

Deep Multi-View Learning using Neuron-Wise Correlation-Maximizing Regularizers

Kui Jia, Jiehong Lin, Mingkui Tan, and Dacheng Tao

Abstract—Many machine learning problems concern with discovering or associating common patterns in data of multiple views or modalities. Multi-view learning is of the methods to achieve such goals. Recent methods propose deep multi-view networks via adaptation of generic Deep Neural Networks (DNNs), which concatenate features of individual views at intermediate network layers (i.e., *fusion layers*). In this work, we study the problem of multi-view learning in such end-to-end networks. We take a regularization approach via multi-view learning criteria, and propose a novel, effective, and efficient neuron-wise correlation-maximizing regularizer. We implement our proposed regularizers collectively as a *correlation-regularized network layer (CorrReg)*. CorrReg can be applied to either fully-connected or convolutional fusion layers, simply by replacing them with their CorrReg counterparts. By partitioning neurons of a hidden layer in generic DNNs into multiple subsets, we also consider a multi-view feature learning perspective of generic DNNs. Such a perspective enables us to study deep multi-view learning in the context of regularized network training, for which we present control experiments of benchmark image classification to show the efficacy of our proposed CorrReg. To investigate how CorrReg is useful for practical multi-view learning problems, we conduct experiments of RGB-D object/scene recognition and multi-view based 3D object recognition, using networks with fusion layers that concatenate intermediate features of individual modalities or views for subsequent classification. Applying CorrReg to fusion layers of these networks consistently improves classification performance. In particular, we achieve the new state of the art on the benchmark RGB-D object and RGB-D scene datasets. We make the implementation of CorrReg publicly available.

I. INTRODUCTION

MANY machine learning problems concern with discovering or associating common patterns in data of multiple views or modalities. Typical applications include retrieving images from texts or vice versa, combining visual and audio signals for content understanding, and object recognition from visual observations of multiple modalities. Data of different views usually contain complementary information, whose statistical distributions in the high-dimensional measurements of individual views may also be different. Multi-view learning methods aim to exploit information contained in multiple views to better accomplish specified learning tasks. In this work, we take image classification, in particular multi-view or multi-modal object recognition (e.g., recognizing objects from

RGB and depth images), as the primary example to study the problem of multi-view learning.

Given feature observations of different views, existing multi-view learning approaches learn latent space representations in either deterministic [1], [2], [3], [4] or probabilistic manners [5], [6], [7], [8]. The learning objective is to make resulting features of different views at each dimension of the latent space more *related* with each other, where relations may be measured by different metrics/criteria [9], [10], [4]. Among various techniques, Canonical Correlation Analysis (CCA) [1] and its extensions [11], [12], [2] are the most representative ones. For example, given two-view data, CCA learns pairs of linear projections so that in the projected space, features of both views are maximally correlated at corresponding dimensions.

Following the success of deep learning [13], [14], deep multi-view learning methods [2], [15] are also proposed recently for learning deep features from multi-view data. These methods apply multi-view learning criteria (e.g., CCA) on top of multiple single-view deep networks (cf. Figure 1-(a) for an illustration); a two-stage scheme of iterative learning is usually adopted to train the network parameters, where view-specific features are learned until the very top layers, to which either a sequential step of multi-view criteria followed by the objectives of the specified learning tasks or regularized learning objectives that seek a balance between multi-view criteria and the final tasks of interest, are applied. Alternatively, one may design deep architectures that concatenate at intermediate network layers (*fusion layers*) output features of lower, parallel layer streams for individual views [3], [8], [16], followed by upper network layers for specified learning tasks (e.g., image classification, cf. Figure 1-(b) for an illustration). Such end-to-end networks have the advantage that the final tasks of interest are achieved directly at the network outputs. However, output features of the lower, parallel streams in such networks capture view-specific patterns, which may not be aligned in a common space for a ready fusion in the subsequent layers.

