# Local Aggregation for Unsupervised Learning of Visual Embeddings

Chengxu Zhuang
Stanford University

Alex Lin Zhai
Stanford University

Daniel Yamins
Stanford University

## Abstract

*Unsupervised approaches to learning in neural networks are of substantial interest for furthering artificial intelligence, both because they would enable the training of networks without the need for large numbers of expensive annotations, and because they would be better models of the kind of general-purpose learning deployed by humans. However, unsupervised networks have long lagged behind the performance of their supervised counterparts, especially in the domain of large-scale visual recognition. Recent developments in training deep convolutional embeddings to maximize non-parametric instance separation and clustering objectives have shown promise in closing this gap. Here, we describe a method that trains an embedding function to maximize a metric of local aggregation, causing similar data instances to move together in the embedding space, while allowing dissimilar instances to separate. This aggregation metric is dynamic, allowing soft clusters of different scales to emerge. We evaluate our procedure on several large-scale visual recognition datasets, achieving state-of-the-art unsupervised transfer learning performance on object recognition in ImageNet, scene recognition in Places 205, and object detection in PASCAL VOC.*

## 1. Introduction

Deep convolutional neural networks (DCNNs) have achieved great success on many tasks across a variety of domains, such as vision [37, 60, 24, 23, 7], audition [26, 21, 11, 47], and natural language processing [68, 29, 9, 38]. However, most successful DCNNs are trained in a supervised fashion on labelled datasets [37, 60, 24, 10, 26], requiring the costly collection of large numbers of annotations. There is thus substantial interest in finding methods that can train DCNNs solely using unlabeled data, which are often readily available. Over many decades of work, substantial progress has been achieved using a wide variety of unsupervised learning approaches [6, 67, 71, 36, 13, 15, 70, 48, 63, 49]. Nevertheless, unsupervised networks are still typically significantly lower performing than their supervised counterparts, and are rarely used in real-world

applications [6, 47, 7].

In contrast to the inefficiency of unsupervised learning in artificial neural networks, humans and non-human primates develop powerful and domain-general visual systems with very few labels [42, 5, 65, 1, 22, 4, 61]. Although the mechanisms underlying the efficiency of biological learning still remain largely unknown [5], researchers reliably report that infants as young as three months can group perceptually similar stimuli [46], even for stimulus types that the infants have never seen before. Moreover, this ability arises long before these infants appear to have an explicit concept of object category [46, 54, 31, 8]. These findings suggest that biological unsupervised learning may take advantage of inherent visual similarity, without requiring sharp boundaries between stimulus categories.

Inspired by these results, we propose a novel unsupervised learning algorithm through local non-parametric aggregation in a latent feature space. First, we non-linearly embed inputs in a lower-dimensional space via a neural network. We then iteratively identify close neighbors surrounding each example in the embedding space, while optimizing the embedding function to strengthen the degree of local aggregation. Our procedure, which we term Local Aggregation (LA), causes inputs that are naturally dissimilar to each other to move apart in the embedding space, while allowing inputs that share statistical similarities to arrange themselves into emergent clusters. By simultaneously optimizing this soft clustering structure and the non-linear embedding in which it is performed, our procedure exposes subtle statistical regularities in the data. The resulting representation in turn robustly supports downstream tasks.

Here, we illustrate the LA procedure in the context of large-scale visual learning. Training a standard convolution neural network with LA using images from ImageNet [10] significantly outperforms current state-of-art unsupervised algorithms on transfer learning to classification tasks on both ImageNet and the Places 205 dataset [72]. In addition, LA shows consistent improvements as the depth of the embedding function increases, allowing it to achieve $60.2\%$ top-1 accuracy on ImageNet classification. This is, as far as we know, the first time an unsupervised model has surpassed the milestone AlexNet network trained directly on

the supervised task. We also show that, through further fine-tuning, LA trained models obtain state-of-the-art results on the PASCAL object detection task.

The remainder of this paper is organized as follows: in section 2, we discuss related work; in section 3, we describe the LA method; in section 4, we show experimental results; in section 5, we present analyses illustrating how this algorithm learns and justifying key parameter choices.

## 2. Related Work

Unsupervised learning methods span a very broad spectrum of approaches going back to the roots of artificial neural networks [53, 39, 2, 58, 27, 41, 28, 30, 25], and are too numerous to fully review here. However, several recent works have achieved exciting progress in unsupervised representation learning [6, 67, 71]. Although the LA method draws inspiration from these works, it differs from them in some important conceptual ways.

**DeepCluster.** DeepCluster [6] (DC) trains a DCNN in a series of iterative rounds. In each round, features from the penultimate layer of the DCNN from the previous round are clustered, and the cluster assignments are used as self-generated supervision labels for further training the DCNN using standard error backprogation. Like DC, LA also uses an iterative training procedure, but the specific process within each iteration differs significantly. First, unlike the clustering step of DC where all examples are divided into mutually-exclusive clusters, our method identifies neighbors separately for each example, allowing for more flexible statistical structures than a partition. Indeed, as shown in Section 5.2, the use of individual semantic neighbor identifiers rather than global clustering is important for performance improvement. Secondly, the optimization step of LA differs from that of DC by optimizing a different objective function. Specifically, DC optimizes the cross-entropy loss between predicted and ground truth cluster labels, requiring an additional and computationally expensive linear readout layer. Moreover, due to arbitrary changes in the cluster label indices across iterative rounds, this additional readout layer needs to be frequently recomputed. In contrast, LA employs an objective function that directly optimizes a local soft-clustering metric, requiring no extra readout layer and only a small amount of additional computation on top of the feature representation training itself. These differences lead both to better final performance and substantially improved training efficiency.

**Instance Recognition.** The Instance Recognition task [67] (IR) treats each example as its own "category" and optimizes the DCNN representation to output an embedding in which all examples are well-separated from each other. LA uses a similar embedding framework, but achieves significantly better performance by pursuing a distinct optimization goal. Specifically, while IR optimizes for equally

separating representations of all examples, LA encourages a balance between separation and clustering on a per-example basis, as measured by the local aggregation criterion. For this reason, the LA approach can be thought of as a principled hybrid between the DC and IR approaches.

**Self-supervised "missing-data" tasks.** These tasks build representations by hiding some information about each example input, and then optimizing the network to predict the hidden information from the visible information that remains. Examples include context prediction [13], colorization of grayscale images [13], inpainting of missing portions of images [52], and the Split-Brain method [71]. However, it is ultimately unclear whether these tasks are perfectly aligned with the needs of robust visual representation. Indeed, it has been found that deeper networks better minimizing the loss functions used in such tasks gain little transfer learning performance on object recognition tasks [14]. Moreover, most missing-data tasks rely on structures that are specific to visual data, making them potentially less general than the embedding/clustering concepts used in DC, IR or our LA method.

**Generative models**. Another broad class of unsupervised learning algorithm, often termed deep generative models, focuses on reconstructing input images from a bottlenecked latent representation. The networks trained by these algorithms use the latent representations for other tasks, including object recognition. These learning methods include classical ones such as Restricted Boltzman Machines [27, 41] as well as more recent ones such as Variational Auto-Encoders [35] and Generative Adversarial Networks [15, 20]. Although the features learned by generative models have been put to a wide variety of exciting uses [40, 69, 12, 43, 33], their power as latent representations for downstream visual tasks such as object recognition has yet to be fully realized.

