

Joint Unsupervised Learning of Deep Representations and Image Clusters

Jianwei Yang, Devi Parikh, Dhruv Batra
Virginia Tech

{jw2yang, parikh, dbatra}@vt.edu

Abstract

In this paper, we propose a recurrent framework for **Joint Unsupervised LEarning (JULE)** of deep representations and image clusters. In our framework, successive operations in a clustering algorithm are expressed as steps in a recurrent process, stacked on top of representations output by a Convolutional Neural Network (CNN). During training, image clusters and representations are updated jointly: image clustering is conducted in the forward pass, while representation learning in the backward pass. Our key idea behind this framework is that good representations are beneficial to image clustering and clustering results provide supervisory signals to representation learning. By integrating two processes into a single model with a unified weighted triplet loss and optimizing it end-to-end, we can obtain not only more powerful representations, but also more precise image clusters. Extensive experiments show that our method outperforms the state-of-the-art on image clustering across a variety of image datasets. Moreover, the learned representations generalize well when transferred to other tasks. The source code can be downloaded from <https://github.com/jwyang/joint-unsupervised-learning>.

1. Introduction

We are witnessing an explosion in visual content. Significant recent advances in machine learning and computer vision, especially via deep neural networks, have relied on supervised learning and availability of copious annotated data. However, manually labelling data is a time-consuming, laborious, and often expensive process. In order to make better use of available unlabeled images, clustering and/or unsupervised learning is a promising direction.

In this work, we aim to address image clustering and representation learning on unlabeled images in a unified framework. It is a natural idea to leverage cluster ids of images as supervisory signals to learn representations and in turn the representations would be beneficial to image clustering. At a high-level view, given a collection of n_s unlabeled images

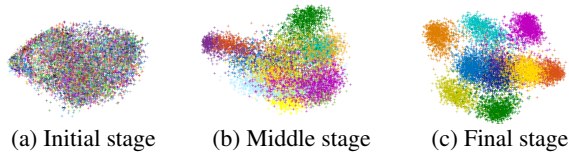


Figure 1: Clustering outputs for MNIST [34] test set at different stages of the proposed method. We conduct PCA on the image representations and then choose the first three dimensions for visualization. Different colors correspond to different clusters. Samples are grouped together gradually and more discriminative representations are obtained.

$\mathbf{I} = \{I_1, \dots, I_{n_s}\}$, the global objective function for learning image representations and clusters can be written as:

$$\operatorname{argmin}_{\mathbf{y}, \theta} \mathcal{L}(\mathbf{y}, \theta | \mathbf{I}) \quad (1)$$

where $\mathcal{L}(\cdot)$ is a loss function, \mathbf{y} denotes the cluster ids for all images, and θ denotes the parameters for representations. If we hold one in $\{\mathbf{y}, \theta\}$ to be fixed, the optimization can be decomposed into two alternating steps:

$$\operatorname{argmin}_{\mathbf{y}} \mathcal{L}(\mathbf{y} | \mathbf{I}, \theta) \quad (2a)$$

$$\operatorname{argmin}_{\theta} \mathcal{L}(\theta | \mathbf{I}, \mathbf{y}) \quad (2b)$$

Intuitively, (2a) can be cast as a conventional clustering problem based on fixed representations, while (2b) is a standard supervised representation learning process.

In this paper, we propose an approach that alternates between the two steps – updating the cluster ids given the current representation parameters and updating the representation parameters given the current clustering result. Specifically, we cluster images using agglomerative clustering[17] and represent images via activations of a Convolutional Neural Network (CNN).

The reason to choose agglomerative clustering is three-fold: 1) it begins with an over-clustering, which is more reliable in the beginning when a good representation has not yet been learned. Intuitively, clustering with representations from a CNN initialized with random weights are not

reliable, but nearest neighbors and over-clusterings are often acceptable; 2) These over-clusterings can be merged as better representations are learned; 3) Agglomerative clustering is a recurrent process and can naturally be interpreted in a recurrent framework.

Our final algorithm is fairly intuitive. We start with an initial over-clustering, update CNN parameters (2b) using image cluster labels as supervisory signals, then merge clusters (2a) and iterate until we reach a stopping criterion. An outcome of the proposed framework is illustrated in Fig. 1. Initially, there are 1,762 clusters for MNIST test set (10k samples), and the representations (image intensities) are not that discriminative. After several iterations, we obtain 17 clusters and more discriminative representations. Finally, we obtain 10 clusters which are well-separated by the learned representations and interestingly correspond primarily to the groundtruth category labels in the dataset, even though the representation is learnt in an unsupervised manner. To summarize, the major contributions of our work are:

- 1 We propose a simple but effective end-to-end learning framework to jointly learn deep representations and image clusters from an unlabeled image set;
- 2 We formulate the joint learning in a recurrent framework, where merging operations of agglomerative clustering are expressed as a forward pass, and representation learning of CNN as a backward pass;
- 3 We derive *a single loss function* to guide agglomerative clustering and deep representation learning, which makes optimization over the two tasks seamless;
- 4 Our experimental results show that the proposed framework outperforms previous methods on image clustering and learns deep representations that can be transferred to other tasks and datasets.

2. Related Work

Clustering Clustering algorithms can be broadly categorized into hierarchical and partitional approaches [25]. Agglomerative clustering is a hierarchical clustering algorithm that begins with many small clusters, and then merges clusters gradually [17, 31, 13]. As for partitional clustering methods, the most well-known is K-means [39], which minimizes the sum of square errors between data points and their nearest cluster centers. Related ideas form the basis of a number of methods, such as expectation maximization (EM) [8, 40], spectral clustering [43, 67, 52], and non-negative matrix factorization (NMF) based clustering [9, 1, 66].

Deep Representation Learning Many works use raw image intensity or hand-crafted features [55, 10, 20, 19, 46, 24] combined with conventional clustering methods. Recently, representations learned using deep neural networks have

presented significant improvements over hand-designed features on many computer vision tasks, such as image classification [30, 51, 54, 49], object detection [15, 14, 21, 47], etc. However, these approaches rely on supervised learning with large amounts of labeled data to learn rich representations. A number of works have focused on learning representations from unlabeled image data. One class of approaches cater to reconstruction tasks, such as auto-encoders [45, 22, 58, 29, 35], deep belief networks (DBN) [33], etc. Another group of techniques learn discriminative representations after fabricating supervisory signals for images, and then finetune them supervisedly for downstream applications [12, 11, 60]. Unlike our approach, the fabricated supervisory signal in these previous works is not updated during representation learning.

Combination A number of works have explored combining image clustering with representation learning. In [56], the authors proposed to learn a non-linear embedding of the undirected affinity graph using stacked autoencoder, and then ran K-means in the embedding space to obtain clusters. In [57], a deep semi-NMF model was used to factorize the input into multiple stacking factors which are initialized and updated layer by layer. Using the representations on the top layer, K-means was implemented to get the final results. Unlike our work, they do not jointly optimize for the representation learning and clustering.

To connect image clustering and representation learning more closely, [64] conducted image clustering and codebook learning iteratively. However, they learned codebook over SIFT feature [37], and did not learn deep representations. Instead of using hand-crafted features, Chen [2] used DBN to learn representations, and then conducted a non-parametric maximum margin clustering upon the outputs of DBN. Afterwards, they fine-tuned the top layer of DBN based on clustering results. A more recent work on jointly optimizing two tasks is found in [61], where the authors trained a task-specific deep architecture for clustering. The deep architecture is composed of sparse coding modules which can be jointly trained through back propagation from a cluster-oriented loss. However, they used sparse coding to extract representations for images, while we use a CNN. Instead of fixing the number of clusters to be the number of categories and predicted labels based on softmax outputs, we predict the labels using agglomerative clustering based on the learned representations. In our experiments we show that our approach outperforms [61].

3. Approach

3.1. Notation

We denote an image set with n_s images by $I = \{I_1, \dots, I_{n_s}\}$. The cluster labels for this image set are $\mathbf{y} = \{y_1, \dots, y_{n_s}\}$. θ are the CNN parameters, based on which we obtain deep representations $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_{n_s}\}$

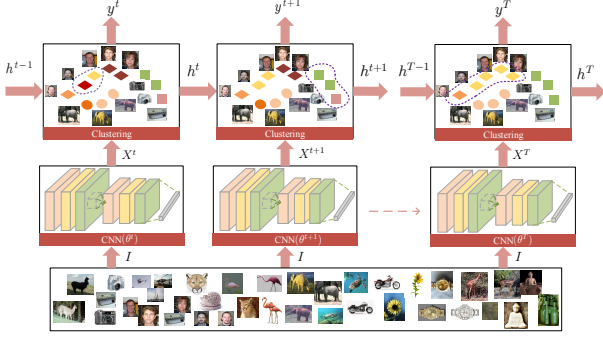


Figure 2: Proposed recurrent framework for unsupervised learning of deep representations and image clusters.

from I . Given the predicted image cluster labels, we organize them into n_c clusters $\mathcal{C} = \{\mathcal{C}_1, \dots, \mathcal{C}_{n_c}\}$, where $\mathcal{C}_i = \{\mathbf{x}_k | y_k = i, \forall k \in 1, \dots, n_s\}$. $\mathcal{N}_i^{K_s}$ are the K_s nearest neighbours of \mathbf{x}_i , and $\mathcal{N}_{\mathcal{C}_i}^{K_c}$ is the set of K_c nearest neighbour clusters of \mathcal{C}_i . For convenience, we sort clusters in $\mathcal{N}_{\mathcal{C}_i}^{K_c}$ in descending order of affinity with \mathcal{C}_i so that the nearest neighbour $\operatorname{argmax}_{\mathcal{C} \in \mathcal{C}^t} \mathcal{A}(\mathcal{C}_i, \mathcal{C})$ is the first entry $\mathcal{N}_{\mathcal{C}_i}^{K_c}[1]$. Here, \mathcal{A} is a function to measure the affinity (or similarity) between two clusters. We add a superscript t to $\{\theta, \mathbf{X}, \mathbf{y}, \mathcal{C}\}$ to refer to their states at timestep t . We use \mathcal{Y} to denote the sequence $\{\mathbf{y}^1, \dots, \mathbf{y}^T\}$ with T timesteps.

3.2. Agglomerative Clustering

As background, we first briefly describe conventional agglomerative clustering [17, 31]. The core idea in agglomerative clustering is to merge two clusters at each step until some stopping conditions. Mathematically, it tries to find two clusters \mathcal{C}_a and \mathcal{C}_b by

$$\{\mathcal{C}_a, \mathcal{C}_b\} = \operatorname{argmax}_{\mathcal{C}_i, \mathcal{C}_j \in \mathcal{C}, i \neq j} \mathcal{A}(\mathcal{C}_i, \mathcal{C}_j) \quad (3)$$

There are many methods to compute the affinity between two clusters [17, 31, 41, 70, 68]. More details can be found in [25]. We now describe how the affinity is measured by \mathcal{A} in our approach.

3.3. Affinity Measure

First, we build a directed graph $G = \langle \mathcal{V}, \mathcal{E} \rangle$, where \mathcal{V} is the set of vertices corresponding to deep representations \mathbf{X} for I , and \mathcal{E} is the set of edges connecting vertices. We define an affinity matrix $\mathbf{W} \in \mathbb{R}^{n_s \times n_s}$ corresponding to the edge set. The weight from vertex \mathbf{x}_i to \mathbf{x}_j is defined by

$$\mathbf{W}(i, j) = \begin{cases} \exp(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{\sigma^2}), & \text{if } \mathbf{x}_j \in \mathcal{N}_i^{K_s} \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

where $\sigma^2 = \frac{a}{n_s K_s} \sum_{\mathbf{x}_i \in \mathbf{X}} \sum_{\mathbf{x}_j \in \mathcal{N}_i^{K_s}} \|\mathbf{x}_i - \mathbf{x}_j\|_2^2$. This way to build up a directed graph can be found in many pre-

vious works such as [70, 68]. Here, a and K_s are two pre-defined parameters (their values are listed in Table 2). After constructing a directed graph for samples, we then adopt the graph degree linkage in [68] to measure the affinity between cluster \mathcal{C}_i and \mathcal{C}_j , denoted by $\mathcal{A}(\mathcal{C}_i, \mathcal{C}_j)$.

3.4. A Recurrent Framework

Our key insight is that agglomerative clustering can be interpreted as a recurrent process in the sense that it merges clusters over multiple timesteps. Based on this insight, we propose a recurrent framework to combine the image clustering and representation learning processes.

As shown in Fig. 2, at the timestep t , images I are first fed into the CNN to get representations \mathbf{X}^t and then used in conjunction with previous hidden state \mathbf{h}^{t-1} to predict current hidden state \mathbf{h}^t , i.e., the image cluster labels at timestep t . In our context, the output at timestep t is $\mathbf{y}^t = \mathbf{h}^t$. Hence, at timestep t

$$\mathbf{X}^t = f_r(I|\theta^t) \quad (5a)$$

$$\mathbf{h}^t = f_m(\mathbf{X}^t, \mathbf{h}^{t-1}) \quad (5b)$$

$$\mathbf{y}^t = f_o(\mathbf{h}^t) = \mathbf{h}^t \quad (5c)$$

where f_r is a function to extract deep representations \mathbf{X}^t for input I using the CNN parameterized by θ^t , and f_m is a merging process for generating \mathbf{h}^t based on \mathbf{X}^t and \mathbf{h}^{t-1} .

