# DEEP SUBCLASS LINEAR DISCRIMINANT ANALYSIS FOR MULTIMODAL FEATURE SPACE LEARNING

*Abin Jose*[†*], *Shen Yan*[‡*], *Mi Zhang*[‡], *Jens-Rainer Ohm*[†]

[†]Institut fur Nachrichtentechnik, RWTH Aachen University, Aachen, Germany
[‡]Michigan State University, East Lansing, USA
{jose,ohm}@ient.rwth-aachen.de,{yanshen6, mizhang}@msu.edu

## ABSTRACT

In this work, we target a known problem in representation learning that is: beyond coarse classification, how can we better model fine-grained categorization? To address this problem, we introduce Deep Subclass Linear Discriminant Analysis (DeepSDA), which utilizes intra-class variation and inter-class similarity during training. We could achieve multimodal classification by maximizing the ratio of between-subclass scatter matrix and within-subclass scatter matrix. We maximize the eigenvalues along the discriminative eignevector directions. Hence the deep neural network is able to learn more discriminative representation space and thus has higher class separation in the linearly separable latent space. We show that DeepSDA leads to significant improvements on diverse fine-grained categorization and attribute learning benchmarks.

***Index Terms***— subclass linear discriminant analysis, deep learning, multimodal optimization, fine-grained categorization, attribute distribution.

## 1. INTRODUCTION

Linear Discriminant Analysis (LDA) [1] is a classical method for finding the optimal linearly separable boundaries between classes in numerous areas such as document recognition [2, 3], face recognition [4, 5] and statistics [6, 7]. LDA aims to separate two or more classes in a supervised manner and is derived based on the assumption that all classes are sampled from unimodal Gaussian distributions with the same covariance. Unfortunately, in many real-world scenarios, data distributions are not ideally unimodal Gaussian [8]. In such case, the LDA projections may not be able to find the optimal decision boundaries in the feature space.

Zhu [9] et al. proposed Subclass Discriminant Analysis (SDA) approach, in which they derive the criteria which is able to find the appropriate division of each class into subclasses based on nearest neighborhood clustering. The new between-class scatter matrix is then used to construct the SDA objective. However, for data with Gaussian homoscedastic subclass structure, SDA cannot generate discriminant subspace such that the Bayes error is minimized. Mixture Subclass Discriminant Analysis [10, 11] (MSDA) approach tries to modify the objective function of SDA by utilizing a novel partitioning procedure to aid discrimination of data with Gaussian homoscedastic subclass structure. In both MSDA and SDA approaches, there is large overlap between models of the subclasses. For reducing this overlap, separability oriented subclass discriminant analysis (SSDA) was proposed by Wan et al. [12]. This approach

employs hierarchical clustering to subdivide a class into subclasses. By finding the subclasses for each class and maximizing Fisher's separation objective function using the re-defined within-subclass and between-subclass scatter matrices, SSDA is expected to find subclasses that better fit to the multimodal distributed data.

*As its core*, all of the above approaches require computing all possible subclasses or subclass pairs in each iteration. This requirement makes it unclear how to use it in mini-batch training with deep neural networks. However, the very first attempt of combining LDA approach for deep feature space optimization was proposed by Dorfer et al. in 2015 [13]. This approach learns linearly separable latent representations in an end-to-end manner. Instead of maximizing the likelihood of target labels as in conventional cross entropy loss for individual samples, the eigenvalues along the discriminant eigenvector directions are maximized. DeepLDA learns linearly separable hidden representations with similar discriminative power in all directions of the LDA space.

Even though, DeepLDA approach generates a low-dimensional feature space, the unimodal distribution of data is assumed which may not be suitable for attribute based-classification which was addressed in MagNet [14] approach. MagNet approach proposes a multimodal-based objective based on Neighborhood Component Analysis (NCA) [15] rather than LDA. This approach adaptively assess similarity across the different classes, and achieve local discrimination by penalizing class distribution overlap. This method aims to construct representations which are good at coarse-level classification as well as maintain more fine-grained information. However, the feature dimensions are still high. RepMet approach [16] extends MagNet approach by storing the centroids as representations that are learnable, rather than statically calculating using $K$-means.

We address the problem of attribute learning by modeling the representation space by considering the subclass level scatter matrices. Our approach, Deep Subclass Linear Discriminant Analysis (DeepSDA), thus incorporates multimodality into classification in contrast to unimodality. The novelty of our work is that by considering intra-subclass variation and inter-subclass similarity, the underlying statistics of multimodal distribution of the data can be included during training and therefore we could obtain more discriminative and compact representation in the LDA space. This is achieved by modeling the within-class and between-class scatter matrices in subclass level. We evaluate DeepSDA on three tasks such as image classification, fine-grained categorization, and attribute learning. Experimental results show competitive performance compared to state-of-the-art metric learning approaches.