## 3. Methods

Our overall objective is to learn an embedding function $f_\theta$ (realized via a neural network) that maps images $\mathbf{I} = \{\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_N\}$ to features $\mathbf{V} = \{\mathbf{v}_1, \mathbf{v}_2, ..., \mathbf{v}_N\}$ with $\mathbf{v}_i = f_\theta(\mathbf{x}_i)$ in a compact $D$-dimension representation space where similar images are clustered while dissimilar images are separated. To achieve this objective, we design an iterative procedure to bootstrap the aggregation power of a deep non-linear embedding function. More specifically, at any given stage during training the embedding function, we dynamically identify two sets of neighbors for an $\mathbf{x}_i$ and its embedding $\mathbf{v}_i$: *close neighbors* $\mathbf{C}_i$ and *background neighbors* $\mathbf{B}_i$. Intuitively, close neighbors are those whose embeddings should be made similar to $\mathbf{v}_i$, while background neighbors are used to set the distance scale with respect to which the judgement of closeness should be measured. To help better understand these two sets, we provide a

schematic illustration in Fig. 1, and describe the details of how they are defined mathematically in section 3.1. Using $\mathbf{B}_i$ and $\mathbf{C}_i$, we then define the level of *local aggregation* $L(\mathbf{C}_i, \mathbf{B}_i | \boldsymbol{\theta}, \mathbf{x}_i)$ near each input $\mathbf{x}_i$, which characterizes the relative level of closeness within $\mathbf{C}_i$, compared to that in $\mathbf{B}_i$. The parameters $\boldsymbol{\theta}$ of the neural network realizing the embedding function are then tuned over the course of training to maximize $L(\mathbf{C}_i, \mathbf{B}_i | \boldsymbol{\theta}, \mathbf{x}_i)$.

### 3.1. Neighbor Identification

We first describe how the neighbor types $\mathbf{B}_i$ and $\mathbf{C}_i$ are defined. *Nearest-neighbor based identification for* $\mathbf{B}_i$: At any given step of optimization, the background neighbors for a given embedded point $\mathbf{v}_i$ are simply defined as the $k$ closest embedded points $\mathcal{N}_k(\mathbf{v}_i)$ within $\mathbf{V}$, where distance is judged using the cosine distance on the embedding space. The number $k$ of background neighbors to be used is a hyperparameter of the algorithm. *Robustified clustering-based identification for* $\mathbf{C}_i$: To identify close neighbors, we first apply an unsupervised clustering algorithm on all embedded points $\mathbf{V}$ to cluster the representations into $m$ groups $\mathbf{G} = \{G_1, G_2, ..., G_m\}$. Let $g(\mathbf{v}_i)$ denote the cluster label of $\mathbf{v}_i$ in this clustering result, i.e. $i \in G_{g(\mathbf{v}_i)}$. In the simplest version of our procedure, we then define $\mathbf{C}_i$ to be the set $G_{g(\mathbf{v}_i)}$. However, because clustering can be a noisy and somewhat arbitrary process, we compute multiple clusterings under slightly different conditions, and then aggregate neighbors across these multiple clusterings to achieve more stable results. Specifically, let $\{\mathbf{G}^{(j)}\}$ be clusters for $H$ distinct clusterings, where $\mathbf{G}^{(j)} = \{G_1^{(j)}, G_2^{(j)}, ..., G_{m^{(j)}}^{(j)}\}$ with $j \in \{1, 2, ..., H\}$, and $\{g^{(j)}\}$ defined accordingly. We then define $\mathbf{C}_i = \bigcup_{j=1}^{H} G_{g^{(j)}(\mathbf{v}_i)}^{(j)}$. The number $m$ of clusters and number $H$ of clusterings are hyperparameters of the algorithm. In this work, we use $k$-means clustering as the standard unsupervised algorithm.

Intuitively, background neighbors are an unbiased sample of nearby points that (dynamically) set the scale at which "close-ness" should be judged; while close neighbors are those that are especially nearby, relative to those in other clusters. The mathematical definitions above represent just one specific way to formalize these ideas, and many alternatives are possible. In Section 5.2, we show that our choices are not arbitrary by exploring the consequences of making alternate decisions.

### 3.2. Local Aggregation Metric

Given the definition of $\mathbf{B}_i$ and $\mathbf{C}_i$, we describe the formulation of our local aggregation metric, $L(\mathbf{C}_i, \mathbf{B}_i | \boldsymbol{\theta}, \mathbf{x}_i)$. We build our formulation upon a non-parametric softmax operation proposed by Wu et al. in [67]. In that work, the authors define the probability that an arbitrary feature $\mathbf{v}$ is

recognized as the $i$-th image to be:

$$P(i|\mathbf{v}) = \frac{\exp(\mathbf{v}_i^T \mathbf{v}/\tau)}{\sum_{j=1}^{N} \exp(\mathbf{v}_j^T \mathbf{v}/\tau)} \quad (1)$$

where $\tau \in [0, 1]$ is a fixed scale hyperparameter, and where both $\{\mathbf{v}_i\}$ and $\mathbf{v}$ are projected onto the L2-unit sphere in the $D$-dimensional embedding space (e.g. normalized such that $\|\mathbf{v}\|_2 = 1$).

Following equation 1, given an image set $\mathbf{A}$, we then define the probability of feature $\mathbf{v}$ being recognized as an image in $\mathbf{A}$ as:

$$P(\mathbf{A}|\mathbf{v}) = \sum_{i \in \mathbf{A}} P(i|\mathbf{v}) \quad (2)$$

Finally, we formulate $L(\mathbf{C}_i, \mathbf{B}_i | \boldsymbol{\theta}, \mathbf{x}_i)$ as the negative log-likelihood of $\mathbf{v}_i$ being recognized as a close neighbor (e.g. is in $\mathbf{C}_i$), given that $\mathbf{v}_i$ is recognized as a background neighbor (e.g. is in $\mathbf{B}_i$):

$$L(\mathbf{C}_i, \mathbf{B}_i | \boldsymbol{\theta}, \mathbf{x}_i) = -\log \frac{P(\mathbf{C}_i \cap \mathbf{B}_i | \mathbf{v}_i)}{P(\mathbf{B}_i | \mathbf{v}_i)} \quad (3)$$

The loss to be minimized is then:

$$\mathcal{L}_i = L(\mathbf{C}_i, \mathbf{B}_i | \boldsymbol{\theta}, \mathbf{x}_i) + \lambda \|\boldsymbol{\theta}\|_2^2 \quad (4)$$

where $\lambda$ is a regularization hyperparameter.

