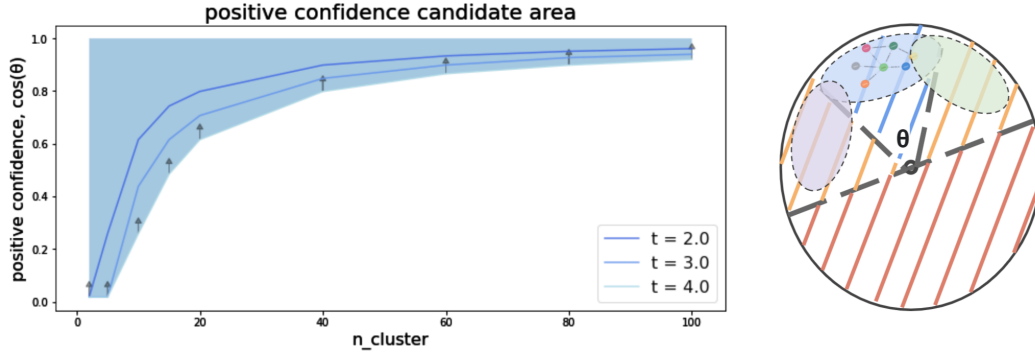


Figure 1: The left shows the positive confidence curve, $\cos \theta$, using different values of tolerant parameter t , and the blue area denotes the similarities that would be taken as positive when $t = 4.0$. The right illustrates the threshold θ , which is calculated based on the number of clusters and tolerant parameter t . The larger t means the higher tolerance for the overlapping (ambiguous) area between two different clusters.

- 2 As we take the cosine similarity (dot product between $L2$ normalized vectors),
 3 we use a threshold angle θ to filter out noisy information from ambiguous instance pairs. That is, if
 4 the similarity between an instance pair is lower than $\cos \theta$, we would regard it as a ambiguous pair.
- 5 To cluster data samples into groups, the feature vectors extracted from samples should be separable
 6 enough on the embedding space. Applying $L2$ normalization is actually projecting a feature vector
 7 into a point on unit sphere. For simplification, we consider the unit 2-sphere as the embedding sphere.
 8 Ideally, the embedding features are uniformly scattered on the surface of the sphere. We assume each
 9 cluster takes over the equal size of surface area and measure the surface area by taking each cluster
 10 as a spherical cap. Then for any given instance, we consider a spherical cap where the pole is its
 11 embedding point on the sphere, and regard the points on the cap as its positive peers. In other words,
 12 the threshold angle θ is calculated as:

$$\begin{aligned} 4\pi r^2 \cdot t &= n_cluster \cdot 2\pi r^2 (1 - \cos \theta) \\ \Rightarrow \cos \theta &= 1 - \frac{2t}{n_cluster} \end{aligned} \quad (1)$$

- 13 Here t is a *tolerant parameter* controlling the tolerance for the intra-cluster overlapping area (Fig. 1
 14 right). The tolerant parameter is set to 4 in all experiments as we empirically find it leads to the best
 15 performance in overall. See Sec. C.2 for the ablation study on the impact of t .

16 B Gradient Derivation

- 17 In this section, we present our derivation for the gradients of the two considered cluster-level
 18 contrastive losses, \mathcal{L}_{clu}^{clu} and \mathcal{L}_{clu} , with respect to a classifier for cluster c , w^c . The notations for
 19 derivations are summarized in Table 1. For convenience, we reprint below the expressions for each.

$$\mathcal{L}_{clu}^{clu} = - \sum_{\substack{c=1, \\ k \in \{a,b\}}}^C \log \frac{\exp(\mathbf{y}_a^c \cdot \mathbf{y}_b^c)}{\sum_{j=1}^C \mathbb{1}_{[i \neq c]} \left[\exp(\mathbf{y}_k^c \cdot \mathbf{y}_a^j) + \exp(\mathbf{y}_k^c \cdot \mathbf{y}_b^j) \right]} \quad (2)$$

$$\mathcal{L}_{clu} = - \log \frac{\sum_{i=1}^N \exp(\mathbf{y}_{i^a, pos} \cdot \mathbf{y}_{i^b, pos} + \mathbf{y}_{i^a, neg} \cdot \mathbf{y}_{i^b, neg})}{\sum_{\forall (i,j,k,m) \in \Phi} \exp(\mathbf{y}_{i^k, pos} \cdot \mathbf{y}_{j^m, neg})}, \quad (3)$$

Table 1: Notations of gradient derivations.