To enjoy the advantage of end-to-end learning while collaboratively benefiting from different views, multi-view learning criteria could be exploited to improve the relations between resulting features of the lower, parallel streams. Under the framework of regularized function learning, this amounts to training network parameters by penalizing objectives of the main learning tasks with correlation-maximizing regularization at fusion layers (cf. Figure 1-(b)). In this work, we are interested in CCA criteria since they have long been the main workhorse for multi-view learning [11], [17], [10]. Empirical results also show that CCA based approaches outperform alter-

K. Jia and J. Lin are with the School of Electronic and Information Engineering, South China University of Technology, Guangzhou, China. Email: kuijia@scut.edu.cn, lin.jiehong@mail.scut.edu.cn.

M. Tan is with the School of Software Engineering, South China University of Technology, Guangzhou, China. Email: mingkuitan@scut.edu.cn.

D. Tao is with the School of Information Technologies, The University of Sydney, New South Wales, Australia. Email: dacheng.tao@sydney.edu.au.

native ones in the context of deep multi-view representation learning [15]. Directly using CCA as the regularizer makes network training very expensive, which is also incompatible with the mini-batch based stochastic gradient descent (SGD), the commonly used algorithm in deep network training (cf. Section II for a discussion). Inspired by batch normalization [18], we propose in this paper a novel *neuron-wise correlation-maximizing regularizer*, and implement the proposed regularizers collectively as a *correlation-regularized network layer* (*CorrReg*). CorrReg can be applied to either fully-connected (FC) or convolutional (conv) fusion layers, simply by replacing these layers with their CorrReg versions (cf. Figure 3 for an illustration). CorrReg fusion has the same computational complexity as the plain one does, which is significantly lower than that of CCA regularization.

We note that the fusion layer in a multi-view network of Figure 1-(b) is *computationally* equivalent to any hidden layer in a generic deep neural network (DNN), by partitioning neurons of the hidden layer into multiple subsets (cf. Figure 1-(c) for an illustration). This equivalence suggests a multi-view feature learning perspective of generic DNNs: when considering hidden neurons of generic DNNs as pattern detectors that characterize different patterns of the input data [19], features learned at each neuron subset of a hidden layer could be considered as a specific *view* of the input data. Ideally, each of such feature views is to learn salient or discriminative patterns of the input data, which should also be generalizable to unseen data. However, there exists an issue of *overfitting*, a phenomenon that specific subsets of layer neurons are trained to be co-adapted to certain patterns in the training data, but cannot generalize well on the held-out test data [20]. As suggested by traditional learning theory [21], the risk of overfitting could be severe for generic DNNs, since modern DNNs have large model capacities and are usually over-parameterized in the sense that they can be trained to fit randomly labelled datasets [22]. Such a risk for generic DNNs can be addressed implicitly by SGD training [23], [24], [25], and explicitly by additional regularization [20], [18]. Our proposed CorrReg takes the second regularization approach: improving correlations of neuron subsets reduces co-adaptation to possibly noisy or irrelevant, subset-specific patterns, and thus alleviates the problem of overfitting. Such a connection with generic DNNs enables us to study CorrReg in the context of regularized network training, and to compare with modern regularization techniques (e.g., Dropout [20]) for deep multi-view representation learning.

We finally summarize our contributions as follows.

- We study in this work the problem of learning deep representations from multi-view data in end-to-end networks. We take a regularization approach via multi-view learning criteria, and propose a novel, effective, and efficient neuron-wise correlation-maximizing regularizer. We implement our proposed regularizers collectively as a correlation-regularized network layer (CorrReg), which will be made publicly available. CorrReg can be applied to either FC or conv based fusion layers, simply by replacing these layers with their CorrReg versions. CorrReg fusion layer has the same computational complexity as a