**Discussion**. Because the definition of $L(\mathbf{C}_i, \mathbf{B}_i | \boldsymbol{\theta}, \mathbf{x}_i)$ is somewhat involved, we describe a simple conceptual analysis that illustrates the intuition for why we chose it as a measure of local aggregation. Letting $\mathbf{C}_i^c$ denote the complement of $\mathbf{C}_i$ in $\mathbf{I}$, we have $P(\mathbf{B}_i | \mathbf{v}_i) = P(\mathbf{C}_i^c \cap \mathbf{B}_i | \mathbf{v}_i) + P(\mathbf{C}_i \cap \mathbf{B}_i | \mathbf{v}_i)$. Thus, from equation 3, we see that $L(\mathbf{C}_i, \mathbf{B}_i | \boldsymbol{\theta}, \mathbf{x}_i)$ is minimized when $P(\mathbf{C}_i \cap \mathbf{B}_i | \mathbf{v}_i)$ is maximized and $P(\mathbf{C}_i^c \cap \mathbf{B}_i | \mathbf{v}_i)$ is minimized. It is easy to understand the meaning of minimizing $P(\mathbf{C}_i^c \cap \mathbf{B}_i | \mathbf{v}_i)$: this occurs as the distances between $\mathbf{v}_i$ and its non-close background neighbors are maximized. The consequences of maximizing $P(\mathbf{C}_i \cap \mathbf{B}_i | \mathbf{v}_i)$ are a bit more subtle. As shown empirically in [66] (albeit in the supervised context), as long as the scaling parameter $\tau \ll 1$, maximizing $P(\mathbf{A}|\mathbf{v}_i)$ for any set $\mathbf{A}$ causes the emergence of natural "sub-categories" in (the embeddings of) $\mathbf{A}$, and encourages $\mathbf{v}_i$ to move closer to *one* of these sub-categories rather than their overall average. This empirical result can be intuitively understood by recognizing the fact that $\exp(\mathbf{v}_i^T \mathbf{v}/\tau)$ increases exponentially when $\mathbf{v}_i^T \mathbf{v}$ approaches 1, suggesting that $P(\mathbf{A}|\mathbf{v}_i)$ will approach 1 when $\mathbf{A}$ includes a small cluster of features that are all very close to $\mathbf{v}$. Putting these observations together, the optimized representation space created by minimizing $L(\mathbf{C}_i, \mathbf{B}_i | \boldsymbol{\theta}, \mathbf{x}_i)$ is, intuitively, like that shown in Fig. 1: a set of embedded points that have formed into small clusters at a distribution of natural scales.
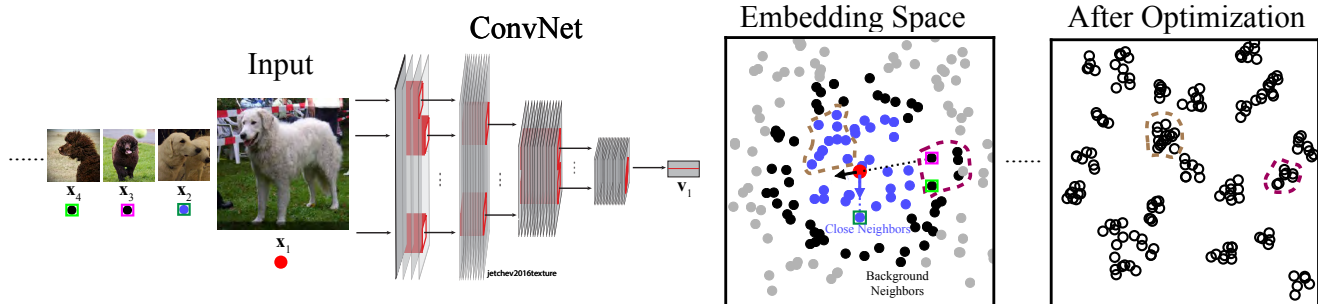
Figure 1. Illustration of the Local Aggregation (LA) method. For each input image, we use a deep neural network to embed it into a lower dimension space ("Embedding Space" panel). We then identify its close neighbors (blue dots) and background neighbors (black dots). The optimization seeks to push the current embedding vector (red dot) closer to its close neighbors and further from its background neighbors. The blue arrow and black arrow are examples of influences from different neighbors on the current embedding during optimization. The "After Optimization" panel illustrates the typical structure of the final embedding after training.

### 3.3. Memory Bank

As defined above, the neighbor identification procedures and the loss function implicitly describe computations involving all the embedded features $\mathbf{V}$, which soon becomes intractable for large datasets. To address this issue, we follow [67, 66] and maintain a running average for $\mathbf{V}$, which is called the *memory bank*, denoted $\bar{\mathbf{V}} = \{\bar{\mathbf{v}}_1, \bar{\mathbf{v}}_2, ..., \bar{\mathbf{v}}_N\}$. Similarly to [67, 66], we initialize the memory bank with random $D$-dimensional unit vectors and then update its values by mixing $\bar{\mathbf{v}}_i$ and $\mathbf{v}_i$ during training as follows:

$$\bar{\mathbf{v}}_i \leftarrow (1 - t)\bar{\mathbf{v}}_i + t\mathbf{v}_i \qquad (5)$$

where $t \in [0, 1]$ is a fixed mixing hyperparameter. With the help of $\bar{\mathbf{V}}$, we can then rewrite the neighbor identification procedures and equation 1 by replacing the feature sets $\mathbf{V}$ with $\bar{\mathbf{V}}$. In particular for $\mathbf{C}_i$, the cluster label function $g$ is applied to $\bar{\mathbf{v}}_i$ by index identification, ensuring the chosen cluster includes the index $i$ itself. After this replacement, it is no longer necessary to recompute $\mathbf{V}$ before every step to identify (good approximations of) $\mathbf{C}_i$ and $\mathbf{B}_i$.

## 4. Results

In this section, we describe tests of the LA method on visual representation learning and compare its performance to that of other methods.

### 4.1. Experiment Settings

We first list key parameters used for network training. Following [67], we set parameter $\tau = 0.07$, $D = 128$, $\lambda = 0.0001$, and $t = 0.5$. For all network structures, we use SGD with momentum of 0.9 and batch size 128. Initial learning rates are set to 0.03, and dropped by a factor of 10 when validation performances saturate, typically leading to training for 200 epochs with two learning rate drops. Most of these parameters are taken from [67], as our conceptual framework is similar, but a further hyper-parameter search might lead to better results, given that our optimization goal differs substantially.

As a warm start for our models, we begin training using the IR loss function for the first 10 epochs, before switching over to using the LA method. Following the methods of [6], for AlexNet [37] and VGG16 [60] architectures, we add batch normalization (BN) layers [32] after all convolution and fully-connected layers, before ReLu operations, to allow a higher learning rate and a faster convergence speed. Though adding BN is known to improve convergence speed but not typically to lead to higher final ImageNet performance levels using supervised training regimes, it is unclear whether this remains true when using unsupervised training methods. Importantly, the potentially competitive IR method [67] did not originally include BN in their AlexNet and VGG16, so to ensure that we have fairly compared that method to LA or DC, we also train AlexNet and VGG16 with BN on the IR task. For all structures, we replace the final category readout layer with a linear layer with $D$ output units, followed by a L2-normalization operation to ensure that the output is a unit vector.

We set $k = 4096$ for computing $\mathbf{B}_i$ using the nearest neighbors procedure. In computing $\mathbf{C}_i$, we use $k$-means [44] implemented in Faiss [34] as the standard unsupervised clustering algorithm, generating multiple clusterings for robustness via different random initializations. Using the notation of Section 3, AlexNet is trained with $H = 3, m = 30000$, VGG16 is trained with $H = 6, m = 10000$, all ResNet structures are trained with $H = 10, m = 30000$. We justify all parameter choices and intuitively explain why they are optimal in Section 5.2. All code for reproducing our training is available at: [WEBSITE_WITHHOLD].

## 4.2. Transcription Learning Results

After fully training networks on ImageNet, we then test the quality of the learned visual representations by evaluating transfer learning to other tasks, including ImageNet classification on held-out validation images, scene classification on Places205 [72], and object detection on PASCAL VOC 2007 [17]. For classification tasks, we also report K-nearest neighbor (KNN) classification results using the embedding features, acquired via a method similar to that in [67]. Specifically, we take top $K$ nearest neighbors $\mathcal{N}_K$ for the feature $\mathbf{v}$ either (for ImageNet) from the saved memory bank or (for Places) from the computed network outputs for center crops of training images. Their labels are then weighted by $\exp(\mathbf{v}_i^T \mathbf{v}/\tau)$ and combined to get final predictions. We report results for $K = 200$ as in [67].