Notations	Dimensionality	Descriptions
N	$\in \mathbb{R}$	Mini-batch size
C	$\in \mathbb{R}$	Number of clusters
H	$\in \mathbb{R}$	Hidden dimension of cluster head g_C
a, b	$\in \mathbb{R}$	The first / second random augmentation
w^c	$\in \mathbb{R}^H$	Weights of the last layer for cluster c
\mathbf{H}_k	$\in \mathbb{R}^{N \times H}$	Hidden feature matrix of a mini-batch using augmentation k
$\mathbf{H}_{k,pos}, \mathbf{H}_{k,neg}$	$\in \mathbb{R}^{N \times H}$	Positive / Negative feature matrix of a mini-batch using augmentation k
\mathbf{h}_{i^k}	$\in \mathbb{R}^H$	i th row of feature matrix \mathbf{h}_k
$\mathbf{h}_{i^k,pos}, \mathbf{h}_{i^k,neg}$	$\in \mathbb{R}^C$	i th row of positive / negative assignment matrix $\mathbf{h}_{k,pos} / \mathbf{h}_{k,neg}$
\mathbf{Y}_k	$\in \mathbb{R}^{N \times C}$	Assignment matrix of a mini-batch using augmentation k
$\mathbf{Y}_{k,pos}, \mathbf{Y}_{k,neg}$	$\in \mathbb{R}^{N \times C}$	Positive / Negative assignment matrix of a mini-batch using augmentation k
\mathbf{y}_k^c	$\in \mathbb{R}^N$	c th column of assignment matrix \mathbf{Y}_k
\mathbf{y}_{i^k}	$\in \mathbb{R}^C$	i th row of assignment matrix \mathbf{Y}_k
$\mathbf{y}_{i^k,pos}, \mathbf{y}_{i^k,neg}$	$\in \mathbb{R}^C$	i th row of positive / negative assignment matrix $\mathbf{Y}_{k,pos} / \mathbf{Y}_{k,neg}$
$y_{i^k,pos}^c, y_{i^k,neg}^c$	$\in \mathbb{R}$	c th element of positive / negative assignment $\mathbf{y}_{i^k,pos} / \mathbf{y}_{i^k,neg}$

20 where

$$\Phi = \{(i, j, k, m) \mid \forall i, j \in \{1, \dots, N\} \text{ and } \forall k, m \in \{a, b\}\}.$$

21 B.1 Cluster-based: Cluster-level Contrastive

22 Recall that contrastive losses is to maximize the similarity between positive pairs and minimize the
 23 similarity of negative ones, For convenience, we consider pairwise relations and divide the gradient
 24 into two parts: (1) positive-pair term and (2) negative-pair term.

25 We start from discussing the stabilization of model generalization, which means the model could
 26 produce similar results no matter whatever kind of augmentation an image has been applied. Both
 27 CC [1], which proposed cluster-based clustering loss, and our model leverage a dual contrastive
 28 framework. Such dual contrastive framework incorporates instance-level contrastive learning, so high
 29 similarities between the instances augmented from the same image are encouraged. As a result, we
 30 have the assumption of generalization below.

31 **Assumption 1 (Generalization)** *Let a, b be any two random augmentations. The statements below*
 32 *hold:*

$$\begin{cases} \mathbf{y}_a^j \approx \mathbf{y}_b^j \approx \hat{\mathbf{y}}^j & \text{for } j \in \{1, \dots, C\} \\ \mathbf{H}_a \approx \mathbf{H}_b \approx \hat{\mathbf{H}} \end{cases}$$

33

34 There is another assumption that helps our derivation. We assume the cluster head would produce
 35 confident results, which means there exists a high prediction score with respect to some specific cluster
 36 c' . To verify this assumption, we study on the confidence of our cluster assignment in Sec. C.1.

37 **Assumption 2 (Confident Cluster Assignment)** *There exists some $c' \in \{1, \dots, C\}$ such that*

$$y_i^{c'} \approx 1 \quad \text{and} \quad y_i^j \approx 0, \forall j \neq c'$$

38

39 This implies that

$$\mathbf{y}_a^c \cdot \mathbf{y}_a^j \approx \mathbf{y}_a^c \cdot \mathbf{y}_b^j \approx 0 \quad \text{for } j \neq c, j \in \{1, \dots, C\} \quad (4)$$

40 With the approximations above, the gradients of cluster-level clustering loss in terms of posi-
 41 tive/negative pairs can be derived as:

$$\begin{aligned}
 \left. \frac{\partial \mathcal{L}_{clu}^{clu}}{\partial \mathbf{w}^c} \right|_{pos} &= -2 \frac{\partial}{\partial \mathbf{w}^c} \log [\exp(\mathbf{y}_a^c \cdot \mathbf{y}_b^c)] \\
 &\approx -2 \frac{\partial}{\partial \mathbf{w}^c} (\hat{\mathbf{y}}^c \cdot \hat{\mathbf{y}}^c) \quad (\text{with Asm. 1}) \\
 &= -4 \hat{\mathbf{H}}^\top \hat{\mathbf{y}}^c
 \end{aligned} \tag{5}$$