The rest of the paper is organized as follows. A definition of

---

[*]Equal Contribution.

DeepSDA objective is given in Section 2. The optimization algorithm is given in Section 3. Experimental setup and results are discussed in Section 4. Concluding remarks and future work are discussed in Section 5.

## 2. WITHIN AND BETWEEN-SUBCLASS SCATTER

DeepSDA aims to learn a linearly separable latent representation space based on the multimodal distribution within the class in an end-to-end manner. Similar to DeepLDA [13], DeepSDA model the feature space to minimize the within class scatter and maximize the between class scatter. The main difference between the proposed approach and the DeepLDA [13] approach is that we consider between-subclass scatter matrix $S_b$ and within-subclass scatter matrix $S_w$.

To achieve multi-model categorization, DeepSDA initially clusters the data points in each class using the $K$-means algorithm with the number of centroids $K$. It then takes a sampled mini-batch of size $M \times D$ as input, where $M$ is the number of sampled clusters, and $D$ is the number of samples sampled from each cluster $I_m$. Given an anchor cluster $I_1$ in the class $c \in C$, DeepSDA selects $M - 1$ hard negative clusters from the set of remaining classes $\{C - c\}$. A mini batch of size $M \times D$ is then selected. The between-subclass scatter matrix $S_b$ and within-subclass scatter matrix $S_w$ used in the DeepSDA objective are defined as:

$$S_w = \frac{1}{M \times D} \sum_{m=1}^{M} \sum_{d=1}^{D} (x_{md} - \mu_m)(x_{md} - \mu_m)^T, \quad (1)$$

$$S_b = \sum_{m=1}^{M} (\mu_m - \mu)(\mu_m - \mu)^T \quad (2)$$

where $x_{md}$ is the input feature $d$ in the cluster $I_m$ of class $C$, $u_m$ is the center of cluster $I_m$ and $\mu$ is the global center of the image date.

Fig. 1 (a) and Fig. 1 (b) illustrate the intra-subclass scatter optimization and between-subclass scatter optimization during training. We consider three-class classification for the ease of illustration. The green, blue, red colors denote class 1, class 2, class 3 separately. The dark and shallow color denotes subclass 1 and subclass 2, which is initially clustered by the K-means algorithm. During training, samples within each cluster aim to be more compact, resulting in minimized intra-subclass variance. At the same time, the between-class clusters that have similar property (i.e., classes which share same attributes) are pulled away from the global center, resulting in increased inter-subclass variance but with local similarity. Thus, we could achieve multimodal classification and also bring the features of classes which share similar attributes closer.

## 3. OPTIMIZATION WITH DEEPSDA

### 3.1. Optimization Objective

To train a deep neural network with DeepSDA objective, instead of minimizing the cross entropy across different classes, the features extracted from the neural network is used to compute the subclass level scatter matrices $S_b$ and $S_w$. The resulting eigenvalue problem is given by:

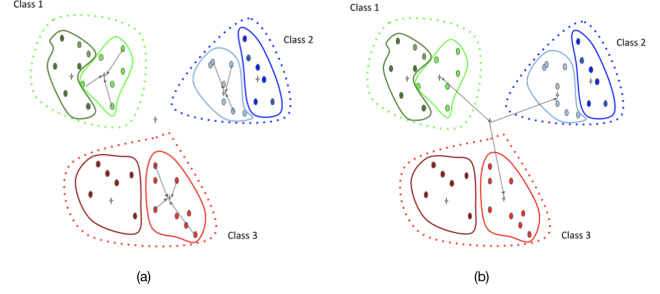$$S_b e_i = v_i (S_w + \lambda I) e_i \quad (3)$$



**Fig. 1**: (a) Within-subclass scatter and (b) between-subclass scatter optimization during training. Dark color indicates attribute 1 and shallow color indicates attribute 2.