**Object Recognition.** To evaluate transfer learning for the ImageNet classification task, we fix network weights learned during the unsupervised procedure, add a linear readout layer on top of each layer we want to evaluate, and train the readout using cross-entropy loss together with L2 weight decay. We use SGD with momentum of 0.9, batch size 128, and weight decay 0.0001. Learning rate is initialized at 0.01 and dropped by a factor of 10 when performance saturates, typically leading to 90 training epochs with two learning rate drops. We report 10-crop validation performances to ensure comparability with [6]. Performance results in Table 1 show that LA significantly outperforms other methods with all architectures, especially in deeper architectures. LA-trained AlexNet reaches 42.4%, which is 1.4% higher than previous state-of-the-art. Improvements over previous unsupervised state-of-the-art are substantially larger for VGG16 (+4.9%), ResNet-18 (+3.7%), and ResNet-50 (+6.2%). In particular, LA-trained ResNet-50 achieves **60.2**% top-1 accuracy on ImageNet classification, surpassing AlexNet trained directly on the supervised task. Using KNN classifiers, LA outperforms the IR task by a large margin with all architectures. There is a consistent performance increase for the LA method both from overall deeper architectures, and from earlier layers to deeper layers within an architecture. Most alternative training methods (e.g. [51, 45, 13, 70]) do not benefit significantly from increasing depth. For example, ResNet-101 trained using Color [70] can only achieve 39.6% and the best performance using ResNet-101 with unsupervised task is only 48.7% with CPC [50].

**Scene Categorization.** To test the generalization ability of the learned representations to a data distribution distinct from that used in training, we assessed transfer to the Places [72] dataset, which includes $2.45M$ images labelled with 205 scene categories. As in the previous section, we train linear readout layers for the scene categorization task on top of the pretrained ImageNet model, using training procedures and hyper-parameters identical to those used in

| Method | conv1 | conv2 | conv3 | conv4 | conv5 | KNN |
|---|---|---|---|---|---|---|
| AlexNet | | | | | | |
| Random | 11.6 | 17.1 | 16.9 | 16.3 | 14.1 | 3.5 |
| Context [13] | 16.2 | 23.3 | 30.2 | 31.7 | 29.6 | – |
| Color [70] | 13.1 | 24.8 | 31.0 | 32.6 | 31.8 | – |
| Jigsaw [48] | **19.2** | 30.1 | 34.7 | 33.9 | 28.3 | – |
| Count [49] | 18.0 | 30.6 | 34.3 | 32.5 | 25.7 | – |
| SplitBrain [71] | 17.7 | 29.3 | 35.4 | 35.2 | 32.8 | 11.8 |
| IR [67] | 16.8 | 26.5 | 31.8 | 34.1 | 35.6 | 31.3 |
| IR(with BN)* | 18.4 | 30.1 | 34.4 | 39.2 | 39.9 | 34.9 |
| DC [6] | 13.4 | 32.3 | **41.0** | 39.6 | 38.2 | – |
| LA (ours) | 18.7 | **32.7** | 38.1 | **42.3** | **42.4** | **38.1** |
| VGG16 | | | | | | |
| IR | 16.5 | 21.4 | 27.6 | 35.1 | 39.2 | 33.9 |
| IR(with BN)* | 13.2 | 18.7 | 27.3 | 39.8 | 50.4 | 42.1 |
| DC* | **18.2** | **27.5** | **41.5** | **51.3** | 52.7 | – |
| LA (ours) | 14.3 | 23.4 | 28.3 | 44.5 | **57.6** | **46.6** |
| ResNet-18 | | | | | | |
| IR | 16.0 | **19.9** | 29.8 | 39.0 | 44.5 | 41.0 |
| DC* | **16.4** | 17.2 | 28.7 | 44.3 | 49.1 | – |
| LA (ours) | 9.1 | 18.7 | **34.8** | **48.4** | **52.8** | **45.0** |
| ResNet-50 | | | | | | |
| IR | 15.3 | 18.8 | 24.9 | 40.6 | 54.0 | 46.5 |
| DC* | **18.9** | **27.3** | 36.7 | **52.4** | 44.2 | – |
| LA (ours) | 10.2 | 23.3 | **39.3** | 49.0 | **60.2** | **49.4** |

Table 1. ImageNet transfer learning and KNN classifier performance. Numbers within the red box are the best for the given architecture. Performances of most methods using AlexNet are taken from [6, 67]. *: performance number produced by us, please refer to the supplementary material for training details.

ImageNet transfer learning. Results shown in Table 2 illustrate that the LA method surpasses previous methods in transfer learning performance with all architectures, especially with deeper networks. Please refer to the supplementary material for $K$-nearest neighbor classification performance. These result indicate strong generalization ability of the visual representations learned via the LA method.

**Object Detection.** The results presented in Table 1 and 2 illustrate the utility of LA for learning representations for visual categorization tasks. However, visual challenges faced in real life also include other tasks, such as object detection. Therefore, we also evaluate the transfer learning ability of our models to the object detection task in the PASCAL VOC 2007 [17] dataset. The typical PASCAL detection task evaluation procedure [6, 71, 67, 64] fine-tunes unsupervised architectures using the Fast RCNN [19] method. However, Fast RCNN is substantially less computationally efficient than more recently proposed pipelines such as Faster RCNN [55] or Mask RCNN [23], and is less well-supported by validated reference implementations

| Method | conv1 | conv2 | conv3 | conv4 | conv5 |
|---|---|---|---|---|---|
| | AlexNet | | | | |
| Random | 15.7 | 20.3 | 19.8 | 19.1 | 17.5 |
| Context [13] | 19.7 | 26.7 | 31.9 | 32.7 | 30.9 |
| Color [70] | 22.0 | 28.7 | 31.8 | 31.3 | 29.7 |
| Jigsaw [48] | **23.0** | 32.1 | 35.5 | 34.8 | 31.3 |
| SplitBrain [71] | 21.3 | 30.7 | 34.0 | 34.1 | 32.5 |
| IR [67] | 18.8 | 24.3 | 31.9 | 34.5 | 33.6 |
| IR(with BN)* | 21.3 | 33.0 | 36.5 | 39.2 | 38.7 |
| DC [6] | 19.6 | **33.2** | **39.2** | 39.8 | 34.7 |
| LA (ours) | 18.7 | 32.7 | 38.2 | **40.3** | **39.5** |
| | VGG16 | | | | |
| IR | 17.6 | 23.1 | 29.5 | 33.8 | 36.3 |
| IR(with BN)* | 17.3 | 22.9 | 27.3 | 39.3 | 45.8 |
| DC* | **21.5** | **31.6** | **40.9** | **45.2** | 44.2 |
| LA (ours) | 20.1 | 25.9 | 31.9 | 44.0 | **50.0** |
| | ResNet-18 | | | | |
| IR | 17.8 | 23.0 | 30.1 | 37.0 | 38.1 |
| DC* | 16.4 | 22.5 | 30.5 | 40.4 | 41.8 |
| LA (ours) | **18.9** | **26.7** | **36.5** | **44.7** | **45.6** |
| | ResNet-50 | | | | |
| IR | 18.1 | 22.3 | 29.7 | 42.1 | 45.5 |
| DC* | **20.1** | **29.1** | 35.3 | 43.2 | 38.9 |
| LA (ours) | 10.3 | 26.4 | **39.9** | **47.2** | **50.1** |

Table 2. Places transfer learning performance. *: performances produced by us, please refer to the supplement for details.

| Method | $A_{Fast}$ | $A_{Faster}$ | $V_{Fast}$ | $V_{Faster}$ | $R_{Faster}$ |
|---|---|---|---|---|---|
| Supervised | 56.8 | 54.3 | 67.3 | 70.0 | 74.6 |
| Jigsaw [48] | 53.2 | – | – | – | – |
| Video [63] | 47.2 | – | 60.2 | – | – |
| Context [13] | 51.1 | – | 61.5 | – | – |
| Trans [64] | – | – | 63.2 | – | – |
| IR [67] | 48.1 | 53.1 | 60.5 | 65.6 | 65.4 |
| DC | **55.4** | – | **65.9** | – | – |
| LA (ours) | – | 53.5 | – | **68.4** | **69.1** |