$$\begin{aligned}
 \left. \frac{\partial \mathcal{L}_{clu}^{clu}}{\partial \mathbf{w}^c} \right|_{neg} &= 2 \frac{\partial}{\partial \mathbf{w}^c} \log \left\{ \sum_{\substack{j=1 \\ k \in \{a,b\}}}^C \mathbb{1}_{[j \neq c]} \left[\exp(\mathbf{y}_k^c \cdot \mathbf{y}_a^j) + \exp(\mathbf{y}_k^c \cdot \mathbf{y}_b^j) \right] \right\} \\
 &= 2 \frac{\partial}{\partial \mathbf{w}^c} \log \left\{ \sum_{j=1}^C \mathbb{1}_{[i \neq c]} \left[\exp(\mathbf{y}_a^c \cdot \mathbf{y}_a^j) + \exp(\mathbf{y}_b^c \cdot \mathbf{y}_b^j) + 2 \exp(\mathbf{y}_a^c \cdot \mathbf{y}_b^j) \right] \right\} \\
 &\approx 2 \frac{\partial}{\partial \mathbf{w}^c} \log \left\{ 4 \sum_{j=1}^C \mathbb{1}_{[j \neq c]} \exp(\hat{\mathbf{y}}^c \cdot \hat{\mathbf{y}}^j) \right\} \quad (\text{with Asm. 1}) \\
 &= \frac{8 \sum_{j=1}^C \left[\mathbb{1}_{[j \neq c]} \hat{\mathbf{H}}^\top \hat{\mathbf{y}}^j \right]}{4 \sum_{j=1}^C \mathbb{1}_{[j \neq c]} \exp(\hat{\mathbf{y}}^c \cdot \hat{\mathbf{y}}^j)} \\
 &\approx \frac{2}{C-1} \sum_{j=1}^C \left[\mathbb{1}_{[j \neq c]} \hat{\mathbf{H}}^\top \hat{\mathbf{y}}^j \right] \quad (\text{with Eq. 4})
 \end{aligned} \tag{6}$$

42 Sum over Eq. 5 and Eq. 6, the gradient of cluster-based contrastive clustering loss is written as:

$$\begin{aligned}
 \frac{\partial \mathcal{L}_{clu}^{clu}}{\partial \mathbf{w}^c} &= \left. \frac{\partial \mathcal{L}_{clu}^{clu}}{\partial \mathbf{w}^c} \right|_{pos} + \left. \frac{\partial \mathcal{L}_{clu}^{clu}}{\partial \mathbf{w}^c} \right|_{neg} \\
 &= -4 \hat{\mathbf{H}}^\top \hat{\mathbf{y}}^c + \frac{2}{C-1} \sum_{j=1}^C \left[\mathbb{1}_{[j \neq c]} \hat{\mathbf{H}}^\top \hat{\mathbf{y}}^j \right] \\
 &= -2 \hat{\mathbf{H}}^\top \left[2 \hat{\mathbf{y}}^c - \frac{2}{C-1} \sum_{j=1}^C \mathbb{1}_{[j \neq c]} \hat{\mathbf{y}}^j \right]
 \end{aligned} \tag{7}$$

43 **Discussion.** The classifiers \mathbf{w}^c is updated through $\mathbf{w}^c = \mathbf{w}^c - \eta \frac{\partial \mathcal{L}_{clu}^{clu}}{\partial \mathbf{w}^c}$ every step. Such optimiza-
 44 tion would pull each classifier \mathbf{w}^c closer to the samples that are assigned to cluster c and push it away
 45 from the other samples. However, if there exists a cluster c' such that most samples in a mini-batch
 46 are assigned to cluster c' , we can make the following observations. First, the classifier $\mathbf{w}^{c'}$ would
 47 get closer to most samples in a mini-batch. Second, for the other classifiers $\mathbf{w}^c, c \neq c'$, the strength
 48 of pushing classifiers away from most samples in a mini-batch (which are assigned to cluster c') is
 49 larger than the one that pulls them closer to the samples belonging to their corresponding cluster c .
 50 The observations above imply that the dominant cluster c' is becoming larger while the growth of the
 51 other clusters are being suppressed. Hence, the optimization through the cluster-based contrastive
 52 loss could not alleviate cluster degeneracy.

53 B.2 Instance-based: Cluster-level Contrastive Loss

54 Notice that we construct the communities in the instance-level subspace, where instance distinction is
 55 preserved with the guidance of instance-level contrastive loss. As we have discussed in appx. B.1, the
 56 assumption of model generalization holds, thus implying the statements in the following also hold:

$$\begin{cases} \mathbf{y}_{i^a,*} \approx \mathbf{y}_{i^b,*} \approx \hat{\mathbf{y}}_{i,*} \\ \mathbf{H}_{a,*} \approx \mathbf{H}_{b,*} \approx \hat{\mathbf{H}}_* \end{cases} \text{ for } i \in \{1, \dots, N\}, * \in \{pos, neg\} \quad (8)$$

57 With Asm. 2, we can further have

$$\mathbf{y}_{i^a,pos} \cdot \mathbf{y}_{i^b,pos} \approx \mathbf{y}_{i^a,neg} \cdot \mathbf{y}_{i^b,neg} \approx 1 \quad (9)$$