In a typical Recurrent Neural Network, one would unroll all timesteps at each training iteration. In our case, that would involve performing agglomerative clustering until we obtain the desired number of clusters, and then update the CNN parameters by back-propagation.

In this work, we introduce a *partial unrolling* strategy, i.e., we split the overall T timesteps into multiple periods, and unroll one period at a time. The intuitive reason we unroll partially is that the representation of the CNN at the beginning is not reliable. We need to update CNN parameters to obtain more discriminative representations for the following merging processes. In each period, we merge a number of clusters and update CNN parameters for a fixed number of iterations at the end of the period. An extreme case would be one timestep per period, but it involves updating the CNN parameters too frequently and is thus time-consuming. Therefore, the number of timesteps per period (and thus the number of clusters merged per period) is determined by a parameter in our approach. We elaborate on this more in Sec. 3.6.

3.5. Objective Function

In our recurrent framework, we accumulate the losses from all timesteps, which is formulated as

$$\mathcal{L}(\{\mathbf{y}^1, \dots, \mathbf{y}^T\}, \{\theta^1, \dots, \theta^T\} | I) = \sum_{t=1}^T \mathcal{L}^t(\mathbf{y}^t, \theta^t | \mathbf{y}^{t-1}, I) \quad (6)$$

Here, \mathbf{y}^0 takes each image as a cluster. At timestep t , we find two clusters to merge given \mathbf{y}^{t-1} . In conventional agglomerative clustering, the two clusters are determined by finding the maximal affinity over all pairs of clusters. In this paper, we introduce a criterion that considers not only the affinity between two clusters but also the local structure surrounding the clusters. Assume from \mathbf{y}^{t-1} to \mathbf{y}^t , we merged a cluster \mathcal{C}_i^t and its nearest neighbour. Then the loss at timestep t is a combination of negative affinities, that is,

$$\mathcal{L}^t(\mathbf{y}^t, \boldsymbol{\theta}^t | \mathbf{y}^{t-1}, \mathbf{I}) = -\mathcal{A}(\mathcal{C}_i^t, \mathcal{N}_{\mathcal{C}_i^t}^{K_c}[1]) \quad (7a)$$

$$-\frac{\lambda}{(K_c - 1)} \sum_{k=2}^{K_c} \left(\mathcal{A}(\mathcal{C}_i^t, \mathcal{N}_{\mathcal{C}_i^t}^{K_c}[1]) - \mathcal{A}(\mathcal{C}_i^t, \mathcal{N}_{\mathcal{C}_i^t}^{K_c}[k]) \right) \quad (7b)$$

where λ weighs (7a) and (7b). Note that \mathbf{y}^t , \mathbf{y}^{t-1} and $\boldsymbol{\theta}^t$ are not explicitly presented at the right side, but they determine the loss via the image cluster labels and affinities among clusters. On the right side of the above equation, there are two terms: 1) (7a) measures the affinity between cluster \mathcal{C}_i and its nearest neighbour, which follows conventional agglomerative clustering; 2) (7b) measures the difference between affinity of \mathcal{C}_i to its nearest neighbour cluster and affinities of \mathcal{C}_i to its other neighbour clusters. This term takes the local structure into account. See Sec. 3.5.1 for detailed explanation.

It is hard to simultaneously derive the optimal $\{\mathbf{y}^1, \dots, \mathbf{y}^T\}$ and $\{\boldsymbol{\theta}^1, \dots, \boldsymbol{\theta}^T\}$ that minimize the overall loss in Eq. (6). As aforementioned, we optimize iteratively in a recurrent process. We divide T timesteps into P partially unrolled periods. In each period, we fix $\boldsymbol{\theta}$ and search optimal \mathbf{y} in the forward pass, and then in the backward pass we derive optimal $\boldsymbol{\theta}$ given the optimal \mathbf{y} . Details will be explained in the following sections.

3.5.1 Forward Pass

In forward pass of the p -th ($p \in \{1, \dots, P\}$) partially unrolled period, we update the cluster labels with $\boldsymbol{\theta}$ fixed to $\boldsymbol{\theta}^p$, and the overall loss in period p is

$$\mathcal{L}^p(\mathcal{Y}^p | \boldsymbol{\theta}^p, \mathbf{I}) = \sum_{t=t_p^s}^{t_p^e} \mathcal{L}^t(\mathbf{y}^t | \boldsymbol{\theta}^p, \mathbf{y}^{t-1}, \mathbf{I}) \quad (8)$$

where \mathcal{Y}^p is the sequence of image labels in period p , and $[t_p^s, t_p^e]$ is the corresponding timesteps in period p . For optimization, we follow a greedy search similar to conventional agglomerative clustering. Starting from the time step t_p^s , it finds one cluster and its nearest neighbour to merge so that \mathcal{L}^t is minimized over all possible cluster pairs.

In Fig. 3, we present a toy example to explain the reason why we employ the term (7b). As shown, it is often the case that the clusters are densely populated in some regions while sparse in some other regions. In conventional

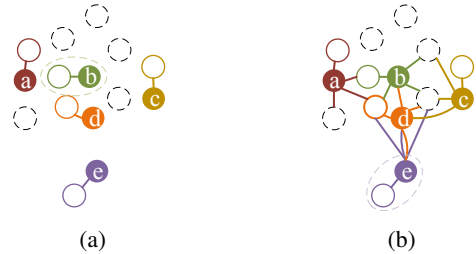


Figure 3: A toy illustration of (a) conventional agglomerative clustering strategy and (b) the proposed one. For simplification, we use a single circle to represent a cluster/sample. In conventional agglomerative clustering, node b and its nearest neighbour are chosen to merge because they are closest to each other; while node e is chosen in our proposed strategy considering the local structure.

agglomerative clustering, it will choose two clusters with largest affinity (or smallest loss) at each time no matter where the clusters are located. In this specific case, it will choose cluster \mathcal{C}_b and its nearest neighbour to merge. In contrast, as shown in Fig. 3(b), our algorithm by adding (7b) will find cluster \mathcal{C}_e , because it is not only close to its nearest neighbour, but also relatively far away from its other neighbours, i.e., the local structure is considered around one cluster. Another merit of introducing (7b) is that it will allow us to write the loss in terms of triplets as explained next.

3.5.2 Backward Pass

In forward pass of the p -th partially unrolled period, we have merged a number of clusters. Let the sequence of optimal image cluster labels be given by $\mathcal{Y}_*^p = \{\mathbf{y}_*^t\}$, and clusters merged in forward pass are denoted by $\{[\mathcal{C}_*^t, \mathcal{N}_{\mathcal{C}_*^t}^{K_c}[1]]\}$, $t \in \{t_p^s, \dots, t_p^e\}$. In the backward pass, we aim to derive the optimal $\boldsymbol{\theta}$ to minimize the losses generated in forward pass. Because the clustering in current period is conditioned on the clustering results of all previous periods, we accumulate the losses of all p periods, i.e.,

$$\mathcal{L}(\boldsymbol{\theta} | \{\mathcal{Y}_*^1, \dots, \mathcal{Y}_*^p\}, \mathbf{I}) = \sum_{k=1}^p \mathcal{L}^k(\boldsymbol{\theta} | \mathcal{Y}_*^k, \mathbf{I}) \quad (9)$$

Minimizing (9) w.r.t $\boldsymbol{\theta}$ leads to representation learning on \mathbf{I} supervised by $\{\mathcal{Y}_*^1, \dots, \mathcal{Y}_*^p\}$ or $\{\mathbf{y}_*^1, \dots, \mathbf{y}_*^{t_p^e}\}$. Based on (7a) and (7b), the loss in Eq. 9 is reformulated to

$$-\frac{\lambda}{K_c - 1} \sum_{t=1}^{t_p^e} \sum_{k=2}^{K_c} \left(\lambda' \mathcal{A}(\mathcal{C}_*^t, \mathcal{N}_{\mathcal{C}_*^t}^{K_c}[1]) - \mathcal{A}(\mathcal{C}_*^t, \mathcal{N}_{\mathcal{C}_*^t}^{K_c}[k]) \right) \quad (10)$$

where $\lambda' = (1 + 1/\lambda)$. (10) is a loss defined on clusters of points, which needs the entire dataset to estimate, making it difficult to use batch-based optimization. However,

Algorithm 1 Joint Optimization on \mathbf{y} and θ

Input:

I : = collection of image data;
 n_c^* : = target number of clusters;

Output:

\mathbf{y}^*, θ^* : = final image labels and CNN parameters;

- 1: $t \leftarrow 0; p \leftarrow 0$
 - 2: Initialize θ and \mathbf{y}
 - 3: **repeat**
 - 4: Update \mathbf{y}^t to \mathbf{y}^{t+1} by merging two clusters
 - 5: **if** $t = t_p^e$ **then**
 - 6: Update θ^p to θ^{p+1} by training CNN
 - 7: $p \leftarrow (p + 1)$
 - 8: **end if**
 - 9: $t \leftarrow t + 1$
 - 10: **until** Cluster number reaches n_c^*
 - 11: $\mathbf{y}^* \leftarrow \mathbf{y}^t; \theta^* \leftarrow \theta^p$
-

we show that this loss can be approximated by a sample-based loss, enabling us to compute unbiased estimators for the gradients using batch-statistics.

The intuition behind reformulation of the loss is that agglomerative clustering starts with each datapoint as a cluster, and clusters at a higher level in the hierarchy are formed by merging lower level clusters. Thus, affinities between clusters can be expressed in terms of affinities between datapoints. We show in the supplement that the loss in (10) can be approximately reformulated as

$$\mathcal{L}(\theta | \mathbf{y}_*^{t_p^e}, I) = -\frac{\lambda}{K_c - 1} \sum_{i,j,k} (\gamma \mathcal{A}(x_i, x_j) - \mathcal{A}(x_i, x_k)) \quad (11)$$

where γ is a weight whose value depends on λ' and how clusters are merged during the forward pass. x_i and x_j are from the same cluster, while x_k is from the neighbouring clusters, and their cluster labels are merely determined by the final clustering result $\mathbf{y}_*^{t_p^e}$. To further simplify the optimization, we instead search x_k in at most K_c neighbour samples of x_i from other clusters in a training batch. Hence, the batch-wise optimization can be performed using conventional stochastic gradient descent method. Note that such triplet losses have appeared in other works [59, 50]. Because it is associated with a weight, we call (35) the weighted triplet loss.

3.6. Optimization

Given an image dataset with n_s samples, we assume the number of desired clusters n_c^* is given to us as is standard in clustering. Then we can build up a recurrent process with $T = n_s - n_c^*$ timesteps, starting by regarding each sample as a cluster. However, such initialization makes the optimization time-consuming, especially when datasets contain

a large number of samples. To address this problem, we can first run a fast clustering algorithm to get the initial clusters. Here, we adopt the initialization algorithm proposed in [69] for fair comparison with their experiment results. Note that other kind of initializations can also be used, e.g. K-means. Based on the algorithm in [69], we obtain a number of clusters which contain a few samples for each (average is about 4 in our experiments). Given these initial clusters, our optimization algorithm learns deep representations and clusters. The algorithm is outlined in Alg. 1. In each partially unrolled period, we perform forward and backward passes to update \mathbf{y} and θ , respectively. Specifically, in the forward pass, we merge two clusters at each timestep. In the backward pass, we run about 20 epochs to update θ , and the affinity matrix W is also updated based on the new representation. The duration of the p -th period is $n_p = \text{ceil}(\eta \times n_c^s)$ timesteps, where n_c^s is the number of clusters at the beginning of current period, and η is a parameter called *unrolling rate* to control the number of timesteps. The less η is, the more frequently we update θ .

4. Experiments

4.1. Image Clustering

We compare our approach with 12 clustering algorithms, including K-means [39], NJW spectral clustering (SC-NJW) [43], self-tuning spectral clustering (SC-ST)[67], large-scale spectral clustering (SC-LS) [3], agglomerative clustering with average linkage (AC-Link)[25], Zeta function based agglomerative clustering (AC-Zell) [70], graph degree linkage-based agglomerative clustering (AC-GDL) [68], agglomerative clustering via path integral (AC-PIC) [69], normalized cuts (N-Cuts) [52], locality preserving non-negative matrix factorization (NMF-LP) [1], NMF with deep model (NMF-D) [57], task-specific clustering with deep model (TSC-D) [61].

For evaluation, we use a commonly used metric: normalized mutual information (NMI) [65]. It ranges in $[0, 1]$. Larger value indicates more precise clustering results.