Table 3. PASCAL VOC 2007 detection mAP. A=AlexNet, V=VGG16, and R=ResNet50. Bold numbers are the best in their columns. Performances with Faster RCNN are produced by us, except that of ResNet50 of IR, which is as reported in [67]. Most numbers using Fast RCNN are taken from [6, 67]. For numbers produced by us, we show the averages of three independent runs. Standard deviations are close to $0.2\%$ in all cases.
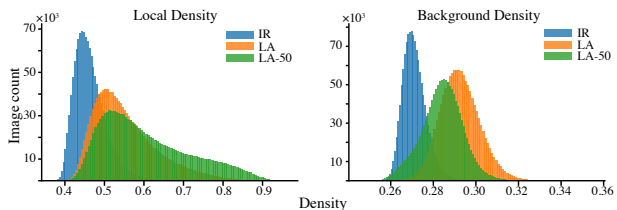


Figure 2. Distributions across all ImageNet training images of local and background densities for feature embeddings. We compare features from ResNet-18 (orange bars) and Resnet-50 (green bars) architectures as trained by the LA method, as well as that of a ResNet-18 architecture trained by the Instance Recognition (IR) method (blue bars). The local and background densities at each embedded vector are estimated by averaging dot products between that vector and, respectively, its top 30 or its 1000th-4096th, nearest neighbors in $\bar{\mathbf{V}}$. See supplementary material for more detail.

in common deep learning frameworks. To ensure training efficiency and correctness, in this work we have used the Faster RCNN pipeline from validated implementations in both TensorFlow and Pytorch. However, because the performance achieved by Faster RCNN can vary somewhat from that of Fast RCNN, direct comparison of these results to numbers generated with Fast RCNN may be misleading. For this reason, we have additionally evaluated models trained with IR and DC using Faster RCNN where possible. For implementation details, please refer to the supplementary material. Results are shown in Table 3, illustrating that the LA method achieves state-of-the-art unsupervised transfer learning for the PASCAL detection task. Interestingly, the performance gaps between the best unsupervised methods and the supervised controls are comparatively smaller for the PASCAL task than for the classification tasks.

## 5. Analysis

### 5.1. Visualizations

In this subsection, we analyze the embedding space through visualizations.

**Density distribution in the embedding space.** The LA optimization objective seeks to minimize the distances be-

tween $\mathbf{v}_i$ and $\mathbf{C}_i$ while maximizing those between $\mathbf{v}_i$ and $\mathbf{B}_i$, intuitively leading to an embedding that is locally dense at some positions but generally sparse across the space. Indeed, Figure 2 shows that the local density of the LA embedding is much higher than that created by the IR method, while the background density is only slightly higher (note differing $x$-axis scales in the figure). Moreover, insofar as deeper networks achieve lower minimums of the LA objective, we expect that their embeddings will exhibit higher local density and lower background density as compared to shallower networks. By comparing the density distributions of the ResNet-18 embedding to that of ResNet-50, Figure 2 shows that this expectation is confirmed. These results help better characterize the LA optimization procedure.

**Success and failure examples.** To help qualitatively illustrate the successes and failures of the LA objective, Figure 3 shows nearest neighbors in the training set for several validation images, both correctly and incorrectly clas-
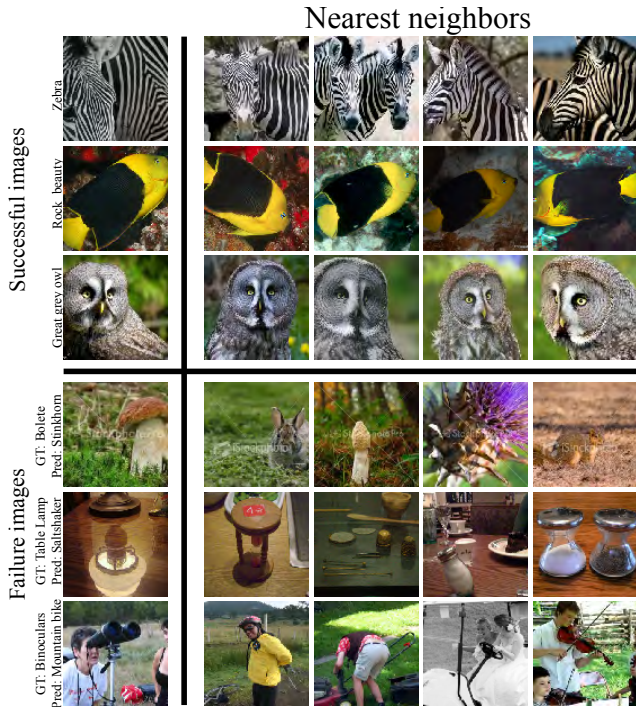
Figure 3. For each of several validation images in the left-most column, nearest neighbors in LA-trained RestNet-50 embedding, with similarity decreasing from left to right. The three top columns are successfully-classified cases, with high KNN-classifier confidence, while the lower three are failure cases, with low KNN-classifier confidence.
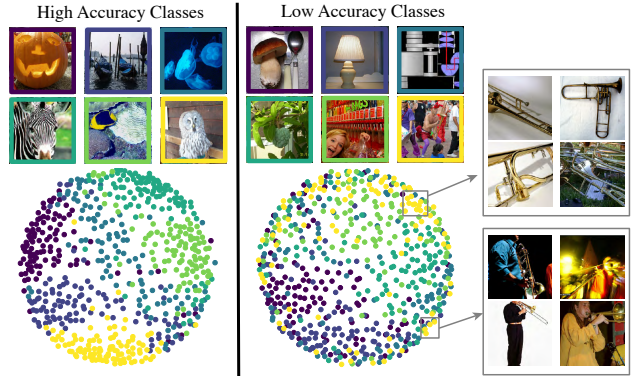


Figure 4. Multi-dimensional scaling (MDS) embedding results for network outputs of classes with high validation accuracy (left panel) and classes with low validation accuracy (right panel). For each class, we randomly choose 100 images of that class from the training set and apply the MDS algorithm to the resulting 600 images. Dots represent individual images in each color-coded category. Gray boxes show examples of images from a single class ("trombone") that have been embedded in two distinct subclusters.

| Choice of $\mathbf{B}_i$ | $\{1, 2, ..., N\}$ | Cluster-based | $\mathcal{N}_{4096}$ |
|---|---|---|---|
| NN performance | 30.2 | 33.2 | **35.7** |

Table 4. Nearest neighbor validation performances of ResNet-18 trained with different choices of $\mathbf{B}_i$. We use $H = 3$ and $m = 1000$ for cluster-based $\mathbf{B}_i$ to make the number of neighbors in $\mathbf{B}_i$ comparable to 4096. In all experiments, we use cluster-based $\mathbf{C}_i$ with $H = 1$ and $m = 10000$.

sified according to the nearest-neighbor classifier. Unsurprisingly, the successful examples show that the LA-trained model robustly groups images belonging to the same category regardless of backgrounds and view points. Interestingly, however, the network shows substantial ability to recognize high-level visual context. This is even more obvious for the failure cases, where it can be seen that the network coherently groups images according to salient characteristics. In fact, most failure cases produced by the LA model appear to be due to the inherently ill-posed nature of the ImageNet category labelling, in which the category label is only one of several potentially valid object types present in the image, and which no unsupervised method could unambiguously resolve. To further illustrate this point, we use the multi-dimensional scaling (MDS) algorithm [3] to visualize part of the embedding space (see Fig. 4). In particular, the LA successfully clusters images with trombones regardless of background, number of trombones, or viewpoint, while it (perhaps inevitably) distinguishes those images from images of humans playing trombones.

## 5.2. Ablations

In this subsection, we empirically justify the design of the LA procedure by ablating or modifying several key features of the procedure. We also provide analyses suggesting intuitive reasons underlying the meaning and influence of parameters on final performance. Please refer to the supplementary material for further analyses.