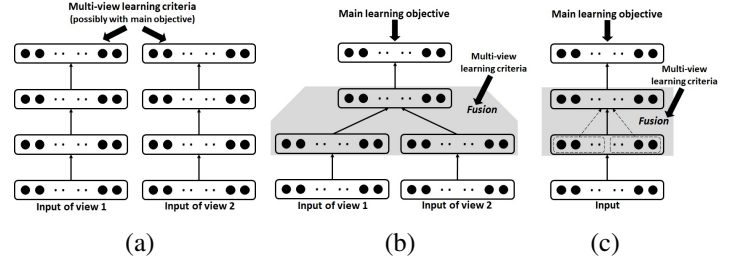


Fig. 1. Two-view illustrations of deep networks for multi-view learning, where *fusion layers* are inside shaded regions that concatenate features of individual views. (a) Deep multi-view features are learned by applying multi-view learning criteria, possibly together with the main learning objective, on top of multiple single-view deep networks; (b) a deep network takes as inputs data of multiple views/modalities for a specified learning task (e.g., multi-modal image classification), where features of individual views are concatenated at an intermediate layer (i.e., the fusion layer), and multi-view learning criteria can be imposed as a regularization on the fusion layer; (c) in generic DNNs, input neurons of a hidden layer can be partitioned into multiple subsets (represented as black circles grouped in different dashed boxes), and features learned at different subsets could be considered as multi-view features of the same input data.

plain one does, which is significantly lower than that of CCA regularization.

- We consider a multi-view feature learning perspective of generic DNNs, by partitioning neurons of a hidden layer in a generic DNN into multiple subsets. Such a connection with generic DNNs enables us to study deep multi-view representation learning in the context of regularized network training, and to compare CorrReg with modern regularization techniques (e.g., Dropout [20]). We note that such a comparison is largely ignored in existing deep multi-view learning methods.
- To investigate the efficacy of CorrReg for regularization of network training, we conduct control experiments on benchmark image classification [26] using generic DNNs [27], [28], [29], [30]. CorrReg consistently improves performance of these networks. To investigate how CorrReg is useful for practical multi-view learning problems, we conduct experiments of RGB-D object/scene recognition and multi-view based 3D object recognition, using networks with fusion layers that concatenate intermediate features of individual modalities or views for subsequent classification. Applying CorrReg to fusion layers of these networks consistently improves classification performance. In particular, we achieve the new state of the art on the benchmark RGB-D object [31] and RGB-D scene [32] datasets.

II. THE PROPOSED NEURON-WISE CORRELATION-MAXIMIZING REGULARIZERS

In this section, we first use generic DNNs to present our proposed regularization method, and explain how our method is applied to their hidden layers. Our method is readily applied to fusion layers of deep networks that take practical multi-view data, which will be introduced in Section III.

We start with a DNN composed of L FC layers. Denote its network parameters as $\Theta = \{\mathbf{W}^l, \mathbf{b}^l\}_{l=1}^L$, where \mathbf{W}^l and \mathbf{b}^l are respectively the weight matrix and bias vector

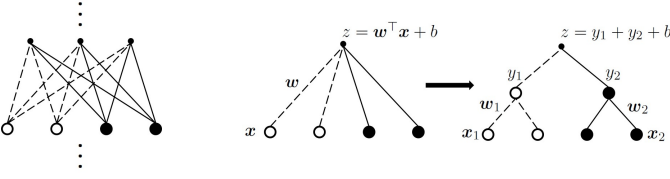


Fig. 2. Two sub-layers are formed by randomly partitioning the input neurons and the corresponding weights of a network layer into two subsets (Left). For a specified output neuron of the layer, two internal features $y_1 = \mathbf{w}_1^\top \mathbf{x}_1$ and $y_2 = \mathbf{w}_2^\top \mathbf{x}_2$ can be computed from the partition (Right). Blank circles and dashed lines represent one sub-layer, and filled circles and solid lines represent the other.

associated with the l^{th} network layer. In the setting of supervised learning, given M training samples $\mathcal{S} = \{\mathbf{s}^i\}_{i=1}^M$ of categorical data, the network parameters in Θ are optimized by minimizing the empirical risk $\frac{1}{M} \sum_{i=1}^M \text{Loss}(\mathbf{s}^i; \Theta)$, where $\text{Loss}(\cdot)$ is a properly chosen loss function, e.g., cross-entropy loss for image classification, and optimization is typically based on SGD or its variants [33]. As discussed in Section I, DNNs of high model capacities are able to learn complex functions but susceptible to overfitting. To remedy, one may apply regularization to reduce their model capacities. Adding regularization to the network training objective results in the following optimization problem

$$\min_{\Theta} \frac{1}{M} \sum_{i=1}^M \text{Loss}(\mathbf{s}^i; \Theta) + \lambda R(\Theta), \quad (1)$$

where $R(\cdot)$ is the regularizer to be specified, and λ is a trade-off parameter.