4.1.1 Datasets

We evaluate the clustering performance on two handwritten digit image datasets (MNIST [34] and USPS¹), two multi-view object image datasets (COIL20 and COIL100 [42]), and four face image datasets (UMist [18], FRGC-v2.0², CMU-PIE [53], Youtube-Face (YTF)) [63]. The number of samples and categories, and image size are listed in Table 1. MNIST consists of training set (60,000) and testing set (10,000). To compare with different approaches, we experiment on the full set (MNIST-full) and testing set (MNIST-test), separately. For face image datasets such as

¹<http://www.cs.nyu.edu/~roweis/data.html>

²http://www3.nd.edu/~cvrl/CVRL/Data_Sets.html

Table 1: Datasets used in our experiments.

Dataset	<i>MNIST</i>	<i>USPS</i>	<i>COIL20</i>	<i>COIL100</i>	<i>UMist</i>	<i>FRGC-v2.0</i>	<i>CMU-PIE</i>	<i>YTF</i>
#Samples	70000	11000	1440	7200	575	2462	2856	10000
#Categories	10	10	20	100	20	20	68	41
Image Size	28×28	16×16	128×128	128×128	112×92	32×32	32×32	55×55

Table 2: Hyper-parameters in our approach.

Hyper-parameter	K_s	a	K_c	λ	γ	η
Value	20	1.0	5	1.0	2.0	0.9 or 0.2

UMist, CMU-PIE, we use the images provided as is without any changes. For FRGC-v2.0 and YTF datasets, we first crop faces and then resize them to a constant size. In FRGC-v2.0 dataset, we randomly choose 20 subjects. As for YTF dataset, we choose the first 41 subjects which are sorted by their names in alphabet order.

4.1.2 Experimental Setup

All the hyper-parameters and their values for our approach are listed in Table 2. In our experiments, K_s is set to 20, the same value to [68]. a and λ are simply set to 1.0. We search the values of K_c and γ for best performance on MNIST-test set. The unrolling rate η for first four datasets is 0.9; and 0.2 for face datasets. The target cluster number n_c^* is set to be the number of categories in each dataset.

We use Caffe [27] to implement our approach. We stacked multiple combinations of convolutional layer, batch normalization layer, ReLU layer and pooling layer. For all the convolutional layers, the number of channels is 50, and filter size is 5×5 with stride = 1 and padding = 0. For pooling layer, its kernel size is 2 and stride is 2. To deal with varying image sizes across datasets, the number of stacked convolutional layers for each dataset is chosen so that the size of the output feature map is about 10×10 . On the top of all CNNs, we append an inner product (*ip*) layer whose dimension is 160. *ip* layer is followed by a L2-normalization layer before being fed to the weighted triplet loss layer or used for clustering. For each partially unrolled period, the base learning rate is set to 0.01, momentum 0.9, and weight decay 5×10^{-5} . We use the *inverse* learning rate decay policy, with *Gamma*=0.0001 and *Power*=0.75. Stochastic gradient descent (SGD) is adopted for optimization.

4.1.3 Quantitative Comparison

We report NMI for different methods on various datasets. Results are averaged from 3 runs. We report the results by re-running the code released by original papers. For those that did not release the code, the corresponding results are

borrowed from the papers. We find the results we obtain are somewhat different from the one reported in original papers. We suspect that these differences may be caused by the different experimental settings or the released code is changed from the one used in the original paper. For all test algorithms, we conduct L2-normalization on the image intensities since it empirically improves the clustering performance. We report our own results in two cases: 1) the straight-forward clustering results obtained when the recurrent process finish, denoted by OURS-SF; 2) the clustering results obtained by re-running clustering algorithm after obtaining the final representation, denoted by OURS-RC. The quantitative results are shown in Table 3. In the table cells, the value before '/' is obtained by re-running code while the value after '/' is that reported in previous papers.

As we can see from Table 3, both OURS-SF and OURS-RC outperform previous methods on all datasets with noticeable margin. Interestingly, we achieved perfect results (NMI = 1) on COIL20 and CMU-PIE datasets, which means that all samples in the same category are clustered into the same group. The agglomerative clustering algorithms, such as AC-Zell, AC-GDL and AC-PIC perform better than other algorithms generally. However, on MNIST-full test, they all perform poorly. The possible reason is that MNIST-full has 70k samples, and these methods cannot cope with such large-scale dataset when using image intensity as representation. However, this problem is addressed by our learned representation. We show that we achieved analogous performance on MNIST-full to MNIST-test set. In most cases, we can find OURS-RC performs better on datasets that have room for improvement. We believe the reason is that OURS-RC uses the final learned representation over the entire clustering process, while OURS-SF starts with image intensity, which indicates that the learned representation is more discriminative than image intensity.³

4.1.4 Generalization Across Clustering Algorithms

We now evaluate if the representations learned by our joint agglomerative clustering and representation learning approach generalize to other clustering techniques. We re-run all the clustering algorithms without any changes of parameters, but using our learned deep representations as fea-

³We experimented with hand-crafted features such as HOG, LBP, spatial pyramid on a subset of the datasets with some of the better clustering algorithms from Table 3, and found that they performed worse.

Table 3: Quantitative clustering performance (NMI) for different algorithms using image intensities as input.

Dataset	<i>COIL20</i>	<i>COIL100</i>	<i>USPS</i>	<i>MNIST-test</i>	<i>MNIST-full</i>	<i>UMist</i>	<i>FRGC</i>	<i>CMU-PIE</i>	<i>YTF</i>
K-means [39]	0.775	0.822	0.447	0.528	0.500	0.609	0.389	0.549	0.761
SC-NJW [43]	0.860/0.889	0.872/0.854	0.409/0.690	0.528/0.755	0.476	0.727	0.186	0.543	0.752
SC-ST [67]	0.673/0.895	0.706/0.858	0.342/0.726	0.445/0.756	0.416	0.611	0.431	0.581	0.620
SC-LS [3]	0.877	0.833	0.681	0.756	0.706	0.810	0.550	0.788	0.759
N-Cuts [52]	0.768/0.884	0.861/0.823	0.382/0.675	0.386/0.753	0.411	0.782	0.285	0.411	0.742
AC-Link [25]	0.512	0.711	0.579	0.662	0.686	0.643	0.168	0.545	0.738
AC-Zell [70]	0.954/0.911	0.963/0.913	0.774/0.799	0.810/0.768	0.017	0.755	0.351	0.910	0.733
AC-GDL [68]	0.945/0.937	0.954/0.929	0.854/0.824	0.864/0.844	0.017	0.755	0.351	0.934	0.622
AC-PIC [69]	0.950	0.964	0.840	0.853	0.017	0.750	0.415	0.902	0.697
NMF-LP [1]	0.720	0.783	0.435	0.467	0.452	0.560	0.346	0.491	0.720
NMF-D [57]	0.692	0.719	0.286	0.243	0.148	0.500	0.258	0.983/0.910	0.569
TSC-D [61]	-/0.928	-	-	-	-/0.651	-	-	-	-
OURS-SF	1.000	0.978	0.858	0.876	0.906	0.880	0.566	0.984	0.848
OURS-RC	1.000	0.985	0.913	0.915	0.913	0.877	0.574	1.00	0.848

Table 4: Quantitative clustering performance (NMI) for different algorithms using our learned representations as inputs.

Dataset	<i>COIL20</i>	<i>COIL100</i>	<i>USPS</i>	<i>MNIST-test</i>	<i>MNIST-full</i>	<i>UMist</i>	<i>FRGC</i>	<i>CMU-PIE</i>	<i>YTF</i>
K-means [39]	0.926	0.919	0.758	0.908	0.927	0.871	0.636	0.956	0.835
SC-NJW [43]	0.915	0.898	0.753	0.878	0.931	0.833	0.625	0.957	0.789
SC-ST [67]	0.959	0.922	0.741	0.911	0.906	0.847	0.651	0.938	0.741
SC-LS [3]	0.950	0.905	0.780	0.912	0.932	0.879	0.639	0.950	0.802
N-Cuts [52]	0.963	0.900	0.705	0.910	0.930	0.877	0.640	0.995	0.823
AC-Link [25]	0.896	0.884	0.783	0.901	0.918	0.872	0.621	0.990	0.803
AC-Zell [70]	1.000	0.989	0.910	0.893	0.919	0.870	0.551	1.000	0.821
AC-GDL [68]	1.000	0.985	0.913	0.915	0.913	0.870	0.574	1.000	0.842
AC-PIC [69]	1.000	0.990	0.914	0.909	0.907	0.870	0.553	1.000	0.829
NMF-LP [1]	0.855	0.834	0.729	0.905	0.926	0.854	0.575	0.690	0.788

tures. The results are shown in Table 4. It can be seen that all clustering algorithms obtain more precise image clusters by using our learned representation. Some algorithms like K-means, AC-Link that performed very poorly with raw intensities perform much better with our learned representations, and the variance in performance across all clustering algorithms is much lower. These results clearly demonstrate that our learned representation is not over-fitting to a single clustering algorithm, but generalizes well across various algorithms. Interestingly, using our learned representation, some of the clustering algorithms perform even better than AC-GDL we build on in our approach.

4.2. Transferring Learned Representation

4.2.1 Cross-Dataset Clustering

Table 5: NMI performance across COIL20 and COIL100.

Layer	<i>data</i>	<i>top(ip)</i>	top-1	top-2
COIL20 → COIL100	0.924	0.927	0.939	0.934
COIL100 → COIL20	0.944	0.949	0.957	0.951

Table 6: NMI performance across MNIST-test and USPS.

Layer	<i>data</i>	<i>top(ip)</i>	top-1	top-2
MNIST-test → USPS	0.874	0.892	0.907	0.908
USPS → MNIST-test	0.872	0.873	0.886	-

In this section, we study whether our learned representations generalize across datasets. We train a CNN based on our approach on one dataset, and then cluster images from another (but related) dataset using the image features extracted via the CNN. Specifically, we experiment on two dataset pairs: 1) multi-view object datasets (COIL20 and COIL100); 2) hand-written digit datasets (USPS and MNIST-test). We use the representation learned from one dataset to represent another dataset, followed by agglomerative clustering. Note that because the image sizes or channels are different across datasets, we resize the input images and/or expand the channels before feeding them to CNN. The experimental results are shown in Table 5 and 6. We use the representations from top *ip* layer and also the *convolutional* or *pooling* layers (top-1, top-2) close to top layer for image clustering. In two tables, compared with

Table 7: Face verification results on LFW.

#Samples	10k	20k	30k	50k	100k
Supervised	0.737	0.746	0.748	0.764	0.770
OURS	0.728	0.743	0.750	0.762	0.767

directly using raw image from the *data* layer, the clustering performance based on learned representations from all layers improve, which indicates that the learned representations can be transferred across these datasets. As perhaps expected, the performance on target datasets is worse compared to learning on the target dataset directly. For COIL20 and COIL100, a possible reason is that they have different image categories. As for MNIST and USPS, the performance beats OURS-SF, but worse than OURS-RC. We find transferring representation learned on MNIST-test to USPS gets close performance to OURS-RC learned on USPS.

4.2.2 Face Verification

We now evaluate the performance of our approach by applying it to face verification. In particular, the representation is learned on Youtube-Face dataset and evaluated on LFW dataset [23] under the restricted protocol. For training, we randomly choose about 10k, 20k, 30k, 50k, 100k samples from YTF dataset. All these subsets have 1446 categories. We implement our approach to train CNN model and cluster images on the training set. Then, we remove the L2-normalization layer and append a softmax layer to fine-tune our unsupervised CNN model *based on the predicted image cluster labels*. Using the same training samples and CNN architecture, we also train a CNN model with a softmax loss supervised by the groundtruth labels of the training set. According to the evaluation protocol in [23], we run 10-fold cross-validation. The cosine similarity is used to compute the similarity between samples. In each of 10 cross-validations, nine folds are used to find the optimal threshold, and the remaining one fold is used for evaluation. The average accuracy is reported in Table. 7. As shown, though no groundtruth labels are used for representation learning in our approach, we obtain analogous performance to the supervised learning approach. Our approach even (slightly) beats the supervised learning method in one case.

4.3. Image Classification

Recently, unsupervised representation learning methods are starting to achieve promising results for a variety of recognition tasks [5, 4, 26, 36]. We are interested in knowing whether the proposed method can also learn useful representation for image classification. We experiment with CIFAR-10 [28]. We follow the pipeline in [5], and base our experiments on their publicly available code. In this pipeline, codebook with 1600 codes is build upon 6×6

Table 8: Image classification accuracy on CIFAR-10.

#Samples	K-means [5]	conv1	conv2	conv1&2
5k	62.81%	63.05%	63.10%	63.50%
10k	68.01%	68.30%	68.46%	69.11%
25k	74.01%	72.83%	72.93%	75.11%
50k (full set)	76.59%	74.68%	74.68%	78.55%

ZCA-whitened image patches, and then used to code the training and testing samples by extracting 1,600-d feature from each of 4 image quadrants. Afterwards, a linear SVM [6] is applied for image classification on 6,400-d feature. In our approach, the only difference is that we learn a new representation from 6×6 patches, and then use these new representations to build the codebook with 1,600 codes. The CNN architecture we use contains two convolutional layers, each of which is combined with a ReLu and a pooling layer, followed by an inner product layer. Both convolutional layers have 50×3 filters with $\text{pad} = 1$. The kernel size of pooling layer is 2, and the stride is 2. To save on training time, 40k randomly extracted patches are extracted from 50k training set and used in all the experiments.