**Dynamic Locality for Background Neighbors.** We chose a nearest-neighbor based procedure for identifying $\mathbf{B}_i$ to embody the idea of dynamically rescaling the local background against which closeness is judged. We tested two ablations of our procedure that isolate the relevance of this choice, including (i) simply using all inputs for background, or (ii) using a fixed clustering-based identification procedure. (See supplement for details on how these were defined.) Experiments show that the local dynamic nearest-neighbor procedure is substantially more performant than either ablation (see Table 4). The desirability of a local rather than global background measurement is consistent with the observation that the density of features varies widely across the embedding space (see Figure 2). That the dynamic nature of the computation of the background is

| $\mathbf{C}_i$ | $\{i\}$ | $\mathcal{N}_{k'}$ | $(1, 10\text{k})$ | $(3, 10\text{k})$ | $(10, 10\text{k})$ | $(10, 30\text{k})$ |
|---|---|---|---|---|---|---|
| NN | 33.9 | 0.1 | 35.7 | 36.2 | 36.1 | **37.9** |

Table 5. Nearest neighbor validation performances of ResNet-18 trained with different choices of $\mathbf{C}_i$. All experiments use $\mathcal{N}_{4096}$ as $\mathbf{B}_i$. $\{i\}$ means $\mathbf{C}_i$ only includes $\mathbf{v}_i$ itself. $(1, 10\text{k})$ means clustering-based $\mathbf{C}_i$ with $H = 1$ and $m = 10000$. Other pairs have similar meanings. See the supplementary material for details.

useful is illustrated by the comparison of results from computing neighbors in an online fashion from $\mathbf{v}_i$, relative to the cluster-based procedure depending only on $\bar{\mathbf{V}}$.

**Robust Clustering for Close Neighbors.** We also sought to understand the importance of the specific clustering procedure for defining close neighbors $\mathbf{C}_i$. One alternative to using cluster-based identification would be to instead identify "especially close" neighbors as those within a neighborhood $\mathcal{N}_{k'}$, for some $k' \ll k$. Using this in the definition of $\mathbf{C}_i$ is equivalent to optimizing the embedding to bring especially close neighbors closer together, while somewhat further away neighbors are moved apart. While this approach would have been a conceptually simpler way to define local aggregation than the cluster-based definition of close neighbors, it turns out to be substantially less effective in producing a useful representation (see Table 5).

Given the need for cluster-based identification, a variety of alternative approaches to $k$-means are theoretically possible, including DBSCAN [16], Affinity Propagation [18], spectral methods [59, 62], and gaussian mixtures [56, 57]. However, our present context is strongly constrained by the requirement that the clustering algorithm scale well to large datasets, effectively limiting the options to $k$-means and DBSCAN. Unfortunately, DBSCAN is known to perform poorly in settings with high ambient dimensions or highly variable density distributions [16], both of which are characteristics of the embedding space we work with here (see Figure 2). Indeed, we find that replacing $k$-means with DBSCAN leads to trivial representations, across a wide variety of parameter settings (see supplement for details).

The robust clustering procedure described in Section 3.1 has several hyperparameters, including number of clusters $m$ and number of clusterings $H$. To intuitively understand their effect, we performed a set of network characterization experiments (see supplement for details). These experiments indicated that two basic factors were of importance in creating clusterings that lead to good representations: the *skewness* of the cluster of close neighbors around its intended target, as measured by the distance from the cluster center to the embedded vector $\mathbf{v}_i$, and the *size* of the cluster, as measured by its cardinality as a set. We found that (i) clusterings of close neighbors with lower skewness were robustly associated with better performance, indicating that skewness should be minimized whenever possible; and (ii)

there was an optimal size for the set of close neighbors that scaled with the representation capacity (i.e. depth) of the underlying network. Both of these facts are consistent with a picture in which the ideal embedding is one in which each category is equally likely to occur and in which each example of each category is equally "representative" – e.g. in which clusters of points corresponding to natural categories occupy isotropic spheres of equal size. Networks of smaller capacity that cannot completely achieve the optimal distribution will (poorly) approximate the optimal embedding by fracturing their embeddings of single categories into subsets that maintain isotropy by reducing the relative size of clusters, each containing only part of the true category. These considerations help explain the optimal settings for parameters $H$ and $m$: higher $H$ (i.e. more clusterings) will tend to produce more isotropic clusters, as outliers due to randomness are averaged out. However, increasing $H$ beyond a point set by the capacity of the network will lead to clusters of too large a size for the network to handle (see supplement Figure 1, from A to B, or from B to C). This negative influence can be shown in Table 5 by the slight performance drop from $(3, 10\text{k})$ to $(10, 10\text{k})$. Increasing $m$ (e.g. the number of clusters) can then compensate by decreasing the neighborhood size without increasing cluster anisotropy (see supplement Figure 1, from C to D). This conpensation can be shown in Table 5 by the performance increase from $(10, 10\text{k})$ to $(10, 30\text{k})$. More experiments detailing these conclusions are shown in the supplementary material.

## 6. Discussion

In this work, we have introduced a local aggregation (LA) objective for learning feature embeddings that seeks to discover a balance between bringing similar inputs together and allowing dissimilar inputs to move apart, embodying a principled combination of several key ideas from recent advances in unsupervised learning. We have shown that when applied to DCNNs, the LA objective creates representations that are useful for transfer learning to a variety of challenging visual tasks. We also analyze aspects of our procedure, giving an intuition for how it works.

In future work we hope to improve the LA objective along a variety of directions, including incorporating non-local manifold learning-based priors for detecting similarity, improving identification of dissimilarity via measures of representational change over multiple steps of learning, and extending to the case of non-deterministic embedding functions. We also seek to apply the LA objective beyond the image processing domain, including to video and audio signals. Finally, we hope to compare the LA procedure to biological vision systems, both in terms of the feature representations learned and the dynamics of learning during visual development.

## References

[1] J. Atkinson. The developing visual brain. 2002. 1

[2] H. B. Barlow. Unsupervised learning. *Neural computation*, 1(3):295–311, 1989. 2

[3] I. Borg and P. Groenen. Modern multidimensional scaling: Theory and applications. *Journal of Educational Measurement*, 40(3):277–280, 2003. 7

[4] J. A. Bourne and M. G. Rosa. Hierarchical development of the primate visual cortex, as revealed by neurofilament immunoreactivity: early maturation of the middle temporal area (mt). *Cerebral cortex*, 16(3):405–414, 2005. 1

[5] O. Braddick and J. Atkinson. Development of human visual function. *Vision research*, 51(13):1588–1609, 2011. 1

[6] M. Caron, P. Bojanowski, A. Joulin, and M. Douze. Deep clustering for unsupervised learning of visual features. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 132–149, 2018. 1, 2, 4, 5, 6, 12

[7] J. Carreira and A. Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6299–6308, 2017. 1

[8] L. B. Cohen and M. S. Strauss. Concept acquisition in the human infant. *Child development*, pages 419–424, 1979. 1

[9] A. Conneau, H. Schwenk, L. Barrault, and Y. Lecun. Very deep convolutional networks for natural language processing. *arXiv preprint arXiv:1606.01781*, 2, 2016. 1

[10] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. 2009. 1, 13

[11] L. Deng, G. Hinton, and B. Kingsbury. New types of deep neural network learning for speech recognition and related applications: An overview. In *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 8599–8603. IEEE, 2013. 1

[12] E. L. Denton, S. Chintala, R. Fergus, et al. Deep generative image models using a laplacian pyramid of adversarial networks. In *Advances in neural information processing systems*, pages 1486–1494, 2015. 2

[13] C. Doersch, A. Gupta, and A. A. Efros. Unsupervised visual representation learning by context prediction. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1422–1430, 2015. 1, 2, 5, 6

[14] C. Doersch and A. Zisserman. Multi-task self-supervised visual learning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2051–2060, 2017. 2