In this work, we are interested in regularizing network training using CCA based multi-view learning criteria [2], [15]. More specifically, for a specified l^{th} network layer and a training sample \mathbf{s} , denote as $\mathbf{x} \in \mathbb{R}^{n_l}$ the feature vector of \mathbf{s} computed all the way up from the input layer to the l^{th} layer. The l^{th} layer computes $f(\mathbf{z}) = f(\mathbf{W}^\top \mathbf{x} + \mathbf{b}) \in \mathbb{R}^{n_{l+1}}$, where $f(\cdot)$ is an element-wise nonlinear activation function such as ReLU [34], $\mathbf{W} \in \mathbb{R}^{n_l \times n_{l+1}}$ and $\mathbf{b} \in \mathbb{R}^{n_{l+1}}$ are the weight matrix and bias vector respectively, and where we omit the superscript l to make the following notations of better clarity. By randomly partitioning n_l input neurons/dimensions of the l^{th} layer into two subsets¹, we get feature subvectors $\mathbf{x}_1 \in \mathbb{R}^{n_{l_1}}$ and $\mathbf{x}_2 \in \mathbb{R}^{n_{l_2}}$ from \mathbf{x} , and the corresponding weight submatrices $\mathbf{W}_1 \in \mathbb{R}^{n_{l_1} \times n_{l+1}}$ and $\mathbf{W}_2 \in \mathbb{R}^{n_{l_2} \times n_{l+1}}$ from \mathbf{W} . We simply have $\mathbf{W}^\top \mathbf{x} = \mathbf{W}_1^\top \mathbf{x}_1 + \mathbf{W}_2^\top \mathbf{x}_2$. Discussions in Section I suggest that \mathbf{x}_1 and \mathbf{x}_2 can be analogously considered as two feature views of the input data. The co-adaptation of features in each view is likely to cause network training to pay more attention to the extraction of view-specific patterns, rather than the category related patterns that are desired to be learned by network training. To address this issue and benefit more from both views of the features, one may use CCA criteria to regularize network training, which aim to increase the feature correlations between the

two views. Given the two-view features $\mathbf{X}_1 = [\mathbf{x}_1^1, \dots, \mathbf{x}_1^M]$ and $\mathbf{X}_2 = [\mathbf{x}_2^1, \dots, \mathbf{x}_2^M]$ of training set \mathcal{S} , applying CCA to the l^{th} network layer amounts to optimizing \mathbf{W}_1 and \mathbf{W}_2 by

$$\begin{aligned} & \max_{\mathbf{W}_1, \mathbf{W}_2} \frac{1}{M} \text{tr}(\mathbf{W}_1^\top \mathbf{X}_1 \mathbf{X}_2^\top \mathbf{W}_2) \\ \text{s.t. } & \frac{1}{M} \mathbf{W}_1^\top \mathbf{X}_1 \mathbf{X}_1^\top \mathbf{W}_1 = \frac{1}{M} \mathbf{W}_2^\top \mathbf{X}_2 \mathbf{X}_2^\top \mathbf{W}_2 = \mathbf{I}, \end{aligned} \quad (2)$$

where \mathbf{I} is an identity matrix of compatible size. The data matrices have been assumed centered for simplicity. Applying the above problem to the l^{th} network layer also assumes implicitly that $n_{l+1} \leq \min(n_{l_1}, n_{l_2})$.