Classification accuracies on test set with different settings are shown in Table 8. We vary the number of training samples and evaluate the performance for representations from different layers. As we can see, the combination of representations from the first and second convolutional layer achieve the best performance. We also use the representation output by inner product layer to learn the codebook. However, it performs poorly. A possible reason is that it discards spatial information of image patches, which may be important for learning a codebook. When using 400k randomly extracted patches to learn the codebook, [5] achieved 77.9%. However, it is still lower than what we achieved. This performance also beats several other methods listed in [4, 16, 26, 36].

5. Conclusion

In this paper, we have proposed an approach to jointly learn deep representations and image clusters. In our approach, we combined agglomerative clustering with CNNs and formulate them as a recurrent process. We used a partially unrolling strategy to divide the timesteps into multiple periods. In each period, we merged clusters step by step during the forward pass and learned representation in the backward pass, which are guided by a single weighted triplet-loss function. The extensive experiments on image clustering, deep representation transfer learning and image classification demonstrate that our approach can obtain more precise image clusters and discriminative representations that generalize well across many datasets and tasks.

6. Acknowledgements

This work was supported in part by the Paul G. Allen Family Foundation, Google, and Institute for Critical Technology and Applied Science (ICTAS) at Virginia Tech through awards to D. P.; and by a National Science Foundation CAREER award, an Army Research Office YIP award, an Office of Naval Research grant, an AWS in Education Research Grant, and GPU support by NVIDIA to D. B. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the U.S. Government or any sponsor.

References

- [1] D. Cai, X. He, X. Wang, H. Bao, and J. Han. Locality preserving nonnegative matrix factorization. In *IJCAI*, volume 9, pages 1010–1015, 2009. 2, 5, 7, 15
- [2] G. Chen. Deep learning with nonparametric clustering. *arXiv preprint arXiv:1501.03084*, 2015. 2
- [3] X. Chen and D. Cai. Large scale spectral clustering with landmark-based representation. In *AAAI*, 2011. 5, 7, 14, 15, 16
- [4] A. Coates and A. Y. Ng. Selecting receptive fields in deep networks. In *NIPS*, pages 2528–2536, 2011. 8
- [5] A. Coates, A. Y. Ng, and H. Lee. An analysis of single-layer networks in unsupervised feature learning. In *International conference on artificial intelligence and statistics*, pages 215–223, 2011. 8
- [6] C. Cortes and V. Vapnik. Support-vector networks. *Machine learning*, 20(3):273–297, 1995. 8
- [7] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *CVPR*, volume 1, pages 886–893. IEEE, 2005. 14
- [8] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the royal statistical society. Series B (methodological)*, pages 1–38, 1977. 2
- [9] C. Ding, T. Li, M. Jordan, et al. Convex and semi-nonnegative matrix factorizations. *IEEE TPAMI*, 32(1):45–55, 2010. 2
- [10] C. Doersch, A. Gupta, and A. A. Efros. Mid-level visual element discovery as discriminative mode seeking. In *NIPS*, pages 494–502, 2013. 2
- [11] C. Doersch, A. Gupta, and A. A. Efros. Unsupervised visual representation learning by context prediction. In *ICCV*, pages 1422–1430, 2015. 2
- [12] A. Dosovitskiy, J. T. Springenberg, M. Riedmiller, and T. Brox. Discriminative unsupervised feature learning with convolutional neural networks. In *NIPS*, pages 766–774, 2014. 2
- [13] Y. Gdalyahu, D. Weinshall, and M. Werman. Self-organization in vision: stochastic clustering for image segmentation, perceptual grouping, and image database organization. *IEEE TPAMI*, 23(10):1053–1074, 2001. 2
- [14] R. Girshick. Fast r-cnn. In *ICCV*, pages 1440–1448, 2015. 2
- [15] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, pages 580–587. IEEE, 2014. 2
- [16] I. J. Goodfellow, A. Courville, and Y. Bengio. Spike-and-slab sparse coding for unsupervised feature discovery. *arXiv preprint arXiv:1201.3382*, 2012. 8
- [17] K. C. Gowda and G. Krishna. Agglomerative clustering using the concept of mutual nearest neighbourhood. *Pattern recognition*, 10(2):105–112, 1978. 1, 2, 3
- [18] D. B. Graham and N. M. Allinson. Characterising virtual eigensignatures for general purpose face recognition. In *Face Recognition*, pages 446–456. Springer, 1998. 5
- [19] D. Han and J. Kim. Unsupervised simultaneous orthogonal basis clustering feature selection. In *IEEE CVPR*, pages 5016–5023, 2015. 2
- [20] B. Hariharan, J. Malik, and D. Ramanan. Discriminative decorrelation for clustering and classification. In *ECCV*, pages 459–472. Springer, 2012. 2
- [21] K. He, X. Zhang, S. Ren, and J. Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. In *ECCV*, pages 346–361. Springer, 2014. 2
- [22] G. E. Hinton and R. R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507, 2006. 2, 17
- [23] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. Technical report, Technical Report 07-49, University of Massachusetts, Amherst, 2007. 8
- [24] H.-C. Huang, Y.-Y. Chuang, and C.-S. Chen. Affinity aggregation for spectral clustering. In *CVPR*, pages 773–780. IEEE, 2012. 2
- [25] A. K. Jain, M. N. Murty, and P. J. Flynn. Data clustering: a review. *ACM computing surveys (CSUR)*, 31(3):264–323, 1999. 2, 3, 5, 7, 15
- [26] Y. Jia, C. Huang, and T. Darrell. Beyond spatial pyramids: Receptive field learning for pooled image features. In *CVPR*, pages 3370–3377. IEEE, 2012. 8
- [27] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. In *Proceedings of the ACM International Conference on Multimedia*, pages 675–678. ACM, 2014. 6
- [28] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images, 2009. 8
- [29] A. Krizhevsky and G. E. Hinton. Using very deep autoencoders for content-based image retrieval. In *ESANN*. Citeseer, 2011. 2
- [30] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, pages 1097–1105, 2012. 2
- [31] T. Kurita. An efficient agglomerative clustering algorithm using a heap. *Pattern Recognition*, 24(3):205–209, 1991. 2, 3
- [32] S. Lazebnik, C. Schmid, and J. Ponce. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *CVPR*, volume 2, pages 2169–2178. IEEE, 2006. 14

- [33] Q. V. Le. Building high-level features using large scale unsupervised learning. In *ICASSP*, pages 8595–8598. IEEE, 2013. 2
- [34] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. 1, 5
- [35] H. Lee, R. Grosse, R. Ranganath, and A. Y. Ng. Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. In *ICML*, pages 609–616. ACM, 2009. 2
- [36] T.-H. Lin and H. Kung. Stable and efficient representation learning with nonnegativity constraints. In *ICML*, pages 1323–1331, 2014. 8
- [37] D. G. Lowe. Object recognition from local scale-invariant features. In *ICCV*, volume 2, pages 1150–1157. IEEE, 1999. 2
- [38] L. Maaten. Learning a parametric embedding by preserving local structure. In *International Conference on Artificial Intelligence and Statistics*, pages 384–391, 2009. 16, 17
- [39] J. MacQueen et al. Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, volume 1, pages 281–297. Oakland, CA, USA., 1967. 2, 5, 7, 15
- [40] G. McLachlan and D. Peel. *Finite mixture models*. John Wiley & Sons, 2004. 2
- [41] J. F. Navarro, C. S. Frenk, and S. D. White. A universal density profile from hierarchical clustering. *The Astrophysical Journal*, 490(2):493, 1997. 3
- [42] S. A. Nene, S. K. Nayar, H. Murase, et al. Columbia object image library (coil-20). Technical report, Technical Report CUCS-005-96, 1996. 5
- [43] A. Y. Ng, M. I. Jordan, Y. Weiss, et al. On spectral clustering: Analysis and an algorithm. *NIPS*, 2:849–856, 2002. 2, 5, 7, 15
- [44] T. Ojala, M. Pietikäinen, and D. Harwood. A comparative study of texture measures with classification based on featured distributions. *Pattern recognition*, 29(1):51–59, 1996. 14
- [45] M. A. Ranzato, F. J. Huang, Y.-L. Boureau, and Y. LeCun. Unsupervised learning of invariant feature hierarchies with applications to object recognition. In *CVPR*, pages 1–8. IEEE, 2007. 2
- [46] K. Rematas, B. Fernando, F. Dellaert, and T. Tuytelaars. Dataset fingerprints: Exploring image collections through data mining. In *IEEE CVPR*, pages 4867–4875, 2015. 2
- [47] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, pages 91–99, 2015. 2
- [48] S. Roweis, G. Hinton, and R. Salakhutdinov. Neighbourhood component analysis. *Advances in Neural Information Processing Systems (NIPS)*, 17:513–520, 2004. 16, 17
- [49] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, pages 1–42, 2014. 2
- [50] F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, pages 815–823, 2015. 5, 14
- [51] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. *arXiv preprint arXiv:1312.6229*, 2013. 2
- [52] J. Shi and J. Malik. Normalized cuts and image segmentation. *IEEE TPAMI*, 22(8):888–905, 2000. 2, 5, 7, 14, 15, 16
- [53] T. Sim, S. Baker, and M. Bsat. The cmu pose, illumination, and expression (pie) database. In *FG*, pages 46–51. IEEE, 2002. 5
- [54] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 2
- [55] S. Singh, A. Gupta, and A. Efros. Unsupervised discovery of mid-level discriminative patches. *ECCV*, pages 73–86, 2012. 2
- [56] F. Tian, B. Gao, Q. Cui, E. Chen, and T.-Y. Liu. Learning deep representations for graph clustering. In *AAAI*, pages 1293–1299, 2014. 2
- [57] G. Trigeorgis, K. Bousmalis, S. Zafeiriou, and B. Schuller. A deep semi-nmf model for learning hidden representations. In *ICML*, pages 1692–1700, 2014. 2, 5, 7
- [58] P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol. Extracting and composing robust features with denoising autoencoders. In *ICML*, pages 1096–1103. ACM, 2008. 2
- [59] J. Wang, Y. Song, T. Leung, C. Rosenberg, J. Wang, J. Philbin, B. Chen, and Y. Wu. Learning fine-grained image similarity with deep ranking. In *CVPR*, pages 1386–1393. IEEE, 2014. 5, 14
- [60] X. Wang and A. Gupta. Unsupervised learning of visual representations using videos. In *ICCV*, pages 2794–2802, 2015. 2
- [61] Z. Wang, S. Chang, J. Zhou, and T. S. Huang. Learning a task-specific deep architecture for clustering. In *arXiv preprint arXiv:1509.00151*, 2015. 2, 5, 7
- [62] S. Wold, K. Esbensen, and P. Geladi. Principal component analysis. *Chemometrics and intelligent laboratory systems*, 2(1-3):37–52, 1987. 16, 17
- [63] L. Wolf, T. Hassner, and I. Maoz. Face recognition in unconstrained videos with matched background similarity. In *CVPR*, pages 529–534. IEEE, 2011. 5
- [64] P. Xie and E. Xing. Integrating image clustering and codebook learning. In *AAAI*, 2015. 2
- [65] W. Xu, X. Liu, and Y. Gong. Document clustering based on non-negative matrix factorization. In *Proceedings of the 26th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 267–273. ACM, 2003. 5
- [66] S. Zafeiriou and M. Petrou. Nonlinear non-negative component analysis algorithms. *IEEE TIP*, 19(4):1050–1066, 2010. 2
- [67] L. Zelnik-Manor and P. Perona. Self-tuning spectral clustering. In *NIPS*, pages 1601–1608, 2004. 2, 5, 7, 15

- [68] W. Zhang, X. Wang, D. Zhao, and X. Tang. Graph degree linkage: Agglomerative clustering on a directed graph. In *ECCV*, pages 428–441. Springer, 2012. 3, 5, 6, 7, 11, 15
- [69] W. Zhang, D. Zhao, and X. Wang. Agglomerative clustering via maximum incremental path integral. *Pattern Recognition*, 46(11):3056–3065, 2013. 5, 7, 14, 15, 16
- [70] D. Zhao and X. Tang. Cyclizing clusters via zeta function of a graph. In *NIPS*, pages 1953–1960, 2009. 3, 5, 7, 15

A. Appendix

A.1. Affinity Measure for Clusters

In this paper, we employ the affinity measure in [68]

$$\begin{aligned} \mathcal{A}(\mathcal{C}_i, \mathcal{C}_j) &= \mathcal{A}(\mathcal{C}_j \rightarrow \mathcal{C}_i) + \mathcal{A}(\mathcal{C}_i \rightarrow \mathcal{C}_j) \\ &= \frac{1}{|\mathcal{C}_i|^2} \mathbf{1}_{|\mathcal{C}_i|}^T \mathbf{W}_{\mathcal{C}_i, \mathcal{C}_j} \mathbf{W}_{\mathcal{C}_j, \mathcal{C}_i} \mathbf{1}_{|\mathcal{C}_i|} \\ &\quad + \frac{1}{|\mathcal{C}_j|^2} \mathbf{1}_{|\mathcal{C}_j|}^T \mathbf{W}_{\mathcal{C}_j, \mathcal{C}_i} \mathbf{W}_{\mathcal{C}_i, \mathcal{C}_j} \mathbf{1}_{|\mathcal{C}_j|} \end{aligned} \quad (12)$$

where \mathbf{W} is the affinity matrix for samples, and $\mathbf{W}_{\mathcal{C}_i, \mathcal{C}_j} \in \mathbb{R}^{|\mathcal{C}_i| \times |\mathcal{C}_j|}$ is the submatrix in \mathbf{W} pointing from samples in \mathcal{C}_i to samples in \mathcal{C}_j , and $\mathbf{W}_{\mathcal{C}_j, \mathcal{C}_i} \in \mathbb{R}^{|\mathcal{C}_j| \times |\mathcal{C}_i|}$ is the one pointing from \mathcal{C}_j to \mathcal{C}_i . $\mathbf{1}_{|\mathcal{C}_i|}$ and $\mathbf{1}_{|\mathcal{C}_j|}$ are two vectors with all $|\mathcal{C}_i|$ and $|\mathcal{C}_j|$ elements be 1, respectively. Therefore, we have $\mathcal{A}(\mathcal{C}_i, \mathcal{C}_j) = \mathcal{A}(\mathcal{C}_j, \mathcal{C}_i)$.