58 Hence, the gradients with regards to positive/negative pairs are:

$$\begin{aligned} \left. \frac{\partial \mathcal{L}_{clu}}{\partial \mathbf{w}^c} \right|_{pos} &= -\frac{\partial}{\partial \mathbf{w}^c} \log \left\{ \sum_{i=1}^N \exp(\mathbf{y}_{i^a,pos} \cdot \mathbf{y}_{i^b,pos} + \mathbf{y}_{i^a,neg} \cdot \mathbf{y}_{i^b,neg}) \right\} \\ &\approx -\frac{\partial}{\partial \mathbf{w}^c} \log \left\{ \sum_{i=1}^N \exp(\hat{\mathbf{y}}_{i,pos} \cdot \hat{\mathbf{y}}_{i,pos} + \hat{\mathbf{y}}_{i^a,neg} \cdot \hat{\mathbf{y}}_{i^b,neg}) \right\} \quad (\text{with Eq. 8}) \\ &\approx \frac{-2 \exp 2 \left[\hat{\mathbf{H}}_{pos}^\top \hat{\mathbf{y}}_{pos}^c + \hat{\mathbf{H}}_{neg}^\top \hat{\mathbf{y}}_{neg}^c \right]}{N \exp 2} \quad (\text{with Eq. 9}) \\ &= -\frac{2}{N} \left[\hat{\mathbf{H}}_{pos}^\top \hat{\mathbf{y}}_{pos}^c + \hat{\mathbf{H}}_{neg}^\top \hat{\mathbf{y}}_{neg}^c \right] \end{aligned} \quad (10)$$

$$\begin{aligned} \left. \frac{\partial \mathcal{L}_{clu}}{\partial \mathbf{w}^c} \right|_{neg} &= \frac{\partial}{\partial \mathbf{w}^c} \log \left\{ \sum_{\forall (i,j,k,m) \in \Phi} \exp(\mathbf{y}_{i^k,pos} \cdot \mathbf{y}_{j^m,neg}) \right\} \\ &\approx \frac{\partial}{\partial \mathbf{w}^c} \log \left\{ 4 \sum_{i=1}^N \sum_{j=1}^N \exp(\hat{\mathbf{y}}_{i,pos} \cdot \hat{\mathbf{y}}_{j,neg}) \right\} \quad (\text{with Eq. 8}) \\ &= \frac{1}{\Psi} \left[\left(\hat{\mathbf{E}} \hat{\mathbf{H}}_{pos} \right)^\top \hat{\mathbf{y}}_{neg}^c + \left(\hat{\mathbf{E}} \hat{\mathbf{H}}_{neg} \right)^\top \hat{\mathbf{y}}_{pos}^c \right] \end{aligned} \quad (11)$$

59 where

$$\Psi \equiv \sum_{i=1}^N \sum_{j=1}^N \exp(\hat{\mathbf{y}}_{i,pos} \cdot \hat{\mathbf{y}}_{j,neg}) \quad (12)$$

$$\hat{\mathbf{E}}_{ij} \equiv \exp(\hat{\mathbf{y}}_{i,pos} \cdot \hat{\mathbf{y}}_{j,neg}), \hat{\mathbf{E}} \in \mathbb{R}^{N \times N} \quad (13)$$

60 Therefore, the total gradient of the proposed instance-based contrastive loss is derived as:

$$\begin{aligned} \frac{\partial \mathcal{L}_{clu}}{\partial \mathbf{w}^c} &= \left. \frac{\partial \mathcal{L}_{clu}}{\partial \mathbf{w}^c} \right|_{pos} + \left. \frac{\partial \mathcal{L}_{clu}}{\partial \mathbf{w}^c} \right|_{neg} \\ &= -\frac{2}{N} \left[\hat{\mathbf{H}}_{pos}^\top \hat{\mathbf{y}}_{pos}^c + \hat{\mathbf{H}}_{neg}^\top \hat{\mathbf{y}}_{neg}^c \right] + \frac{1}{\Psi} \left[\left(\hat{\mathbf{E}} \hat{\mathbf{H}}_{pos} \right)^\top \hat{\mathbf{y}}_{neg}^c + \left(\hat{\mathbf{E}} \hat{\mathbf{H}}_{neg} \right)^\top \hat{\mathbf{y}}_{pos}^c \right] \end{aligned} \quad (14)$$

61 C More Ablation Studies

62 C.1 Confidence Learning Curve

63 We plot the confidence learning curve by averaging the maximum assignment score of each sample
 64 within a mini-batch. As shown Fig. 2, the confidence score are quite high (approximate to 1) after a
 65 few training epochs, which verifies the assumption we made for the gradient derivation.

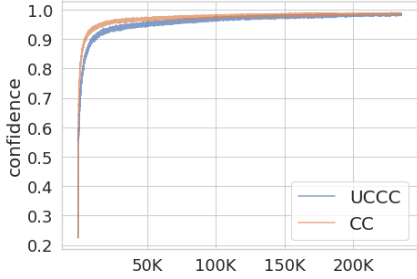


Figure 2: The learning curve of assignment confidence on CIFAR-10.

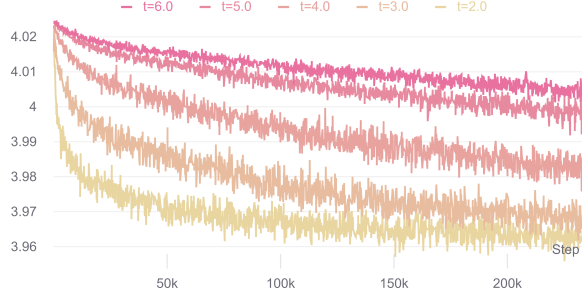


Figure 3: The training curves of our proposed cluster-level contrastive loss using different values of the tolerant parameter t on CIFAR-10.