When directly using (2) as the regularizer $R(\cdot)$ in (1), network training requires solving (2) with stochastic optimization methods. Unfortunately, as pointed out in [9], the objective (2) does not easily admit a stochastic optimization due to the involvement of data covariance matrices in the constraints. One may use batch gradient descent to solve (2). It computes gradients of the correlation objective w.r.t. the CCA projected features $\mathbf{W}_1^\top \mathbf{X}_1$ and $\mathbf{W}_2^\top \mathbf{X}_2$, which in turn will be used through back-propagation to compute the gradients w.r.t. \mathbf{W}_1 and \mathbf{W}_2 , and w.r.t. all the network parameters in the layers below [2]. This is expensive as it involves computation of covariance matrices (of $\mathbf{W}_1^\top \mathbf{X}_1$ and $\mathbf{W}_2^\top \mathbf{X}_2$), their inverse square roots, and also performing matrix singular value decomposition (SVD). One may nevertheless try mini-batch based gradient descent to solve (2), which, however, may produce singular data covariance matrices; [2] also points out that solving (2) by mini-batch based stochastic optimization empirically gives unsatisfactory results. Given these challenges of directly using CCA as the regularizer $R(\cdot)$, we are motivated to find an alternative way to improve the correlations between the two feature views \mathbf{X}_1 and \mathbf{X}_2 .

Inspired by batch normalization [18], we propose to simplify the full CCA regularization in DNNs by considering the following two aspects. Firstly, instead of learning \mathbf{W}_1 and \mathbf{W}_2 to increase correlations at all the n_{l+1} dimensions of the resulting features jointly, we propose to learn $\mathbf{w}_{1,i}$ and $\mathbf{w}_{2,i}$, $i \in \{1, \dots, n_{l+1}\}$, independently for each output neuron of the l^{th} layer, where $\mathbf{w}_{1,i}$ and $\mathbf{w}_{2,i}$ are the i^{th} columns of \mathbf{W}_1 and \mathbf{W}_2 respectively. Note that such a decoupling suggests that the resulting features at the n_{l+1} dimensions could be correlated, but at the same time the constraint of $n_{l+1} \leq \min(n_{l_1}, n_{l_2})$ is also relaxed, enabling its flexible use in DNNs as tasks demand. Secondly, we use mini-batches of m training samples, rather than all the M ones, to approximate the statistics (i.e., mean and variance) necessary for computing correlations. This second simplification is enabled by the independent neuron-wise learning of $\mathbf{w}_{1,i}$ and $\mathbf{w}_{2,i}$: since in the joint case, the size m of mini-batches is required to be big enough to avoid singularity of covariance matrices.

For a specified output neuron of the l^{th} layer, we first introduce the (scalar) random variables Y_1 and Y_2 whose samples are respectively computed as $y_1 = \mathbf{w}_1^\top \mathbf{x}_1$ and $y_2 = \mathbf{w}_2^\top \mathbf{x}_2$, as illustrated in Figure 2, where we omit the neuron index for notational clarity. Given a mini-batch of m training samples, we propose in this work the following

¹For simplicity, we only consider in this work the case of partitioning neurons of a hidden layer in a generic DNN into two subsets.

neuron-wise correlation-maximizing regularizer

$$\text{Corr}(Y_1, Y_2) \approx \frac{\sum_{i=1}^m (y_1^i - \mu_1)(y_2^i - \mu_2)}{\sqrt{\sigma_1^2 \sigma_2^2 + \epsilon}}, \quad (3)$$

where the scalar ϵ is introduced for numerical stability, and

$$\begin{aligned} y_1^i &= \mathbf{w}_1^\top \mathbf{x}_1^i \quad i = 1, \dots, m, \\ y_2^i &= \mathbf{w}_2^\top \mathbf{x}_2^i \quad i = 1, \dots, m, \\ \mu_1 &= \frac{1}{m} \sum_{i=1}^m y_1^i, \quad \mu_2 = \frac{1}{m} \sum_{i=1}^m y_2^i, \\ \sigma_1^2 &= \sum_{i=1}^m (y_1^i - \mu_1)^2, \quad \sigma_2^2 = \sum_{i=1}^m (y_2^i - \mu_2)^2. \end{aligned}$$