According to (12), we can derive

$$\mathcal{A}((\mathcal{C}_m \cup \mathcal{C}_n) \rightarrow \mathcal{C}_i) = \mathcal{A}(\mathcal{C}_m \rightarrow \mathcal{C}_i) + \mathcal{A}(\mathcal{C}_n \rightarrow \mathcal{C}_i) \quad (13)$$

which has also been shown in [68]. Meanwhile,

$$\begin{aligned} \mathcal{A}(\mathcal{C}_i \rightarrow (\mathcal{C}_m \cup \mathcal{C}_n)) &= \beta \mathbf{1}_{|\mathcal{C}_m| + |\mathcal{C}_n|}^T \mathbf{W}_{\mathcal{C}_m \cup \mathcal{C}_n, \mathcal{C}_i} \mathbf{W}_{\mathcal{C}_i, \mathcal{C}_m \cup \mathcal{C}_n} \mathbf{1}_{|\mathcal{C}_m| + |\mathcal{C}_n|} \\ &= \beta \mathbf{1}_{|\mathcal{C}_m|}^T \mathbf{W}_{\mathcal{C}_m, \mathcal{C}_i} \mathbf{W}_{\mathcal{C}_i, \mathcal{C}_m} \mathbf{1}_{|\mathcal{C}_m|} + \beta \mathbf{1}_{|\mathcal{C}_n|}^T \mathbf{W}_{\mathcal{C}_n, \mathcal{C}_i} \mathbf{W}_{\mathcal{C}_i, \mathcal{C}_n} \mathbf{1}_{|\mathcal{C}_n|} \\ &\quad + \beta \mathbf{1}_{|\mathcal{C}_m|}^T \mathbf{W}_{\mathcal{C}_m, \mathcal{C}_i} \mathbf{W}_{\mathcal{C}_i, \mathcal{C}_n} \mathbf{1}_{|\mathcal{C}_n|} + \beta \mathbf{1}_{|\mathcal{C}_n|}^T \mathbf{W}_{\mathcal{C}_n, \mathcal{C}_i} \mathbf{W}_{\mathcal{C}_i, \mathcal{C}_m} \mathbf{1}_{|\mathcal{C}_m|} \end{aligned} \quad (14)$$

where $\beta = 1/(|\mathcal{C}_m| + |\mathcal{C}_n|)^2$.

A.2. Approximated Affinity Measure

During agglomerative clustering, we need to re-compute the affinity between the merged cluster to all other clusters based on 13 and 14 repeatedly. It is simple to compute 13. However, to get $\mathcal{A}(\mathcal{C}_i \rightarrow (\mathcal{C}_m \cup \mathcal{C}_n))$, we need a lot of computations. These time costs become dominant and remarkable when we have a large-scale dataset. To accelerate the computations, we introduce an approximation method. At the right side of (14), we assume samples in \mathcal{C}_m and \mathcal{C}_n have similar affinities to \mathcal{C}_i . This assumption is mild because the condition to merge \mathcal{C}_m and \mathcal{C}_n is that they are similar to each other. In this case, the ratio between $\mathbf{W}_{\mathcal{C}_i, \mathcal{C}_m} \mathbf{1}_{|\mathcal{C}_m|}$ and $\mathbf{W}_{\mathcal{C}_i, \mathcal{C}_n} \mathbf{1}_{|\mathcal{C}_n|}$ is analogy to the ratio between the number of

samples in two set, i.e.,

$$\mathbf{W}_{\mathcal{C}_i, \mathcal{C}_m} \mathbf{1}_{|\mathcal{C}_m|} = \frac{|\mathcal{C}_m|}{|\mathcal{C}_n|} \mathbf{W}_{\mathcal{C}_i, \mathcal{C}_n} \mathbf{1}_{|\mathcal{C}_n|} \quad (15)$$

Based on (15), we have

$$\mathbf{1}_{|\mathcal{C}_m|}^T \mathbf{W}_{\mathcal{C}_m, \mathcal{C}_i} \mathbf{W}_{\mathcal{C}_i, \mathcal{C}_n} \mathbf{1}_{|\mathcal{C}_n|} = \frac{|\mathcal{C}_n|}{|\mathcal{C}_m|} \mathbf{1}_{|\mathcal{C}_m|}^T \mathbf{W}_{\mathcal{C}_m, \mathcal{C}_i} \mathbf{W}_{\mathcal{C}_i, \mathcal{C}_m} \mathbf{1}_{|\mathcal{C}_m|} \quad (16a)$$

$$\mathbf{1}_{|\mathcal{C}_n|}^T \mathbf{W}_{\mathcal{C}_n, \mathcal{C}_i} \mathbf{W}_{\mathcal{C}_i, \mathcal{C}_m} \mathbf{1}_{|\mathcal{C}_m|} = \frac{|\mathcal{C}_m|}{|\mathcal{C}_n|} \mathbf{1}_{|\mathcal{C}_n|}^T \mathbf{W}_{\mathcal{C}_n, \mathcal{C}_i} \mathbf{W}_{\mathcal{C}_i, \mathcal{C}_n} \mathbf{1}_{|\mathcal{C}_n|} \quad (16b)$$

As a result, we can re-formulate (14) to

$$\begin{aligned} \mathcal{A}(\mathcal{C}_i \rightarrow (\mathcal{C}_m \cup \mathcal{C}_n)) &= \frac{1}{(|\mathcal{C}_m|^2 + |\mathcal{C}_n||\mathcal{C}_n|)} \mathbf{1}_{|\mathcal{C}_m|}^T \mathbf{W}_{\mathcal{C}_m, \mathcal{C}_i} \mathbf{W}_{\mathcal{C}_i, \mathcal{C}_m} \mathbf{1}_{|\mathcal{C}_m|} \\ &\quad + \frac{1}{(|\mathcal{C}_m||\mathcal{C}_n| + |\mathcal{C}_n|^2)} \mathbf{1}_{|\mathcal{C}_n|}^T \mathbf{W}_{\mathcal{C}_n, \mathcal{C}_i} \mathbf{W}_{\mathcal{C}_i, \mathcal{C}_n} \mathbf{1}_{|\mathcal{C}_n|} \end{aligned} \quad (17)$$

Therefore, we have

$$\begin{aligned} \mathcal{A}(\mathcal{C}_i \rightarrow (\mathcal{C}_m \cup \mathcal{C}_n)) &= \frac{|\mathcal{C}_m|}{|\mathcal{C}_m| + |\mathcal{C}_n|} \mathcal{A}(\mathcal{C}_i \rightarrow \mathcal{C}_m) \\ &\quad + \frac{|\mathcal{C}_n|}{|\mathcal{C}_m| + |\mathcal{C}_n|} \mathcal{A}(\mathcal{C}_i \rightarrow \mathcal{C}_n) \end{aligned} \quad (18)$$

Consequently, we have

$$\begin{aligned} \mathcal{A}(\mathcal{C}_m \cup \mathcal{C}_n, \mathcal{C}_i) &= \mathcal{A}(\mathcal{C}_m \rightarrow \mathcal{C}_i) + \mathcal{A}(\mathcal{C}_n \rightarrow \mathcal{C}_i) \\ &\quad + \frac{|\mathcal{C}_m|}{|\mathcal{C}_m| + |\mathcal{C}_n|} \mathcal{A}(\mathcal{C}_i \rightarrow \mathcal{C}_m) \\ &\quad + \frac{|\mathcal{C}_n|}{|\mathcal{C}_m| + |\mathcal{C}_n|} \mathcal{A}(\mathcal{C}_i \rightarrow \mathcal{C}_n) \end{aligned} \quad (19)$$

Above approximation provides us a potential way to reduce the computational complexity of agglomerative clustering. Though we computed $\mathcal{A}(\mathcal{C}_i \rightarrow (\mathcal{C} \cup \mathcal{C}_n))$ based on Eq. (14) in all our experiments, we found the approximation version achieves analogy performance while costs much less time than the original one. We further simplify the computation by assuming a constant ratio α between the terms at the right side of Eq. (14):

$$\mathbf{1}_{|\mathcal{C}_m|}^T \mathbf{W}_{\mathcal{C}_m, \mathcal{C}_i} \mathbf{W}_{\mathcal{C}_i, \mathcal{C}_n} \mathbf{1}_{|\mathcal{C}_n|} = \alpha \mathbf{1}_{|\mathcal{C}_m|}^T \mathbf{W}_{\mathcal{C}_m, \mathcal{C}_i} \mathbf{W}_{\mathcal{C}_i, \mathcal{C}_m} \mathbf{1}_{|\mathcal{C}_m|} \quad (20a)$$

$$\mathbf{1}_{|\mathcal{C}_n|}^T \mathbf{W}_{\mathcal{C}_n, \mathcal{C}_i} \mathbf{W}_{\mathcal{C}_i, \mathcal{C}_m} \mathbf{1}_{|\mathcal{C}_m|} = \alpha \mathbf{1}_{|\mathcal{C}_n|}^T \mathbf{W}_{\mathcal{C}_n, \mathcal{C}_i} \mathbf{W}_{\mathcal{C}_i, \mathcal{C}_n} \mathbf{1}_{|\mathcal{C}_n|} \quad (20b)$$

Based on above assumption,

$$\begin{aligned} \mathcal{A}(\mathcal{C}_m \cup \mathcal{C}_n, \mathcal{C}_i) &= \mathcal{A}(\mathcal{C}_m \rightarrow \mathcal{C}_i) + \mathcal{A}(\mathcal{C}_n \rightarrow \mathcal{C}_i) \\ &\quad + \frac{(1 + \alpha)|\mathcal{C}_m|^2}{(|\mathcal{C}_m| + |\mathcal{C}_n|)^2} \mathcal{A}(\mathcal{C}_i \rightarrow \mathcal{C}_m) \\ &\quad + \frac{(1 + \alpha)|\mathcal{C}_n|^2}{(|\mathcal{C}_m| + |\mathcal{C}_n|)^2} \mathcal{A}(\mathcal{C}_i \rightarrow \mathcal{C}_n) \end{aligned} \quad (21)$$

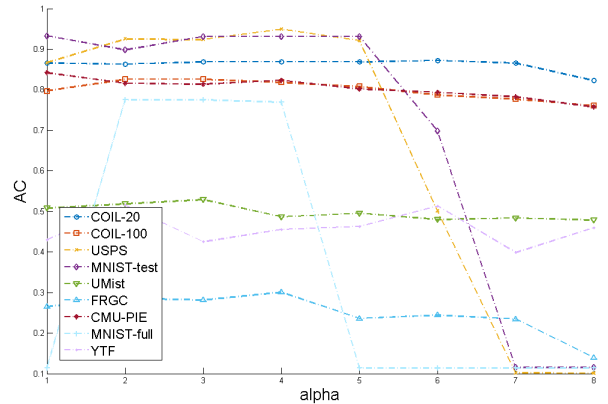
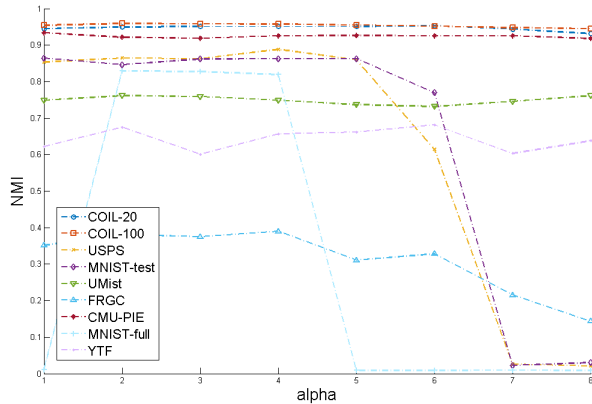


Figure 4: Performance of agglomerative clustering with approximations. Left one is NMI metric, and right one is AC metric. The first column is without acceleration. For the other columns from left to right, $\alpha = \{-0.2, -0.1, 0, 0.1, 0.2, 0.3, 0.5\}$.