66 C.2 Tolerant Parameter

67 To observe the influence of the tolerant parameter t for positive confidence
 68 threshold, we conduct experiments varying t from $\{2,3,4,5,6\}$.
 69 Table 2 shows the proposed model achieves the best performance
 70 when $t = 4$. Additionally, if the tolerant parameter is too large, the
 71 performances on three considered metrics might drop significantly.
 72 As the large t implies the low threshold value to construct positive
 73 peer relations, information from the noisy peers may cause the model
 74 hard to distinguish samples from different clusters.

Table 2: The influence of different tolerant parameter t on CIFAR-10 dataset.

Metrics	NMI	ARI	ACC
$t = 2$	70.8	64.6	79.0
$t = 3$	73.1	70.1	83.0
$t = 4$	74.8	71.0	84.2
$t = 5$	74.0	70.4	83.8
$t = 6$	69.9	64.4	77.8

75 According to the loss curves (Fig. 3), we also find that larger values of t usually require more training
 76 steps for model convergence.

77 C.3 Qualitative Analysis

78 We conducted a qualitative analysis to examine how well the clustering result is on CIFAR-10,
 79 ImageNet-10, and large-scale tiny-ImageNet dataset.

80 **CIFAR-10 Dataset.** The confusion matrix between the ground truth labels and classification results
 81 is shown in Fig. 4. A perfect classification would only place items on the diagonal line. The figure
 82 shows that the proposed model finds the right cluster for most images except for some error-prone
 83 classes such as *birds*, *cats*, and *dogs*.

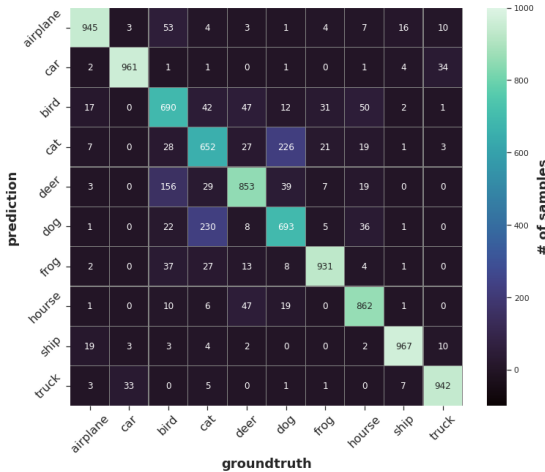
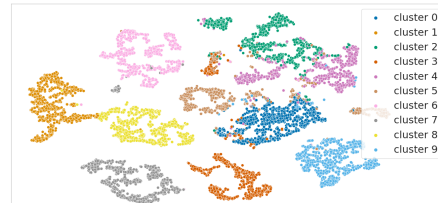


Figure 4: The confusion matrix on CIFAR-10.



(a) The clustering result of UCCC.



(b) Ground-truth labels.

Figure 5: t-SNE visualization on CIFAR-10.

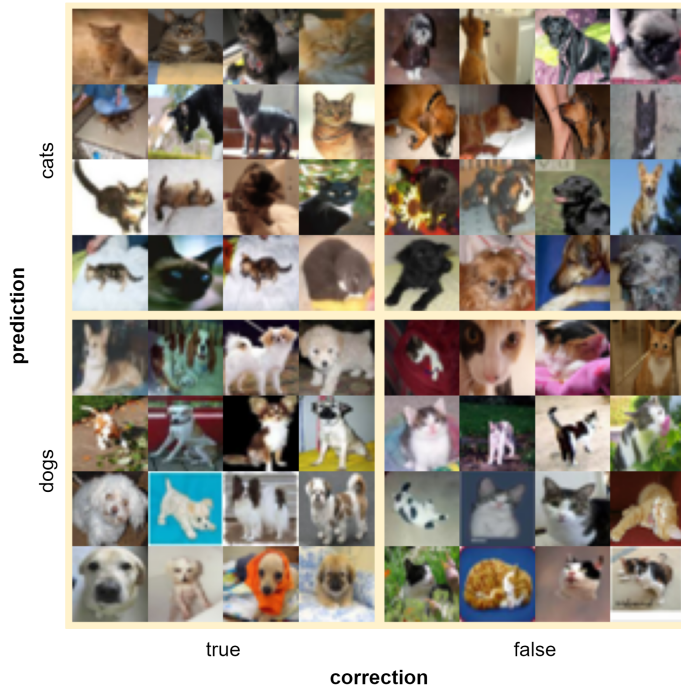


Figure 6: Example images from the cluster *dogs* and *cats* on CIFAR-10. The two left blocks contain the correctly classified images and those misclassified images are demonstrated in the right blocks.

84 We further perform t-SNE visualization in the instance-level embedding space. Fig. 5 compares the
 85 visualization results of predicted clusters and the ground-truth labels, where different colors indicated
 86 different labels/clusters. As shown in Fig. 5b, it is observed that the extracted features from class
 87 *dogs* and class *cat* (denoted in navy and light pink respectively) are hardly distinguishable. This
 88 explains why many dog images are misclassified as cats, and vice versa (Fig. 6). Nonetheless, the
 89 result in Fig. 5a still proves the efficacy of the instance-level contrastive learning since the features
 90 from different clusters are mostly well separated.