where $e$ denotes the resulting eigenvectors and $v$ denotes the corresponding eigenvalues. The identity matrix in (3) is used to stabilize the small eigenvalues as proposed in [13]. Once the optimization is finished, each eigenvalue $v_i$ quantifies the amount of discriminative variance or separation in the direction of the corresponding eigenvector $e_i$. Hence, the goal is to maximize the individual eigenvalues which reflects the separation in the respective eigenvector directions in LDA space. Here the objective function is similar to the one proposed in [13]. However, our $S_w$ and $S_b$ are calculated in sub-class level as given in (1) and (2). This helps to achieve multi-modal classification taking into consideration the attribute concentration within the classes. The optimization objective focuses on maximizing the eigenvalues that are smaller than the mean of the smallest $C - 1$ eigenvalues thus increasing the class separability along the eigenvector directions. The optimization objective is defined as follows:

$$\mathcal{W}^* = \arg\max_{\mathcal{W}} \frac{1}{L} \sum_{l=1}^{L} v_l, \quad (4)$$

$$\{v_1, ..., v_L\} = \{v_j | v_j < \frac{1}{C-1} \sum_{c=1}^{C-1} v_c\} \quad (5)$$

$\mathcal{W}$ is the learnable parameters of the neural network. In Eq. 5, we select $L$ smallest eigenvalues that are smaller than the mean of the $C - 1$ smallest eigenvalues, and maximize them according to Eq. 4. The intuition behind this formulation is to learn a net parametrization that pushes as much discriminative variance as possible into all of the $C - 1$ available feature dimensions.

### 3.2. Nearest Neighbourhood Sampling

Conventional LDA approach comes with a potential drawback that large batch size is needed to calculate meaningful scatter matrices. This leads to memory issues in practical situations. One intuitive solution is using random sampling to build the mini-batch. However, random sampling completely ignores the finer assignment of target neighborhood structure, which is regarded as the additional prior information during training. If sampled clusters do not keep the true target neighborhood structure, the clusters sharing similar attribute information is possible to be optimized in different directions. Hence, we cannot obtain sufficient contextual insight of the neighborhood structure. In order to both liberate from the unimodality assumption and unreasonable prior target neighborhood assignments,

we adopt a nearest neighborhood sampling based approach first proposed in [14]. The key idea of the nearest neighborhood sampling based approach is to facilitate discriminative training by performing approximate nearest neighbour retrieval given a selected anchor cluster. Specifically, we first retrieve the nearest top $M - 1$ hard negative clusters. The hard negative clusters refer to clusters that do not belong to the anchor cluster class. We then randomly sample examples from the sampled hard negative clusters. In this way, at each iteration during training, the local neighbourhoods that share similar information are optimized together. Thus, the local similarity across different classes can be captured and this significantly improves training efficiency. We summarize the nearest neighborhood sampling approach which is in line with the MagNet [14] approach. However the main difference from the MagNet approach here is that, the feature space we optimize is the low-dimensional LDA space.

1. Sample an anchor cluster $I_1 \sim p_I(.)$, where we choose $p_I(I) \propto L_I$ .

2. Select one or more nearest clusters from each remaining classes set $\{C - c\}$, resulting in $M - 1$ hard negative clusters $I_2, ..., I_M$ of $I_1$.

3. For each cluster $I_m, m = 1, ...M$, randomly sample $D$ examples $\boldsymbol{x}_1^m, \boldsymbol{x}_D^m \sim p_{I_m}(.)$.

where we choose $p_I(I) \sim L_I$ (sorted loss in the minibatch). This means that we select one cluster as our anchor, and select $M - 1$ nearest imposter clusters as its local neighborhoods based on the L2 distance w.r.t the anchor cluster centroid in the feature space to construct the mini-batch. As such, as long as $M$ is reasonably large, we are able to sample clusters with similar information across different classes. $p_{I_m}(.)$ denotes the uniform distribution, namely we uniformly sample $D$ samples from each of the $M$ hard negative clusters. The main difference from the nearest neighbourhood sampling proposed in MagNet [14] approach paper is that, we do not do the sampling in the entire feature space but do it more locally. As a result, in addition to intra-cluster variance, the similar clusters are optimized together w.r.t the global center so that the local similarity can be captured.

### 3.3. Training Pipeline

Based on the optimization objective and the nearest neighbourhood sampling, we build an end-to-end training pipeline which combines all the components described above. We summarize the training pipeline in Algorithm 1 below.

---

**Algorithm 1** DeepSDA training algorithm
---
1: Forward the images $X_{\text{img}}$ through the network and compute the global mean $\boldsymbol{\mu}$ and run K-means clustering for each class $c \in C$.