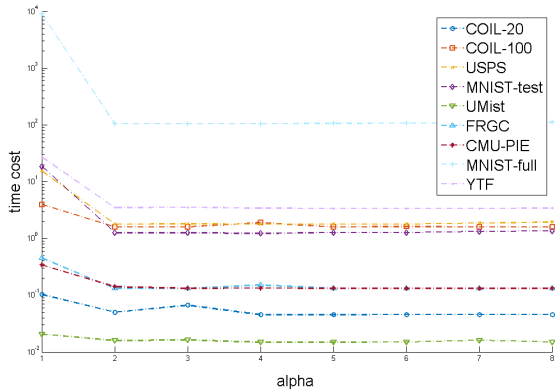


Figure 5: Time cost for different values of α . The first column is the time cost without acceleration. For the other columns from left to right, $\alpha = \{-0.2, -0.1, 0, 0.1, 0.2, 0.3, 0.5\}$.

We test various values for α , which are $\{-0.2, -0.1, 0, 0.1, 0.2, 0.3, 0.5\}$. We show the quantitative comparison in Fig. 4. We use image intensities as input to rule out all random factors. The original AC-GDL algorithm is used as the baseline. By conducting experiments on various datasets, we find a valid range $[-0.2, 0.1]$ for α which helps achieve analogous or even better performance to the one without acceleration. These results indicate that we may do not need to compute the explicit value of affinities to obtain equivalent level performance. Also, to measure how much time we can save by using our approximation, we compare the time cost between original AC-GDL algorithm and accelerated one in Fig. 5. It is clear that our approximation algorithm has much lower computational complexity.

A.3. Cluster-based to Sample-based Loss

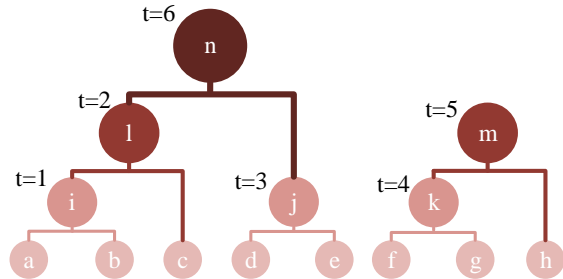


Figure 6: A illustration of agglomerative clustering.

In this part, we explain how to convert cluster-based loss to sample-based loss. Because it depends on specific agglomerative clustering processes, we use a toy example in Fig. 6 for illustration. We set K_c be 2 for simplicity. In Fig. 6, there are six time steps, and thus $T = 6$. We assume they are in a single partial unrolled period. The leaf nodes represent single samples. For simplicity, we omit $\frac{\lambda}{K_c - 1}$ in (10), obtaining the overall loss

$$\mathcal{L}(\theta|\mathcal{Y}_*, I) = - \sum_{t=1}^6 \left(\lambda' \mathcal{A}(C_*^t, \mathcal{N}_{C_*^t}^{K_c}[1]) - \mathcal{A}(C_*^t, \mathcal{N}_{C_*^t}^{K_c}[2]) \right) \quad (22)$$

Given above loss function, we decompose it from first time step ($t = 1$) to the most recent time step ($t = 6$):

- **t=1:** $C_*^1 = C_a, \mathcal{N}_{C_*^1}^2[1] = C_b$ and $\mathcal{N}_{C_*^1}^2[2] = C_c$. We have

$$\mathcal{L}(\theta|\mathcal{y}_*^1, I) = - (\lambda' \mathcal{A}(C_a, C_b) - \mathcal{A}(C_a, C_c)) \quad (23)$$

Clearly, above is sample-based weighted triplet loss

function, where samples C_a and C_b are positive pair and C_a and C_c are negative pair.

- **t=2:** $C_*^2 = C_i$, $\mathcal{N}_{C_*^2}^2[1] = C_c$ and $\mathcal{N}_{C_*^2}^2[2] = C_d$. We have

$$\begin{aligned} \mathcal{L}(\theta|\{\mathbf{y}_*^1, \mathbf{y}_*^2\}, I) &= \mathcal{L}(\theta|\mathbf{y}_*^1, I) \\ &\quad - (\lambda' \mathcal{A}(C_i, C_c) - \mathcal{A}(C_i, C_d)) \end{aligned} \quad (24)$$

Since $C_i = C_a \cup C_b$, we base on Eq. (19) for approximation

$$\begin{aligned} \mathcal{A}(C_i, C_c) &= \mathcal{A}(C_a \rightarrow C_c) + \mathcal{A}(C_b \rightarrow C_c) \\ &\quad + \frac{1}{2} \mathcal{A}(C_c \rightarrow C_a) + \frac{1}{2} \mathcal{A}(C_c \rightarrow C_b) \end{aligned} \quad (25)$$

$$\begin{aligned} \mathcal{A}(C_i, C_d) &= \mathcal{A}(C_a \rightarrow C_d) + \mathcal{A}(C_b \rightarrow C_d) \\ &\quad + \frac{1}{2} \mathcal{A}(C_d \rightarrow C_a) + \frac{1}{2} \mathcal{A}(C_d \rightarrow C_b) \end{aligned} \quad (26)$$

Thus,

$$\begin{aligned} \mathcal{L}(\theta|\{\mathbf{y}_*^1, \mathbf{y}_*^2\}, I) &= -\lambda' \mathcal{A}(C_a, C_b) - (\lambda' - 1) \mathcal{A}(C_a \rightarrow C_c) - \lambda' \mathcal{A}(C_b \rightarrow C_c) \\ &\quad - \left(\frac{\lambda'}{2} - 1\right) \mathcal{A}(C_c \rightarrow C_a) - \frac{\lambda'}{2} \mathcal{A}(C_c \rightarrow C_b) \\ &\quad + \mathcal{A}(C_a \rightarrow C_d) + \mathcal{A}(C_b \rightarrow C_d) \\ &\quad + \frac{1}{2} \mathcal{A}(C_d \rightarrow C_a) + \frac{1}{2} \mathcal{A}(C_d \rightarrow C_b) \end{aligned} \quad (27)$$

At current time step, sample a , b and c belong to the same cluster C_l , while sample d is from another cluster. (27) computes the sample-based weighted triplet loss for samples in C_l and sample d . Except for C_l , the other clusters all have merely one sample. No need to compute triplet loss for them. It should be pointed out that λ' in above loss function should be not less than 2 so that the affinities for all pairs in C_l are enlarged.

- **t=3:** $C_*^3 = C_d$, $\mathcal{N}_{C_*^3}^2[1] = C_e$ and $\mathcal{N}_{C_*^3}^2[2] = C_f$. We have

$$\begin{aligned} \mathcal{L}(\theta|\{\mathbf{y}_*^1, \mathbf{y}_*^2, \mathbf{y}_*^3\}, I) &= \mathcal{L}(\theta|\{\mathbf{y}_*^1, \mathbf{y}_*^2\}, I) \\ &\quad - \lambda' (\mathcal{A}(C_d, C_e) - \mathcal{A}(C_d, C_f)) \end{aligned} \quad (28)$$

Besides the loss $\mathcal{L}(\theta|\{\mathbf{y}_*^1, \mathbf{y}_*^2\}, I)$ for C_l , we also compute the loss for C_j in (28) because it contains two samples, d and e .

- **t=4:** $C_*^4 = C_f$, $\mathcal{N}_{C_*^4}^2[1] = C_g$ and $\mathcal{N}_{C_*^4}^2[2] = C_h$. We have

$$\begin{aligned} \mathcal{L}(\theta|\{\mathbf{y}_*^1, \dots, \mathbf{y}_*^4\}, I) &= \mathcal{L}(\theta|\{\mathbf{y}_*^1, \mathbf{y}_*^2, \mathbf{y}_*^3\}, I) \\ &\quad - (\lambda' \mathcal{A}(C_f, C_g) - \mathcal{A}(C_f, C_h)) \end{aligned} \quad (29)$$

Here, we additionally compute the weighted triplet loss for cluster C_k since it contains two samples.

- **t=5:** $C_*^5 = C_k$, $\mathcal{N}_{C_*^5}^2[1] = C_h$ and $\mathcal{N}_{C_*^5}^2[2] = C_j$. We have

$$\begin{aligned} \mathcal{L}(\theta|\{\mathbf{y}_*^1, \dots, \mathbf{y}_*^5\}, I) &= \mathcal{L}(\theta|\{\mathbf{y}_*^1, \dots, \mathbf{y}_*^4\}, I) - (\lambda' \mathcal{A}(C_k, C_h) - \mathcal{A}(C_k, C_j)) \end{aligned} \quad (30)$$

Because $C_k = C_f \cup C_g$, we have

$$\begin{aligned} \mathcal{A}(C_k, C_h) &= \mathcal{A}(C_f \rightarrow C_h) + \mathcal{A}(C_g \rightarrow C_h) \\ &\quad + \frac{1}{2} \mathcal{A}(C_h \rightarrow C_f) + \frac{1}{2} \mathcal{A}(C_h \rightarrow C_g) \end{aligned} \quad (31)$$

$$\begin{aligned} \mathcal{A}(C_k, C_j) &= \mathcal{A}(C_f \rightarrow C_j) + \mathcal{A}(C_g \rightarrow C_j) \\ &\quad + \frac{1}{2} \mathcal{A}(C_j \rightarrow C_f) + \frac{1}{2} \mathcal{A}(C_j \rightarrow C_g) \end{aligned} \quad (32)$$

Since $C_j = C_d \cup C_e$, we further transform above equation to

$$\begin{aligned} \mathcal{A}(C_k, C_j) &= \frac{1}{2} \mathcal{A}(C_f \rightarrow C_d) + \frac{1}{2} \mathcal{A}(C_f \rightarrow C_e) \\ &\quad + \frac{1}{2} \mathcal{A}(C_g \rightarrow C_d) + \frac{1}{2} \mathcal{A}(C_g \rightarrow C_e) \\ &\quad + \frac{1}{2} \mathcal{A}(C_d \rightarrow C_f) + \frac{1}{2} \mathcal{A}(C_e \rightarrow C_f) \\ &\quad + \frac{1}{2} \mathcal{A}(C_d \rightarrow C_g) + \frac{1}{2} \mathcal{A}(C_e \rightarrow C_g) \end{aligned} \quad (33)$$

Similar to the relation between sample a and c at time steps $t = 1, 2$, sample f and h belong to the same cluster C_m at current time step while they are from different clusters at time step $t = 4$. Based on the approximation, the terms $\mathcal{A}(C_f \rightarrow C_h)$ and $\mathcal{A}(C_h \rightarrow C_f)$ in two time steps will be merged. As a result, the final loss is computed on intra-cluster pairs and inter-cluster pairs sampled from three clusters C_l, C_j and C_m .

- **t=6:** $C_*^6 = C_l$, $\mathcal{N}_{C_*^6}^2[1] = C_j$ and $\mathcal{N}_{C_*^6}^2[2] = C_m$. Thus

$$\begin{aligned} \mathcal{L}(\theta|\{\mathbf{y}_*^1, \dots, \mathbf{y}_*^6\}, I) &= \mathcal{L}(\theta|\{\mathbf{y}_*^1, \dots, \mathbf{y}_*^5\}, I) \\ &\quad - (\lambda' \mathcal{A}(C_l, C_j) - \mathcal{A}(C_l, C_m)) \end{aligned} \quad (34)$$

Similar to the decomposition procedures above, both $\mathcal{A}(C_l, C_j)$ and $\mathcal{A}(C_l, C_m)$ can be transformed to sample-based affinities. Because C_l and C_j are regarded as different clusters previously, sample pairs from both of them are with positive weights in the loss function. However, it will be diminished by positive pairs (with negative weights) at current time step.