91 **ImageNet-10 Dataset.** The confusion matrix and the t-SNE visualization are demonstrated in the
 92 figures below. In Fig. 7, a high concentration of items on the diagonal line confirms the proposed
 93 model correctly groups all samples into 10 classes. Fig. 8 also verifies our clustering result is almost
 94 the same as the ground-truth labels.

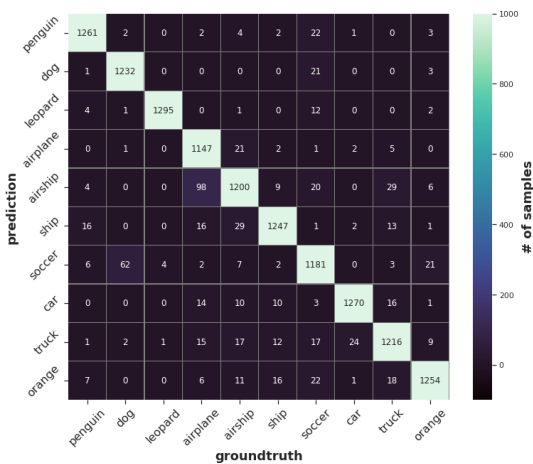
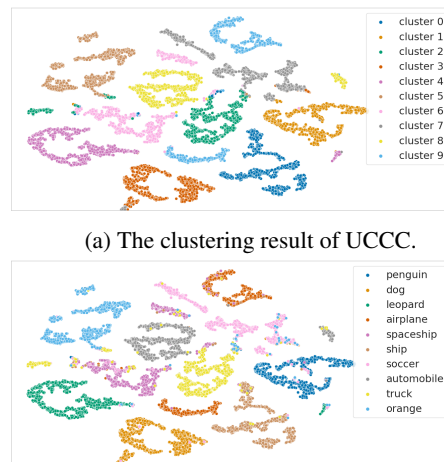


Figure 7: The confusion matrix on ImageNet-10.



(a) The clustering result of UCCC.

(b) Ground-truth labels.

Figure 8: t-SNE visualization on ImageNet-10.

Superclass	Classes				
Animals	goldfish	salamandra	bullfrog	bell toad	American alligator
	fox constrictor	trilobite	scorpion	black widow	tarantula
	centipede	goose	koala	jellyfish	bruin
	brain coral	snail	slug	sea slug	Maine lobster
	crawfish	black stork	king penguin	albatross	dugong
	Chihuahua	Yorkshire terrier	golden retriever	alsatian	standard poodle
tabby	Persian cat	Egyptian cat	puma	lion	
Insects	ladybug	fly	bee	hopper	walkingstick
	roach	mantis	dragonfly	monarch	sulfur butterfly
	holothurian	guinea pig	dog	ox	bison
	bighorn	gazelle	dromedary	orang	chimp
	baboon	African elephant	panda	abacus	judge's robe
	altar	apron	backpack	banister	barbershop
Others	barn	cask	basketball	tub	wagon
	beacon	beaker	beer bottle	bikini	binoculars
	birdhouse	bowtie	brass	broom	pail
	bullet	meat market	taper	cannon	cardigan
	ATM	CD player	chain	chest	Christmas stocking
	cliff dwelling	keypad	candy store	convertible	crane
	dam	desk	board	drumstick	dumbbell
	flagpole	fountain	freight car	frypan	fur coat
	gamask	go-kart	gondola	hourglass	iPod
	ricksha	kimono	lampshade	mower	lifeboat
	limo	magnetic compass	maypole	military uniform	mini
	moving van	nail	neck brace	obelisk	oboe
	organ	parking meter	pay-phone	paling	pill bottle
	plunger	pole	wagon	poncho	pop bottle
	potter's wheel	missile	punch bag	reel	icebox
	remote	rocket	rugby ball	sandal	school bus
	scoreboard	sewing machine	storckel	sock	sombrero
	space heater	spider web	sport car	steel arch bridge	stopwatch
	shades	suspension bridge	bathing trunks	syringe	teapot
	teddy	thatch	torch	tractor	triumphal arch
	trolleybus	turnstile	umbrella	vestment	viaduct
	volleyball	water jug	water tower	wok	wooden spoon
	drop	coral reef	lakeside	coast	acorn
	comic book	plate	alp		
Food	guacamole	icecream	lolly	pretzel	mashed potato
	cauliflower	bell pepper	mushroom	orange	lemon
	banana	poegranate	meatloaf	pizza	potpie

Table 3: Superclass definition for tiny-ImageNet.

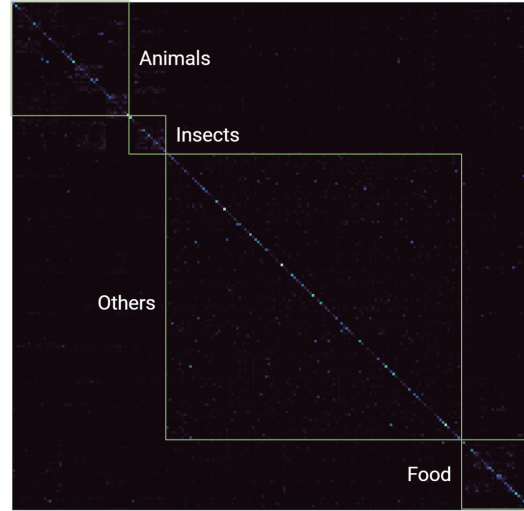


Figure 9: The confusion matrix on tiny-ImageNet.