2: **while** *Not Converge* **do**
3:    Select an anchor cluster $I_1$, and select $M - 1$ hard negative clusters $I_2, ..., I_M$ of $I_1$ from class set $\{C - c\}$.
4:    Randomly choose $D$ samples from each cluster $I_m$ to build mini-batch $B$ of size $M \times D$.
5:    Compute $\boldsymbol{S_w}$ and $\boldsymbol{S_b}$ and do $SVD$ to obtain the C-1 eigenvalues $v_i$.
6:    The mean of the smaller part of $v_i$ is computed using (4).
7:    Iteration $\leftarrow$ Iteration + 1
8: After each epoch, go back to 1 and continue training until convergence.

---

## 4. EXPERIMENTAL RESULTS

We evaluate the performance of DeepSDA and compare it with state-of-the-art approaches on three tasks: image classification, fine-grained categorization, and attribute learning. In particular, we use CIFAR-10 dataset [17], Oxford-102 Flowers [18] and ImageNet Attributes dataset [14] for each of the three tasks. We use the pre-trained ResNet-50 [19] on ImageNet [20] as the feature extractor. For experiments on Oxford-102 Flowers [18] and ImageNet Attributes dataset [14], images are resized to $224 \times 224$ with standard data augmentation. We start training with learning rate 0.01 with SGD optimizer. $D = 5$ is fixed for all the experiments. We only change $K$ and $M$ for different datasets.

### 4.1. Results on Image Classification

Table 1 lists the comparison results on the image classification tasks on CIFAR-10 dataset [17] using Softmax, DeepLDA [13] and DeepSDA. $M$ is set to 10 in this experiment. As shown, models trained with DeepSDA with different numbers of centroids $K$ consistently outperform the unimodal-based DeepLDA [13] approach. In particular, with $K = 2$, DeepSDA achieves the best result, which is also better than softmax based cross entropy minimization.[1]

To shed light on why DeepSDA is able to achieve better performance, we illustrate the individual eigenvalues during the training process in Fig. 2. As shown, during training, the mean value of the eigenvalues is increasing. This indicates that the discrimination power of the learned latent space also increases. We also measure the ratio, r between the maximum and minimum eigenvalues which remains almost constant after the first training iteration. This indicates that the maximum eigenvalue is nearly double the minimum eigenvalue, indicating that the discrimination power along all the eigenvector directions are relatively high. The trivial solution of LDA, gives maximum separation along the first eigenvector direction and the separation along the remaining eigenvector directions are minimum. If the network would have followed the trivial solution of LDA, this ratio would be very low since the minimum eigenvalue will be quite small compared to the maximum eigenvalue, which will decrease the classification accuracy as well.

**Table 1**: Comparison with Softmax and DeepLDA on Image Classification task on CIFAR-10 [17].

| Metric Type | K | Accuracy [%] | Dimensions |
|---|---|---|---|
| Softmax | 1 | 94.04 | 256 |
| DeepLDA [13] | 1 | 92.71 | 9 |
| DeepSDA | 2 | **94.78** | **9** |
|  | 3 | 94.27 | 9 |
|  | 4 | 94.52 | 9 |

To demonstrate the superiority of DeepSDA over unimodal-based approach, we plot the 2D t-SNE [21] representations on CIFAR-10 for Softmax and DeepSDA in Fig. 3 (a) and Fig. 3 (b) separately. For better illustration purpose, $K$ is set to 2. As shown in Fig. 3 (a), given the unimodal limitation, Softmax results in unimodal separation due to enforcement of semantic similarity. In contrast, as shown in Fig. 3 (b), DeepSDA is able to gracefully embrace the intra-class variation and accurately identify subclasses.

---

[1]It should be noted that DeepSDA is able to achieve such accuracy with a much lower dimension (i.e., 9) compared to the 256 dimensional features used in the softmax approach
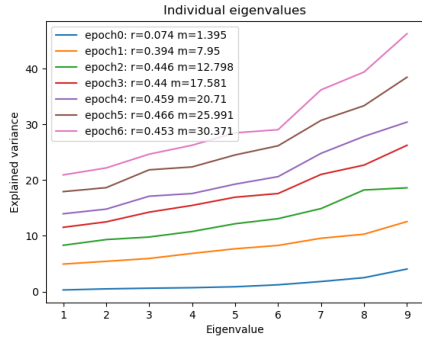
**Fig. 2**: The evolution of individual eigenvalues during CIFAR-10 training. r indicates the ratio and m indicates the mean value. We could obtain the convergence in 6 epochs as we have used pre-trained model.
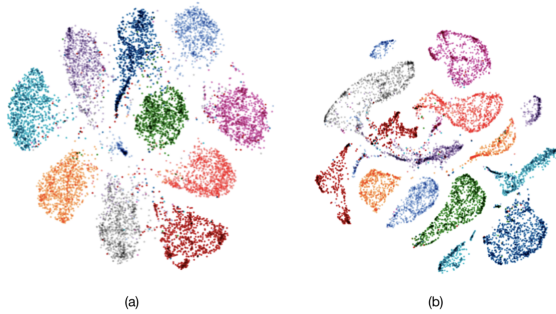