Though we use a toy example to show that the cluster-based loss can be transformed to sample-based loss above, the reduction is general to any possible agglomerative clustering processes because the loss for clusters at high-level

can always be decomposed to the losses on clusters at low-level until it reaches to single samples. The difference among various processes lies on the different weights associated with sample-based affinities. We should know that sample pairs from the same cluster may be with positive weights. One way to avoid this is increase λ' . In our implementation, we aim to increase affinities between samples from the same clusters, while decrease the affinities between samples from different clusters. And the clusters are determined by cluster ids at current step. Therefore, we assign a consistent weight γ to any affinities from the same cluster and 1 to any affinities from different clusters. Because we use SGD for batch optimization, the scales for affinities do not affect much on the performance. It is the signs affect much. Accordingly, at any given time step T , the overall loss is approximated to

$$\mathcal{L}(\theta|\mathbf{y}_*^T, \mathbf{I}) = -\frac{\lambda}{K_c - 1} \sum_{i,j,k} (\gamma \mathcal{A}(\mathbf{x}_i, \mathbf{x}_j) - \mathcal{A}(\mathbf{x}_i, \mathbf{x}_k)) \quad (35)$$

Note that we replace \mathbf{y}_* in (35) by \mathbf{y}_*^T in (35) because it is merely determined by current \mathbf{y}^T , regardless of $\{\mathbf{y}_*^1, \dots, \mathbf{y}_*^{T-1}\}$. As a result, we do not need to record $\{\mathbf{y}_*^1, \dots, \mathbf{y}_*^{T-1}\}$. This simplifies the batch optimization for CNN. Concretely, given a sample \mathbf{x}_i , we randomly select a sample \mathbf{x}_j which belongs to the same cluster, while select neighbours of \mathbf{x}_i that from other clusters to be \mathbf{x}_k . To omit the case that $\mathcal{A}(\mathbf{x}_i, \mathbf{x}_j)$ is much larger than $\mathcal{A}(\mathbf{x}_i, \mathbf{x}_k)$, we also add a margin threshold like the triplet loss function used in [59, 50].

A.4. Detailed CNN Architectures in our Paper

In this paper, the CNN architectures vary from dataset to dataset. As we mentioned in the main paper, we stacked different number of layers for different datasets so that the size of most top layer response map is about 10×10 . In Table 9, we list the architectures for the datasets used in our paper. "conv" means convolutional layer. "bn" means batch normalization layer. "wt-loss" means weighted triplet loss layer. \checkmark means the layer is used, while $-$ means the layer is not used.

A.5. Performance Evaluated by Accuracy

In this section, we evaluate the performance of different algorithms based on clustering accuracy (AC) metric, as a supplement to the NMI metric used in our main paper. As we can see from table 10, the proposed method outperform other methods on all datasets, which has similar trend as evaluated using NMI. Meanwhile, according to table 11, all other clustering algorithms are boosted after using the learned representation as evaluated on AC. These results further prove the proposed method is superior to other clustering algorithms and also learns powerful deep

representations that generalize well across different clustering algorithms.

A.6. Robustness Analysis

We choose the two most important parameters: unfolding rate η and K_s for evaluating the robustness of our approach to variations in these parameters. In these experiments, we set all the other parameters except for the target one to default values listed in Table 2 in the main paper. As we can see from Fig. 7, when the unfolding rate increases, the performance is not affected much for most of the datasets. For K_s , the performance is stable when $K_s \leq 50$ for all datasets. It drops with larger values of K_s for a few datasets. Increasing K_s also result in similar degradation in the agglomerative clustering algorithms we compare to. This suggests that K_s should not be set to very large value in general.

A.7. Reliability Analysis

We evaluate the reliability by measuring the purity of samples at the beginning of our algorithm. Because we use agglomerative clustering, there are very few samples in each cluster at the beginning (average is about 4 in our experiments). Most samples in the same cluster tend to belong to the same category. Quantitatively, for each sample in a dataset, we count the number of samples (K_m) that belong to the same category within its K nearest neighbours, and then compute the precision K_m/K for it. In Fig. 8, we report the average precision across all samples. As we can see, based on raw image data, all datasets have high ratios when K is smaller, and the ratios increase further when using our learned deep representations. Consequently, when K is small, the pseudo-labels are reliable enough to learn plausible deep representations.

A.8. Clustering based on Hand-Crafted Features

We also evaluate the performance of clustering based on image features, instead of image intensities. We choose three different types of datasets for testing: COIL100, MNIST-test and UMist, and three types of clustering algorithms including SC-LS [3], N-Cuts [52] and AC-PIC [69] for comparison since their better performance among all the algorithms. For these three datasets, we use spatial pyramid descriptor [32]⁴, histogram of oriented gradient (HOG) [7]⁵ and local binary pattern (LBP) [44] for representation, respectively. We report the results in Table 12. \downarrow means performance become worse, and \uparrow means it become better. Almost all algorithms perform worse than using original image as input. It indicates hand-crafted features

⁴<http://slazebni.cs.illinois.edu/research/SpatialPyramid.zip>

⁵<http://www.robots.ox.ac.uk/~vgg/research/caltech/phog.html>

Table 9: CNN architectures for different datasets in our paper.

Dataset	<i>COIL20</i>	<i>COIL100</i>	<i>USPS</i>	<i>MNIST-test</i>	<i>MNIST-full</i>	<i>UMist</i>	<i>FRGC</i>	<i>CMU-PIE</i>	<i>YTF</i>
conv1	✓	✓	✓	✓	✓	✓	✓	✓	✓
bn1	✓	✓	✓	✓	✓	✓	✓	✓	✓
relu1	✓	✓	✓	✓	✓	✓	✓	✓	✓
pool1	✓	✓	✓	✓	✓	✓	✓	✓	✓
conv2	✓	✓	–	✓	✓	✓	✓	✓	✓
bn2	✓	✓	–	✓	✓	✓	✓	✓	✓
relu2	✓	✓	–	✓	✓	✓	✓	✓	✓
pool2	✓	✓	–	–	–	✓	✓	✓	✓
conv3	✓	✓	–	–	–	✓	–	–	–
bn3	✓	✓	–	–	–	✓	–	–	–
relu3	✓	✓	–	–	–	✓	–	–	–
pool3	✓	✓	–	–	–	✓	–	–	–
conv4	✓	✓	–	–	–	–	–	–	–
bn4	✓	✓	–	–	–	–	–	–	–
relu4	✓	✓	–	–	–	–	–	–	–
pool4	✓	✓	–	–	–	–	–	–	–
ip1	✓	✓	✓	✓	✓	✓	✓	✓	✓
l2-norm	✓	✓	✓	✓	✓	✓	✓	✓	✓
wt-loss	✓	✓	✓	✓	✓	✓	✓	✓	✓

Table 10: Quantitative clustering performance (AC) for different algorithms using image intensities as input.

Dataset	<i>COIL20</i>	<i>COIL100</i>	<i>USPS</i>	<i>MNIST-test</i>	<i>MNIST-full</i>	<i>UMist</i>	<i>FRGC</i>	<i>CMU-PIE</i>	<i>YTF</i>
K-means [39]	0.665	0.580	0.467	0.560	0.564	0.419	0.327	0.246	0.548
SC-NJW [43]	0.641	0.544	0.413	0.220	0.502	0.551	0.178	0.255	0.551
SC-ST [67]	0.417	0.300	0.308	0.454	0.311	0.411	0.358	0.293	0.290
SC-LS [3]	0.717	0.609	0.659	0.740	0.714	0.568	0.407	0.549	0.544
N-Cuts [52]	0.544	0.577	0.314	0.304	0.327	0.550	0.235	0.155	0.536
AC-Link [25]	0.251	0.269	0.421	0.693	0.657	0.398	0.175	0.201	0.547
AC-Zell [70]	0.867	0.811	0.575	0.693	0.112	0.517	0.266	0.765	0.519
AC-GDL [68]	0.865	0.797	0.867	0.933	0.113	0.563	0.266	0.842	0.430
AC-PIC [69]	0.855	0.840	0.855	0.920	0.115	0.576	0.320	0.797	0.472
NMF-LP [1]	0.621	0.553	0.522	0.479	0.471	0.365	0.259	0.229	0.546
OURS-SF	1.000	0.894	0.922	0.940	0.959	0.809	0.461	0.980	0.684
OURS-RC	1.000	0.916	0.950	0.961	0.964	0.809	0.461	1.000	0.684

Table 11: Quantitative clustering performance (AC) for different algorithms using our learned representations as inputs.

Dataset	<i>COIL20</i>	<i>COIL100</i>	<i>USPS</i>	<i>MNIST-test</i>	<i>MNIST-full</i>	<i>UMist</i>	<i>FRGC</i>	<i>CMU-PIE</i>	<i>YTF</i>
K-means [39]	0.821	0.751	0.776	0.957	0.969	0.761	0.476	0.834	0.660
SC-NJW [43]	0.738	0.659	0.716	0.868	0.972	0.707	0.485	0.776	0.521
SC-ST [67]	0.851	0.705	0.661	0.960	0.958	0.697	0.496	0.896	0.575
SC-LS [3]	0.867	0.735	0.792	0.960	0.973	0.733	0.502	0.802	0.571
N-Cuts [52]	0.888	0.626	0.634	0.959	0.971	0.798	0.504	0.981	0.441
AC-Link [25]	0.678	0.539	0.773	0.955	0.964	0.795	0.495	0.947	0.602
AC-Zell [70]	1.000	0.931	0.879	0.879	0.969	0.790	0.449	1.000	0.644
AC-GDL [68]	1.000	0.920	0.949	0.961	0.878	0.790	0.461	1.000	0.677
AC-PIC [69]	1.000	0.950	0.955	0.958	0.882	0.790	0.438	1.000	0.652
NMF-LP [1]	0.769	0.603	0.778	0.955	0.970	0.725	0.481	0.504	0.575

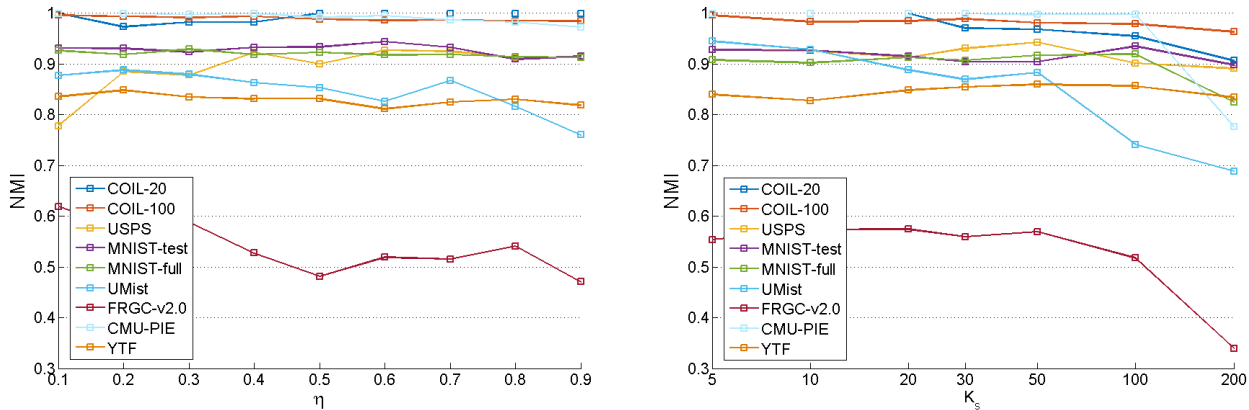


Figure 7: Clustering performance (NMI) with different η (left) and K_s (right).

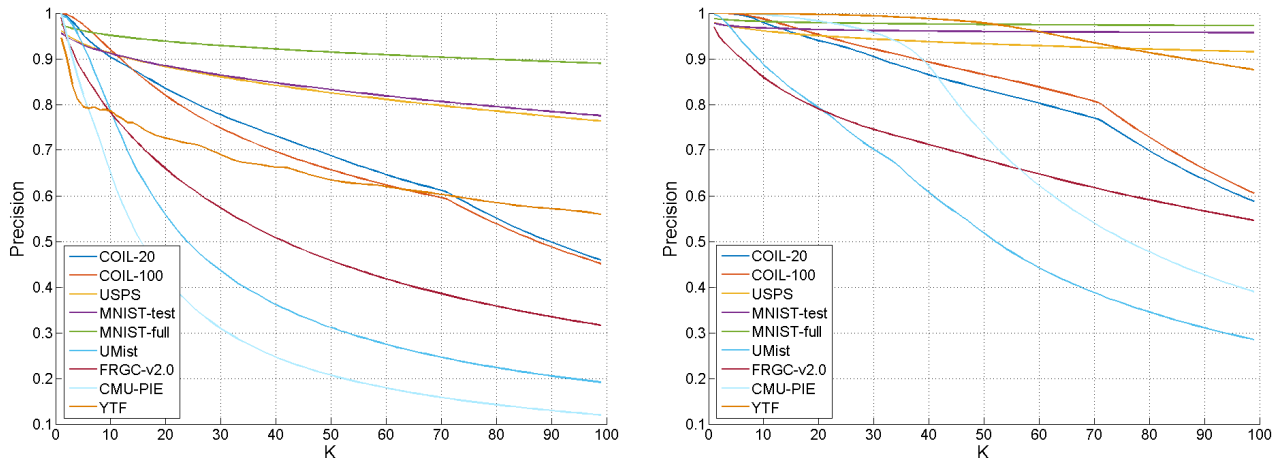


Figure 8: Average purity of K -nearest neighbour for varying values of K . Left is computed using raw image data, while right is computed using our learned representation.

Table 12: Clustering performance (NMI) based on hand-crafted features.

Dataset	COIL100	MNIST-test	UMist	FRGC
SC-LS [3]	0.733↓	0.625↓	0.752↓	0.338↓
N-Cuts [52]	0.722↓	0.423↑	0.420↓	0.238↓
AC-PIC [69]	0.878↓	0.735↓	0.734↓	0.322↓

should be designed dataset by dataset. In contrast, directly learning from image intensities is more straightforward and also achieves better performance.