95 **Tiny-ImageNet Dataset - 200 classes.** In Fig. 9, we mark the superclasses defined in Table 3. The
 96 results show that the misclassified examples tend to be assigned to other clusters from within the
 97 same superclass. Additionally, we demonstrate images from the testing set that were assigned to the
 98 same cluster in Fig. 10, 11, and 12. In particular, Fig. 12 shows some failure cases where different
 99 objects are grouped together due to similar image background.

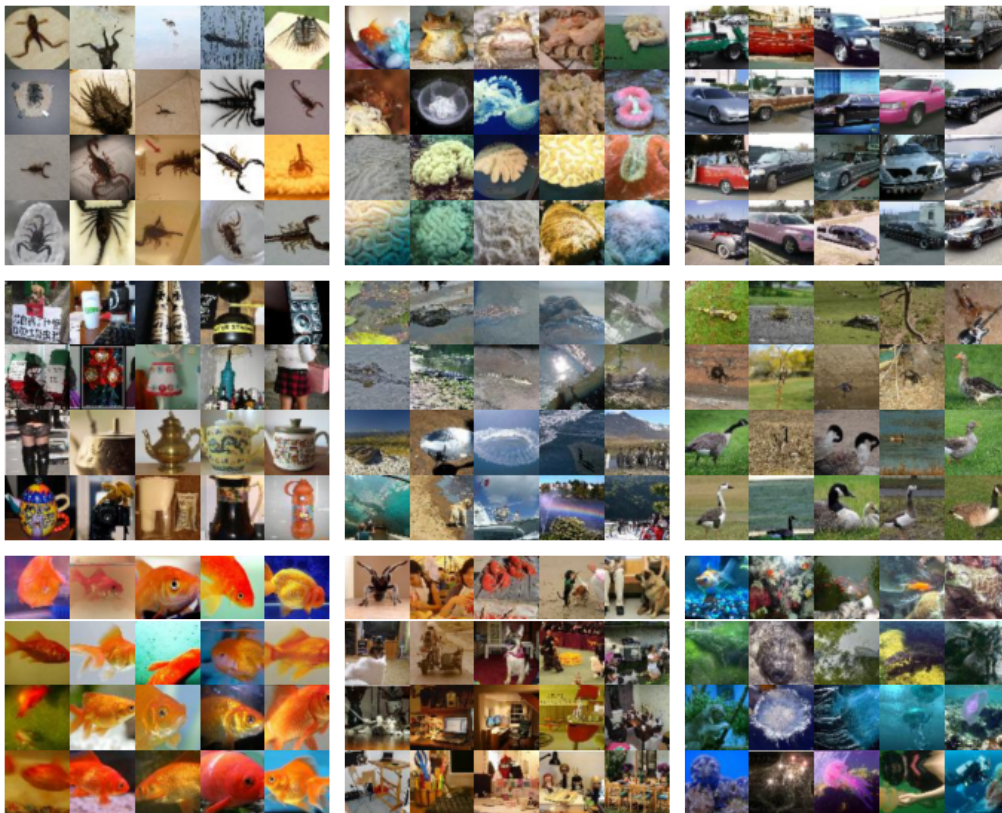


Figure 10: Example clusters of tiny-ImageNet (1).