[15] J. Donahue, P. Krähenbühl, and T. Darrell. Adversarial feature learning. *arXiv preprint arXiv:1605.09782*, 2016. 1, 2

[16] M. Ester, H.-P. Kriegel, J. Sander, X. Xu, et al. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Kdd*, volume 96, pages 226–231, 1996. 8

[17] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010. 5

[18] B. J. Frey and D. Dueck. Clustering by passing messages between data points. *science*, 315(5814):972–976, 2007. 8

[19] R. Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015. 5

[20] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014. 2

[21] A. Hannun, C. Case, J. Casper, B. Catanzaro, G. Diamos, E. Elsen, R. Prenger, S. Satheesh, S. Sengupta, A. Coates, et al. Deep speech: Scaling up end-to-end speech recognition. *arXiv preprint arXiv:1412.5567*, 2014. 1

[22] R. S. Harwerth, E. L. Smith, G. C. Duncan, M. Crawford, and G. K. Von Noorden. Multiple sensitive periods in the development of the primate visual system. *Science*, 232(4747):235–238, 1986. 1

[23] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 1, 5

[24] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1, 12

[25] D. O. Hebb. *The organization of behavior*. na, 1961. 2

[26] G. Hinton, L. Deng, D. Yu, G. Dahl, A.-r. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, B. Kingsbury, et al. Deep neural networks for acoustic modeling in speech recognition. *IEEE Signal processing magazine*, 29, 2012. 1, 12

[27] G. E. Hinton, S. Osindero, and Y.-W. Teh. A fast learning algorithm for deep belief nets. *Neural computation*, 18(7):1527–1554, 2006. 2

[28] G. E. Hinton and R. R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *science*, 313(5786):504–507, 2006. 2

[29] J. Hirschberg and C. D. Manning. Advances in natural language processing. *Science*, 349(6245):261–266, 2015. 1

[30] J. J. Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the national academy of sciences*, 79(8):2554–2558, 1982. 2

[31] J. S. Husaim and L. B. Cohen. Infant learning of ill-defined categories. *Merrill-Palmer Quarterly of Behavior and Development*, pages 443–456, 1981. 1

[32] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015. 4

[33] N. Jetchev, U. Bergmann, and R. Vollgraf. Texture synthesis with spatial generative adversarial networks. *arXiv preprint arXiv:1611.08207*, 2016. 2

[34] J. Johnson, M. Douze, and H. Jégou. Billion-scale similarity search with gpus. *arXiv preprint arXiv:1702.08734*, 2017. 4

[35] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 2

[36] P. Krähenbühl, C. Doersch, J. Donahue, and T. Darrell. Data-dependent initializations of convolutional neural networks. *arXiv preprint arXiv:1511.06856*, 2015. 1

[37] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012. 1, 4

[38] A. Kumar, O. Irsoy, P. Ondruska, M. Iyyer, J. Bradbury, I. Gulrajani, V. Zhong, R. Paulus, and R. Socher. Ask me anything: Dynamic memory networks for natural language processing. In *International Conference on Machine Learning*, pages 1378–1387, 2016. 1

[39] Q. V. Le, M. Ranzato, R. Monga, M. Devin, K. Chen, G. S. Corrado, J. Dean, and A. Y. Ng. Building high-level features using large scale unsupervised learning. *arXiv preprint arXiv:1112.6209*, 2011. 2

[40] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4681–4690, 2017. 2

[41] H. Lee, R. Grosse, R. Ranganath, and A. Y. Ng. Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. In *Proceedings of the 26th annual international conference on machine learning*, pages 609–616. ACM, 2009. 2

[42] T. L. Lewis and D. Maurer. Multiple sensitive periods in human visual development: evidence from visually deprived children. *Developmental Psychobiology: The Journal of the International Society for Developmental Psychobiology*, 46(3):163–183, 2005. 1

[43] C. Li and M. Wand. Precomputed real-time texture synthesis with markovian generative adversarial networks. In *European Conference on Computer Vision*, pages 702–716. Springer, 2016. 2

[44] S. Lloyd. Least squares quantization in pcm. *IEEE transactions on information theory*, 28(2):129–137, 1982. 4

[45] T. Malisiewicz, A. Gupta, and A. A. Efros. Ensemble of exemplar-svms for object detection and beyond. 2011. 5

[46] D. Mareschal and P. C. Quinn. Categorization in infancy. *Trends in cognitive sciences*, 5(10):443–450, 2001. 1

[47] K. Noda, Y. Yamaguchi, K. Nakadai, H. G. Okuno, and T. Ogata. Audio-visual speech recognition using deep learning. *Applied Intelligence*, 42(4):722–737, 2015. 1

[48] M. Noroozi and P. Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *European Conference on Computer Vision*, pages 69–84. Springer, 2016. 1, 5, 6

[49] M. Noroozi, H. Pirsiavash, and P. Favaro. Representation learning by learning to count. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5898–5906, 2017. 1, 5

[50] A. v. d. Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 5

[51] D. Pathak, R. Girshick, P. Dollár, T. Darrell, and B. Hariharan. Learning features by watching objects move. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2701–2710, 2017. 5

[52] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2536–2544, 2016. 2

[53] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015. 2

[54] D. H. Rakison and L. M. Oakes. *Early category and concept development: Making sense of the blooming, buzzing confusion*. Oxford University Press, 2003. 1

[55] S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *Advances in Neural Information Processing Systems (NIPS)*, 2015. 5, 13

[56] D. Reynolds. Gaussian mixture models. *Encyclopedia of biometrics*, pages 827–832, 2015. 8

[57] D. A. Reynolds, T. F. Quatieri, and R. B. Dunn. Speaker verification using adapted gaussian mixture models. *Digital signal processing*, 10(1-3):19–41, 2000. 8

[58] T. D. Sanger. Optimal unsupervised learning in a single-layer linear feedforward neural network. *Neural networks*, 2(6):459–473, 1989. 2

[59] J. Shi and J. Malik. Normalized cuts and image segmentation. *Departmental Papers (CIS)*, page 107, 2000. 8

[60] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 1, 4

[61] O. Sporns, D. R. Chialvo, M. Kaiser, and C. C. Hilgetag. Organization, development and function of complex brain networks. *Trends in cognitive sciences*, 8(9):418–425, 2004. 1

[62] X. Y. Stella and J. Shi. Multiclass spectral clustering. In *null*, page 313. IEEE, 2003. 8

[63] X. Wang and A. Gupta. Unsupervised learning of visual representations using videos. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2794–2802, 2015. 1, 6

[64] X. Wang, K. He, and A. Gupta. Transitive invariance for self-supervised visual representation learning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1329–1338, 2017. 5, 6

[65] J. Wattam-Bell, D. Birtles, P. Nyström, C. Von Hofsten, K. Rosander, S. Anker, J. Atkinson, and O. Braddick. Reorganization of global form and motion processing during human visual development. *Current Biology*, 20(5):411–415, 2010. 1

[66] Z. Wu, A. A. Efros, and S. X. Yu. Improving generalization via scalable neighborhood component analysis. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 685–701, 2018. 3, 4

[67] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3733–3742, 2018. 1, 2, 3, 4, 5, 6

[68] T. Young, D. Hazarika, S. Poria, and E. Cambria. Recent trends in deep learning based natural language processing. *ieee Computational intelligenCe magazine*, 13(3):55–75, 2018. 1

[69] H. Zhang, T. Xu, H. Li, S. Zhang, X. Wang, X. Huang, and D. N. Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5907–5915, 2017. 2

[70] R. Zhang, P. Isola, and A. A. Efros. Colorful image colorization. In *European conference on computer vision*, pages 649–666. Springer, 2016. 1, 5, 6

[71] R. Zhang, P. Isola, and A. A. Efros. Split-brain autoencoders: Unsupervised learning by cross-channel prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1058–1067, 2017. 1, 2, 5, 6