A.9. Visualizing Data in Low Dimension

Projecting high-dimensional data into low-dimensional space can help people to intuitively understand the data. Though the proposed method is aimed to learn deep representations and image clusters, we note that it can be naturally converted to a parametric visualization method for an image dataset by slightly alternating the objective. Instead of updating the affinities among samples based on the learned representations gradually, we consistently use the affinities among raw image data to perform the agglomerative cluster, which then guides representation learning in low-dimensional space. By this way, we can obtain a low-dimensional space (2D or 3D) which can retain the structure of the original data.

We compare three dimension reduction techniques, principle component analysis (PCA) [62], neighbourhood com-

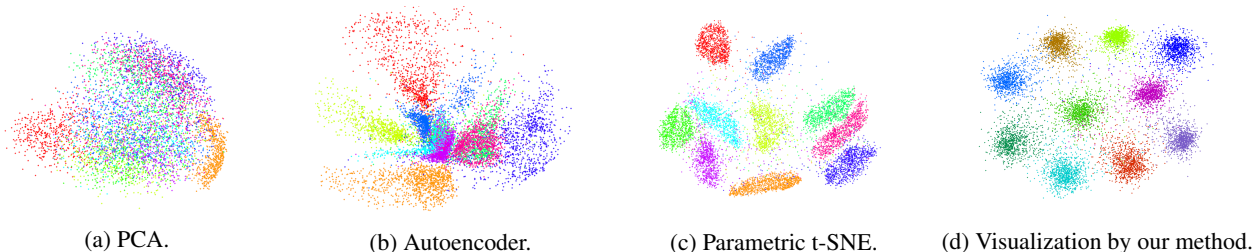


Figure 9: Visualization of 10,000 MNIST test samples in different embedding spaces.

ponents analysis (NCA) [48], and parametric t-SNE [38]. Though both [38] and our visualization method are based on neural networks, there are two main differences: 1) In [38], a Kullback-Leibler divergence between the joint distributions of original data and the embedded data is considered. However, in our method, we employ a weighted triplet loss that directly takes the local structure of embedded data into account; 2) In [38], the authors need to pre-train a stack of RBMs layer-by-layer, and then fine-tune the neural network. Nevertheless, we directly train the neural network from scratch end-to-end.

We perform experiments on MNIST dataset. In MNIST, 60,000 training samples are used to learn the low-dimensional embedding space, and 10,000 test samples are used for evaluation. To train a D -dimensional embedding, we first remove the normalization layer and then stack on the top another linear layer whose dimension is D . To thoroughly explore the local structure in the original data, we merge the clusters with a lower unfolding rate ($\eta = 0.2$). The learning process is stopped when the number of clusters reaches to 10. Though we stop the learning process as such, it should be noted that the stop criterion is not confined. For quantitative analysis, we compute the nearest-neighbor classification error and trustworthiness as in [38].

In Table 13, we show the 1-nearest neighbour classification error on MNIST test dataset. We copy the best results of the compared methods from [38]. As we can see, our method outperforms all three other methods across three different embedding dimensions. These results illustrates that our method can obtain low-dimensional embedding with better generalization ability.

For visualization, it is important to retain the original data structure in the embedding space. For quantitative comparison, we report the trustworthiness of learned low-dimensional embedding in Table 14. Larger value means better preservation of original data structure. As we can see, our method is not as good as parametric t-SNE. These results are explainable. During training, we merely pay attention to the local structure among samples from different clusters, while omitting the relations among samples within

Table 13: 1-nearest neighbor classification error on low-dimensional embedding of MNIST dataset.

Method	2D	10D	30D
PCA [62]	0.782	0.430	0.108
NCA [48]	0.568	0.088	0.073
Autoencoder [22]	0.668	0.063	0.027
Param. t-SNE [38]	0.099	0.046	0.027
OURS	0.067	0.019	0.027

Table 14: Trustworthiness T(12) on low-dimensional embedding of MNIST dataset.

Method	2D	10D	30D
PCA [62]	0.744	0.991	0.998
NCA [48]	0.721	0.968	0.971
Autoencoder [22]	0.729	0.996	0.999
Param. t-SNE [38]	0.927	0.997	0.999
Ours	0.768	0.936	0.975

one cluster. Therefore, the algorithm will learn embeddings that discriminate clusters well but possibly disorder the samples in each cluster. We believe this can be solved by introducing a loss to confine the within-cluster structure. We leave this as a future work for limited space.

A.10. Visualizing Learned Deep Representations

We show the first three principle components of learned representations in Fig. 10 and Fig. 11 at different stages. For comparison, we show the image intensities at the first column. We use different colors for representing different clusters that we predict during the algorithm. At the bottom of each plot, we give the number of clusters at the corresponding stage. At the final stage, the number of cluster is same to the number of categories in the dataset. After a number of iterations, we can learn more discriminative representations for the datasets, and thus facilitate more precise clustering results.

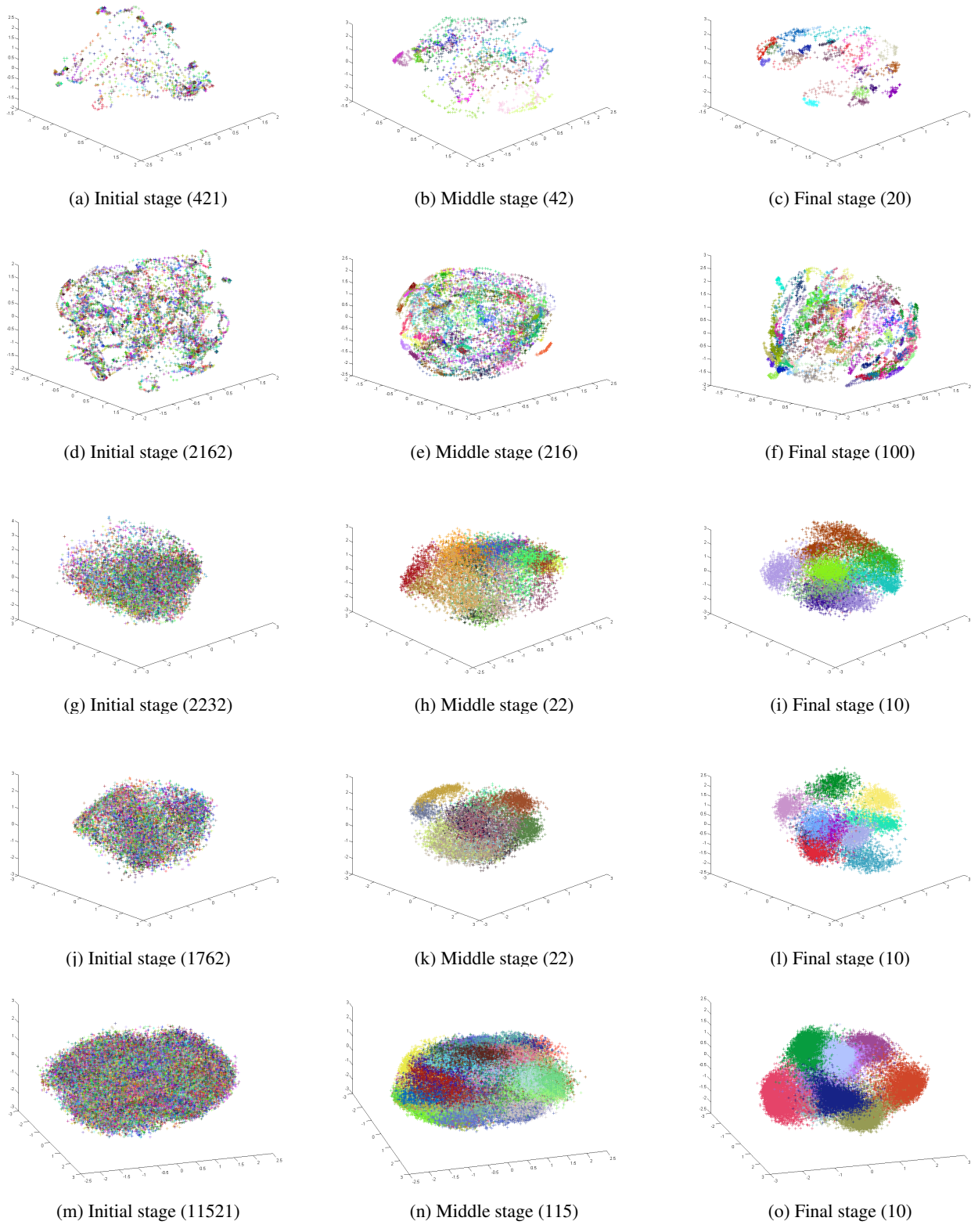


Figure 10: Learned representations at different stages on five datasets. From top to bottom, they are *COIL20*, *COIL100*, *USPS* and *MNIST-test* and *MNIST-full*. The first column are image intensities. For *MNIST-test*, we show another view point different from Fig.1 in the main paper.

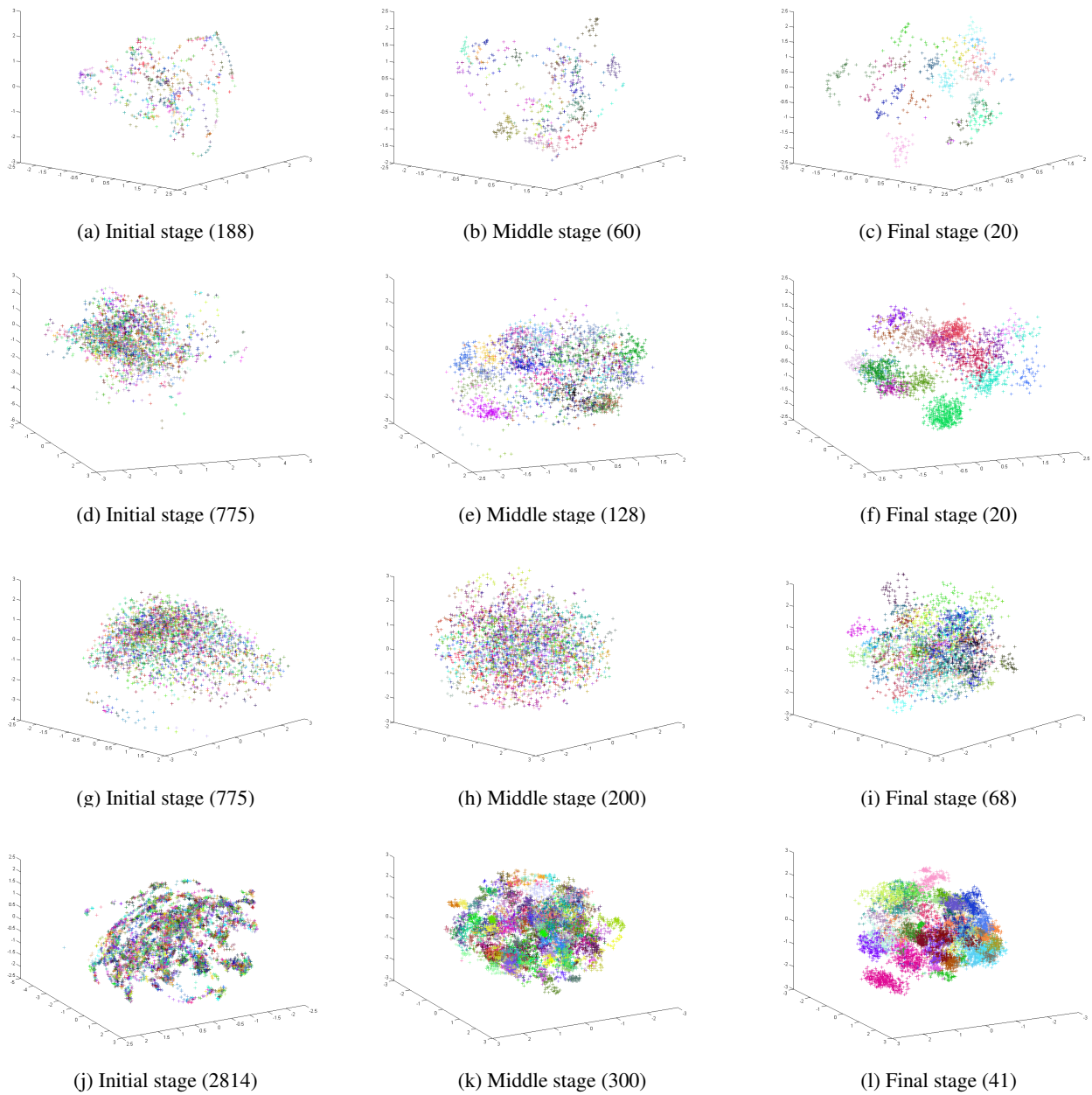


Figure 11: Learned representations as different stages on four datasets. From top to bottom, they are *UMist*, *FRGC*, *CMU-PIE* and *YTF*. The first column are image intensities.