Based on the neuron-wise regularizer (3), we specify the general objective function (1) as the following regularized problem to improve network training

$$\min_{\Theta} \frac{1}{M} \sum_{i=1}^M \text{Loss}(\mathbf{s}^i; \Theta) - \lambda \sum_{g \in \mathcal{G}} \text{Corr}(\mathcal{S}; \theta_g), \quad (4)$$

where g indexes the group \mathcal{G} of neurons in the network layers that are specified to apply regularization, θ_g denotes a subset of network parameters Θ that are involved in the computation of the incoming features of the neuron g , and we have slightly abused the use of notation Corr with that in (3). Note that θ_g and $\theta_{g'}$ may contain overlapped parameters, i.e., those parameters associated with the common layers below g and g' . The main objective (4) can be optimized using SGD or its variants [33], by sampling mini-batches of training samples in iterative steps. Details are presented shortly. Complexity analysis presented in Section II-A1 shows that compared with standard SGD training, our regularizer (3) increases computation cost only by a constant factor.

A. Correlation-Regularized Network Layer

We still use the illustration in Figure 2 as the running example. In the forward pass of a mini-batch of size m , any output neuron of the l^{th} layer that is specified to apply the regularization (3) computes $y_1^i = \mathbf{w}_1^\top \mathbf{x}_1^i$ and $y_2^i = \mathbf{w}_2^\top \mathbf{x}_2^i$, $i = 1, \dots, m$, which give output features of the neuron, before nonlinear activation $f(\cdot)$, as $z^i = y_1^i + y_2^i + b$, $i = 1, \dots, m$, where b is the bias associated with this neuron.² In the backward pass, we need to compute the gradients of the neuron-wise regularizer w.r.t. the weight vectors $\frac{\partial \text{Corr}}{\partial \mathbf{w}_1}$ and $\frac{\partial \text{Corr}}{\partial \mathbf{w}_2}$, and also those w.r.t. the input features $\frac{\partial \text{Corr}}{\partial \mathbf{x}_1^i}$, $\frac{\partial \text{Corr}}{\partial \mathbf{x}_2^i}$, $i = 1, \dots, m$. These gradients can be derived via multi-variable chain rule, and we give their explicit forms in Appendix A. Note that the later ones are used to compute through back-propagation gradients of the regularizer w.r.t. network parameters in the lower layers that are also involved in the computation of $\{z^i\}_{i=1}^m$.

We usually apply (3) to all the n_{l+1} output neurons of the specified l^{th} layer. To make regularized network training efficient, we note that for these n_{l+1} neurons, their respective internal features and statistics, i.e., $\{y_1^i\}_{i=1}^m$, $\{y_2^i\}_{i=1}^m$, μ_1 ,

μ_2 , σ_1^2 , and σ_2^2 , and also the corresponding gradients can be computed independently and in parallel. Given the mini-batch of layer inputs of the two-view features $\mathbf{x}^i = [\mathbf{x}_1^i, \mathbf{x}_2^i]$, $i = 1, \dots, m$, we write gradients of the n_{l+1} neuron-wise regularizers in the compact forms as $\left[\frac{\partial \text{Corr}}{\partial \mathbf{x}_1^i}, \frac{\partial \text{Corr}}{\partial \mathbf{x}_2^i} \right]$, $i = 1, \dots, m$, and $\left[\frac{\partial \text{Corr}}{\partial \mathbf{w}_1}, \frac{\partial \text{Corr}}{\partial \mathbf{w}_2} \right]$, where Corr denotes the n_{l+1} correlation objectives compactly. More specifically, $\frac{\partial \text{Corr}}{\partial \mathbf{x}_1^i}$ sums the gradients of neuron-wise regularizers in the layer w.r.t. the input \mathbf{x}_1^i , and $\frac{\partial \text{Corr}}{\partial \mathbf{w}_1}$ independently computes the gradient of each neuron-wise regularizer w.r.t. its associated weight vector \mathbf{w}_1 ; the same operations apply to $\frac{\partial \text{Corr}}{\partial \mathbf{x}_2^i}$ and $\frac{\partial \text{Corr}}{\partial \mathbf{w}_2}$.