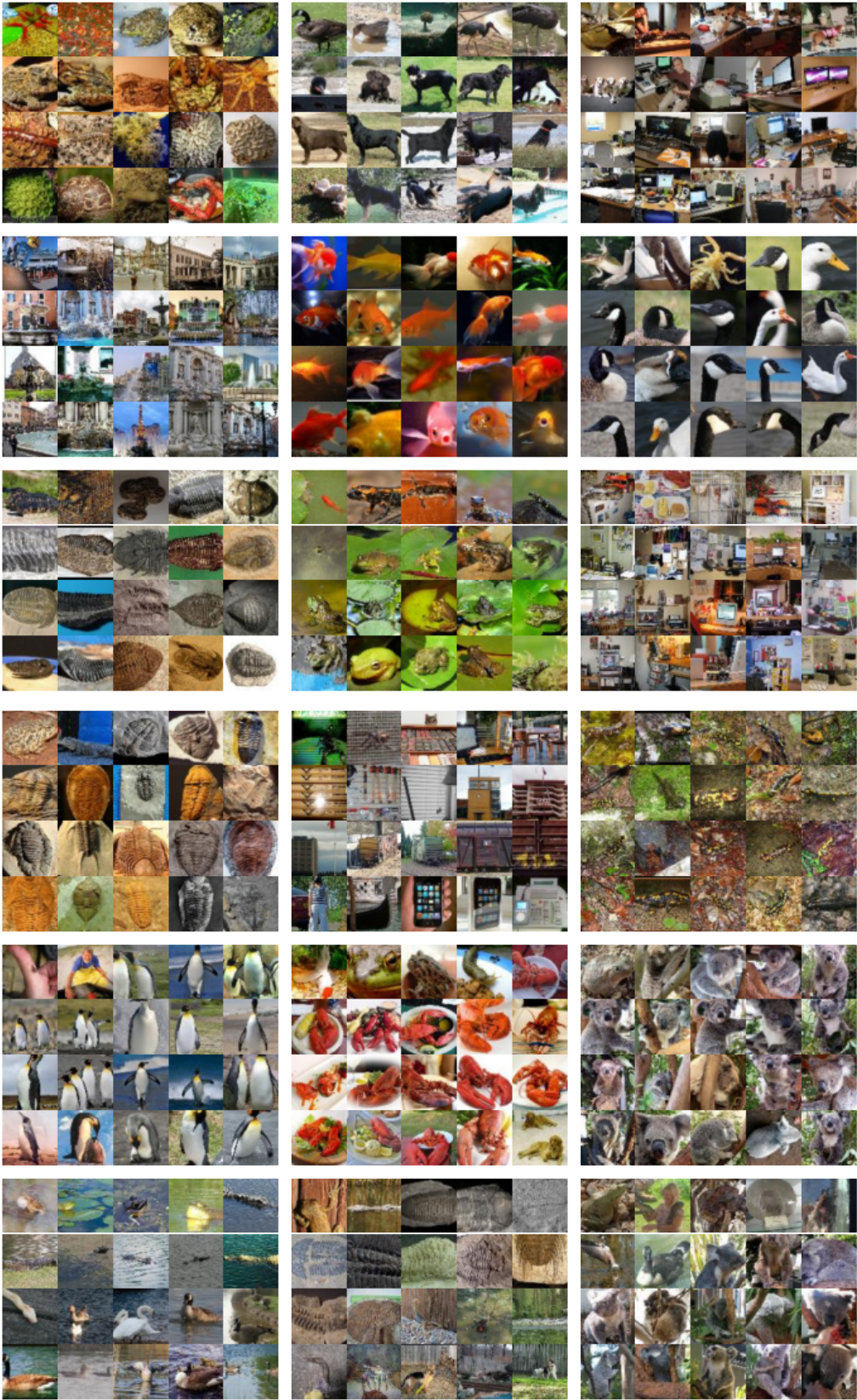


Figure 11: Example clusters of tiny-ImageNet (2).

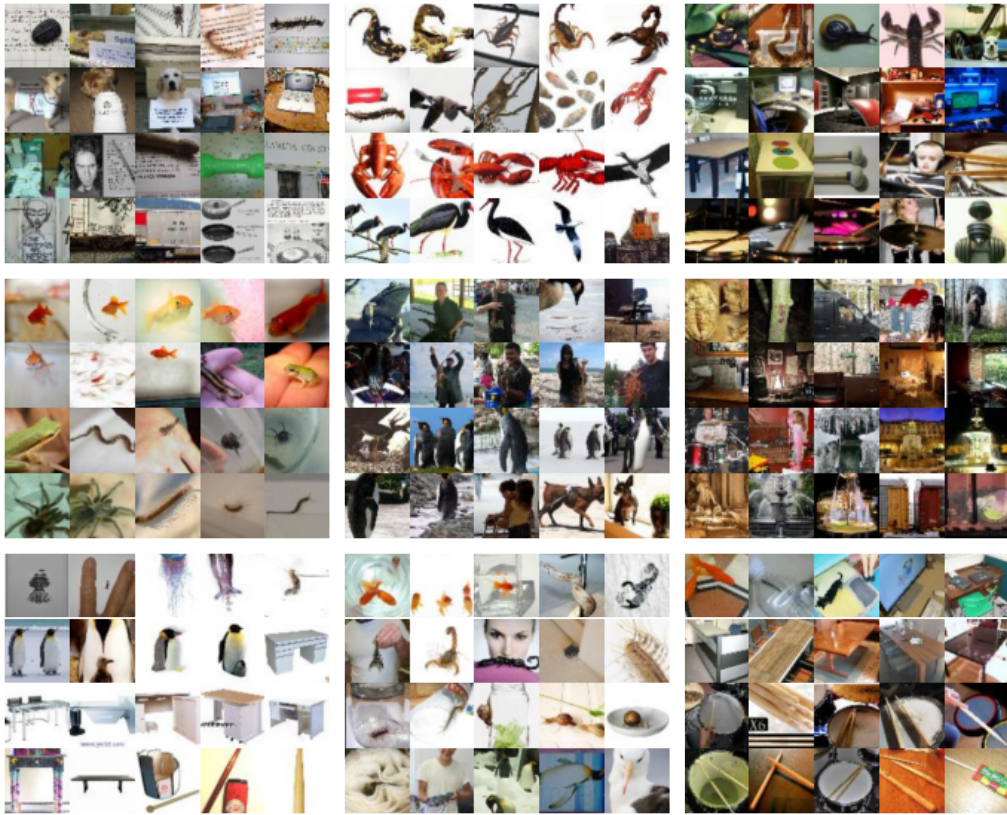


Figure 12: Incorrect clusters of tiny-ImageNet predicted by our model.

100 **References**

- 101 [1] Yunfan Li, Peng Hu, Zitao Liu, Dezhong Peng, Joey Tianyi Zhou, and Xi Peng. Contrastive
102 clustering. *arXiv preprint arXiv:2009.09687*, 2020.