[72] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva. Learning deep features for scene recognition using places database. In *Advances in neural information processing systems*, pages 487–495, 2014. 1, 5, 13

# Supplementary Material

## A. Clustering Combination

In this section we provide an illustration figure for the effects of combining multiple clusterings in Figure 1, which is also mentioned in Section 5.2. Additionally, we show the nearest neighbor validation performances in Table 1 to support our hyper-parameter choices for $H, m$ in different architectures.
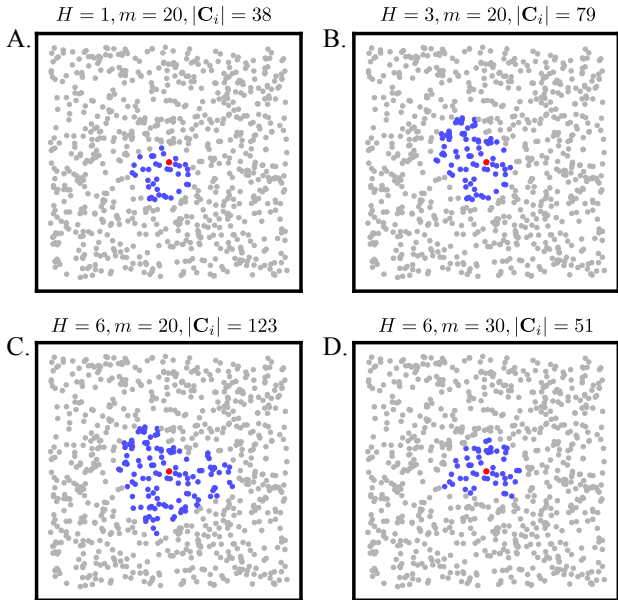


Figure 1. Illustration of the effect of combining across multiple clusterings to achieve robustness. The target embedded vector $\mathbf{v}_i$ is represented by the red dot, while blue dots represent close neighbors $\mathbf{C}_i$ under the specified hyperparameter settings.

| Network Setting | A | V | R-18 | R-50 |
|---|---|---|---|---|
| $(1, 1\text{k})$ | – | – | 35.2 | – |
| $(1, 10\text{k})$ | 30.6 | 38.9 | 35.7 | 40.2 |
| $(1, 20\text{k})$ | – | – | 35.0 | – |
| $(3, 10\text{k})$ | **31.1** | – | 36.2 | – |
| $(6, 10\text{k})$ | 30.4 | **39.7** | 37.3 | 42.4 |
| $(10, 10\text{k})$ | – | – | 36.1 | 42.3 |
| $(10, 30\text{k})$ | – | – | **37.9** | **43.4** |

Table 1. Nearest neighbor validation performances of different architectures trained with different choices of $\mathbf{C}_i$. "A" means "AlexNet". "V" means VGG16. "R" means "ResNet". Similarly to Table 5 in the main text, $(1, 10\text{k})$ means clustering-based $\mathbf{C}_i$ with $H = 1$ and $m = 10000$.

## B. Results Details

### B.1. Transfer Learning Details

Besides the settings listed in the main paper, there are additional settings for data augmentation during our transfer learning training to ImageNet and Places 205 datasets. In general, we use random crop and random horizontal flip as data augmentation techniques during transfer learning for all architectures on both ImageNet and Places 205 datasets, where the specific random crop implementation varies across networks and datasets. For AlexNet on ImageNet and all architectures on Places 205, we use the AlexNet style random crop [26], which is first resizing the image so that its smallest side is 256 and then randomly cropping a $224 \times 224$ patch. For VGG16, ResNet-18, and ResNet-50 on ImageNet, we use the ResNet style random crop [24], which is first randomly choosing a patch whose aspect ratio and area suffice two conditions and then resizing that path to $224 \times 224$. The two sufficed conditions are: its area is at least $20\%$ of the overall area and at most $100\%$ of the overall area; its aspect ratio ranges from $3/4$ to $4/3$. We use the same data augmentation techniques for the same architecture trained with different methods.

### B.2. DeepCluster Results Details

The DeepCluster [6] VGG16, ResNet-18, and ResNet-50 results are produced by us, where the DC-VGG16 network is provided by the authors and the DC-ResNet-18 and DC-ResNet-50 networks are trained by us using the provided source codes.

More specifically, for ResNet-18, two implementations of DC-ResNet-18 network are trained. Both of them modifies the standard ResNet-18 architecture by removing the final pooling and final fully connected layer and then adding additional fully connected layers, where the last layer has 10000 units. One implementation (DC-ResNet-18-A) only has that 10000-unit fully connected layer and the other implementation (DC-ResNet-18-B) has two more 4096-unit fully connected layers before that. We find that DC-ResNet-18-B performs slightly better than DC-ResNet-18-A and thus report the performances of DC-ResNet-18-B in the main paper.

Similarly for ResNet-50, two implementations (DC-ResNet-50-A and DC-ResNet-50-B) are trained. However, we find it impossible to train DC-ResNet-50-B as the $k$-means clustering results always become trivial at the third epoch. So the results reported in the paper are from DC-ResNet-50-A, which should only be slightly worse than DC-ResNet-50-B.

Other hyper-parameters for network training are mostly the same as used in the provided source codes. Meanwhile, all hyper-parameters for transfer learning to ImageNet and Places 205 are also the same as provided, except the data

augmentation techniques which are the same as described in Section B.1.

## B.3. Places KNN Results

We run models on center crops of training images in Places 205 [72] dataset to generate the memory bank $\bar{V}$. We then run the KNN validation similarly to the ImageNet [10] KNN procedure, which is described in the main paper. The results are shown in Table 2.

| Network | KNN |
|---|---|
| IR with BN - A | 36.9 |
| LA - A | **37.5** |
| IR with BN - V | 40.1 |
| LA - V | **41.9** |
| IR - R18 | 38.6 |
| LA - R18 | **40.3** |
| IR - R50 | 41.6 |
| LA - R50 | **42.4** |

Table 2. KNN results for Places 205 dataset. "A" means "AlexNet". "V" means VGG16. "R" means "ResNet".

## B.4. Faster RCNN Details

Our Faster RCNN [55] implementations are based on tf-faster-rcnn. We use SGD with momentum of 0.9, batch size 256, and weight decay 0.0001. Learning rate is initialized as 0.001 and dropped by a factor of 10 after 50000 steps. We train the models for 70000 steps. In particular, we set the number of total RoIs for training the region classifier to be 128 to reproduce the original Faster RCNN results, as indicated by [**?**]. For AlexNet, we fine-tune all layers. For VGG16, we fix "conv1" and "conv2" while fine-tuning others. For ResNet-50, we fix the first convolution layer and the first three blocks while fine-tuning others. Other hyper-parameters are the same as the default settings in tf-faster-rcnn.

## C. Other Hyperparameters

There are several other adjustable hyper-parameters in LA training procedure, such as the updating frequency for the clustering results, the parameter $k$ in $\mathcal{N}_k$ for $\mathbf{B}_i$, and whether doing clustering on $\bar{V}$ or network outputs on center crops of $\mathbf{I}$. In this section, we show results of experiments illustrating the influences of these parameters in Table 3.

| Setting | NN perf. |
|---|---|
| Baseline | 35.7 |
| $k = 2048$ | 35.4 |
| $k = 8192$ | 35.8 |
| center_crop | 35.8 |
| more_freq | 35.7 |

Table 3. Nearest neighbor validation performances for ResNet-18 trained with different settings. "Baseline" uses $H = 1, m = 10000$, and $k = 4096$. Other settings change one of the hyper-parameters while keeping the others the same. "center_crop" represents the experiment with clustering result acquired on the center crops rather than $\bar{V}$. "more_freq" represents the experiment with clustering result updated every 1000 steps.