**Fig. 3**: 2D t-SNE visualization on representations for CIFAR-10 trained with (a) Softmax and (b) DeepSDA. Clearly training with DeepSDA with K=2 generates a feature space which is multimodal.

## 4.2. Results on Fine-Grained Categorization

Table 2 lists the comparison results on the fine-grained categorization task using Softmax, DeepLDA [14], Triplet [22], MagNet [14], RepMet [16] and DeepSDA with $K = 2$. As shown, DeepSDA outperforms state-of-the-art approaches by a large margin. This result again demonstrates the superiority of incorporating multimodal distribution into the neural network optimization.

**Table 2**: Comparison of test set errors of various state-of-the-art approaches on fine-grained visual categorization datasets Oxford-102 Flowers and ImageNet Attributes.

| | Error Rate [%] | |
|---|---|---|
| Approach | Oxford-102 Flowers | ImageNet Attributes |
| Softmax | 11.2 | 14.1 |
| Triplet [22] | 17.0 | 26.8 |
| MagNet [14] | 8.6 | 15.9 |
| RepMet [16] | 11.2 | 13.2 |
| DeepLDA [13] | 7.3 | 13.6 |
| DeepSDA | **3.8** | **11.9** |

As our second experiment in this task, we compare selecting different number of nearest hard negative clusters using DeepSDA and the results are summarized in Table 3. As shown, we observe monotonically improved results by increasing the number of $M$. This is because with a relatively small $M$ (*e.g.* $M = 10$), the anchor cluster is not able to retrieve enough hard negative clusters, resulting in poor local similarity across different classes.

**Table 3**: Effect of the number of nearest hard negative clusters.

| M | Error Rate [%] |
|---|---|
| 10 | 34.5 |
| 30 | 13.2 |
| 60 | **11.9** |
| 90 | 12.0 |

**Table 4**: Mean attribute precision (MAP) as a function of neighborhood size on the ImageNet Attributes dataset. The comparison with different approaches is provided.

| Approach | Modality | MAP@4 | MAP@64 | MAP@256 |
|---|---|---|---|---|
| Softmax | unimodal | 0.61 | 0.31 | 0.29 |
| Triplet [22] | unimodal | 0.50 | 0.29 | 0.27 |
| MagNet [14] | multimodal | 0.70 | 0.40 | 0.32 |
| RepMet [16] | multimodal | 0.81 | **0.52** | **0.39** |
| DeepLDA [13] | unimodal | 0.69 | 0.36 | 0.33 |
| DeepSDA | multimodal | **0.82** | 0.51 | **0.39** |

## 4.3. Results on Attribute Learning

We set $K = 2$ and $M = 60$ for the attribute learning task. Similar to the fine-grained categorization task, DeepSDA also consistently outperforms state-of-the-art approaches. In addition to error rate, we measure Mean Attribute Precision (MAP) as the second metric to measure the performance of DeepSDA. MAP is a measure proposed in [13] which computes the fraction of neighbours sharing the same attribute. The results are reported in Table 4. Training with Softmax, Triplet [22] and DeepLDA [13] objective does not perform well in the attribute learning task since it is a unimodal approach. When the neighbourhood size is increased, local similarity across different classes is very important to obtain good performance on this metric. Hence, the unimodal-based approaches perform worse than multimodal-based approaches. Compared to MagNet [14], DeepSDA obtains better MAP on all the different neighbourhood sizes. This is because features are optimized in the highly compact LDA space which has more discriminative information compared to the original feature space. Our result is comparable to RepMet [16] approach.

## 5. CONCLUSIONS

In this paper, we propose DeepSDA, a novel approach that utilizes the subclass level within-class scatter and between-class scatter to achieve multimodal class separability, which is particularly important for attribute based classification. We evaluated the performance of our algorithm on three diverse tasks. The results suggest that using multimodality approach can successfully train feature space with increased discriminability. By utilizing the underlying statistics of multimodal distribution inside the data and the optimized LDA representations, we obtained superior results. DeepSDA operates on local neighborhoods in the representation space and adaptively defines the similarity that is being optimized to account for the changing representations of the training data. Visualization of multimodal representations further suggests that exploiting the subclass structure can help the feature space to achieve local similarity across different classes. As future work, we target to store the centroid as representations that are learnable, rather than statically calculating it with K-means algorithm.

## 6. REFERENCES

[1] R.A. Fisher, "The use of multiple measurements in taxonomic problems," *Annals of eugenics*, vol. 7, pp. 179–188, 01 1936.

[2] Kari Torkkola, "Linear discriminant analysis in document classification," *IEEE TextDM 2001*, 12 2001.

[3] Chun He, Louisa Lam, and Ching Suen, "Rejection measurement based on linear discriminant analysis for document recognition," *IJDAR*, vol. 14, pp. 263–272, 09 2011.

[4] Peter Belhumeur, Joao Hespanha, and David Kriegman, *Eigenfaces vs. Fisherfaces: Recognition Using Class Specific Linear Projection*, vol. 19, pp. 43–58, 01 2006.

[5] Fatma Chelali, A. Djeradi, and R. Djeradi, "Linear discriminant analysis for face recognition," 04 2009, pp. 1–10.

[6] Edward Altman, Giancarlo Marco, and Franco Varetto, "Corporate distress diagnosis: Comparisons using linear discriminant analysis and neural networks (the italian experience)," *Journal of Banking & Finance*, vol. 18, pp. 505–529, 05 1994.

[7] Andrew Webb, *Introduction to Statistical Pattern Recognition*, pp. 1–31, 07 2003.

[8] Zizhu Fan, Yong Xu, and David Zhang, "Local linear discriminant analysis framework using sample neighbors," *Neural Networks, IEEE Transactions on*, vol. 22, pp. 1119 – 1132, 08 2011.

[9] Manli Zhu and Aleix M. Martínez, "Subclass discriminant analysis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 28, pp. 1274–1286, 2006.

[10] Nikolaos Gkalelis, Vasileios Mezaris, and Ioannis Kompatsiaris, "Mixture subclass discriminant analysis," *Signal Processing Letters, IEEE*, vol. 18, pp. 319 – 322, 06 2011.

[11] Nikolaos Gkalelis, Vasileios Mezaris, Ioannis Kompatsiaris, and Tania Stathaki, "Mixture subclass discriminant analysis link to restricted gaussian model and other generalizations," *Neural Networks and Learning Systems, IEEE Transactions on*, vol. 24, pp. 8–21, 01 2013.

[12] Huan Wan, Hui Wang, Gongde Guo, and Xin Wei, "Separability-oriented subclass discriminant analysis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, pp. 409–422, 2018.

[13] Matthias Dorfer, Rainer Kelz, and Gerhard Widmer, "Deep linear discriminant analysis," *CoRR*, vol. abs/1511.04707, 2015.

[14] Oren Rippel, Manohar Paluri, Piotr Dollár, and Lubomir D. Bourdev, "Metric learning with adaptive density discrimination," *CoRR*, vol. abs/1511.05939, 2015.

[15] Jacob Goldberger, Geoffrey E Hinton, Sam T. Roweis, and Ruslan R Salakhutdinov, "Neighbourhood components analysis," in *Advances in Neural Information Processing Systems 17*, L. K. Saul, Y. Weiss, and L. Bottou, Eds., pp. 513–520. MIT Press, 2005.

[16] Eli Schwartz, Leonid Karlinsky, Joseph Shtok, Sivan Harary, Mattias Marder, Sharath Pankanti, Rogério Schmidt Feris, Abhishek Kumar, Raja Giryes, and Alexander M. Bronstein, "Repmet: Representative-based metric learning for classification and one-shot object detection," *CoRR*, vol. abs/1806.04728, 2018.

[17] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton, "The cifar-10 dataset," *online: http://www. cs. toronto. edu/kriz/cifar. html*, 2014.

[18] M-E. Nilsback and A. Zisserman, "Delving into the whorl of flower segmentation," in *Proceedings of the British Machine Vision Conference*, 2007, vol. 1, pp. 570–579.

[19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," *CVPR*, 2016.

[20] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei, "Imagenet: A large-scale hierarchical image database," *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 248–255, 2009.

[21] Laurens van der Maaten and Geoffrey E. Hinton, "Visualizing data using t-sne," 2008.

[22] Florian Schroff, Dmitry Kalenichenko, and James Philbin, "Facenet: A unified embedding for face recognition and clustering," *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 815–823, 2015.