In other words, we implement our proposed scheme (3) as a *correlation-regularized network layer* (CorrReg): in the forward pass, the computation is the same as a standard network layer, i.e., we do not explicitly compute the internal features $\mathbf{y}_1 = \mathbf{W}_1^\top \mathbf{x}_1$ and $\mathbf{y}_2 = \mathbf{W}_2^\top \mathbf{x}_2$, and instead we directly compute $\mathbf{z} = \mathbf{W}^\top \mathbf{x} + \mathbf{b}$, followed by element-wise nonlinear activation; in the backward pass, we compute the gradients of the correlation-regularized loss w.r.t. layer weights and layer inputs, which we spell out as

$$\begin{aligned} & \frac{1}{m} \sum_{i=1}^m \frac{\partial \text{Loss}(\mathbf{x}^i)}{\partial \mathbf{W}} - \lambda \left[\frac{\partial \text{Corr}}{\partial \mathbf{w}_1}, \frac{\partial \text{Corr}}{\partial \mathbf{w}_2} \right] \\ & \frac{\partial \text{Loss}(\mathbf{x}^i)}{\partial \mathbf{x}^i} - \lambda \left[\frac{\partial \text{Corr}}{\partial \mathbf{x}_1^i}, \frac{\partial \text{Corr}}{\partial \mathbf{x}_2^i} \right] \quad i = 1, \dots, m, \end{aligned}$$

where $\frac{\partial \text{Loss}(\mathbf{x}^i)}{\partial \mathbf{W}}$ and $\frac{\partial \text{Loss}(\mathbf{x}^i)}{\partial \mathbf{x}^i}$ can be computed via standard back-propagation. The gradient of the main loss w.r.t. the layer bias vector \mathbf{b} can be obtained in the same way.

1) *Analysis of Computational Complexity*: We analyze the additional computation cost incurred by imposing CorrReg on a network layer. Consider an l^{th} layer that computes $f(\mathbf{z}) = f(\mathbf{W}^\top \mathbf{x} + \mathbf{b}) \in \mathbb{R}^{n_{l+1}}$ for a mini-batch of m samples, and that has layer parameters $\mathbf{W} \in \mathbb{R}^{n_l \times n_{l+1}}$ and $\mathbf{b} \in \mathbb{R}^{n_{l+1}}$. Assume arithmetic with individual elements has complexity $\mathcal{O}(1)$. In the forward pass, the computation for the layer with or without CorrReg is the same, and its complexity is $\mathcal{O}(mn_l n_{l+1})$. In the backward pass, without using CorrReg the complexity for back-propagating gradients through this layer is $\mathcal{O}(mn_l n_{l+1})$. By simplifying the gradient formulas given in Appendix A and also writing them in matrix forms, we have the same complexity of $\mathcal{O}(mn_l n_{l+1})$ when using CorrReg. In summary, the complexity by imposing CorrReg on a network layer increases only by a constant factor. In contrast, using the CCA objective (2) as the regularizer involves computing inverse square root of covariance matrices of the size $n_{l+1} \times n_{l+1}$, and also performing SVD for matrix of the same size (one may refer to [2] for gradient formulas of CCA objective); it has the overall complexity of $\mathcal{O}(mn_l n_{l+1} + mn_{l+1}^2 + n_{l+1}^3)$ in the backward pass, which is significantly worse than that of CorrReg.

B. Correlation-Regularized Convolutional Networks

Our proposed CorrReg can regularize both FC and conv layers. As discussed above, for an FC layer computing $f(\mathbf{z}) = f(\mathbf{W}^\top \mathbf{x} + \mathbf{b})$, we apply regularization to the input \mathbf{z} of f by

²For each neuron that is applied the regularization, we have ever introduced trainable scalar parameters v_1 and v_2 to re-scale the internal features y_1 and y_2 , i.e., $z = v_1 y_1 + v_2 y_2 + b$, which is similar to the scheme introduced in [35]. In our experiments this alternative scheme does not necessarily improve the performance.

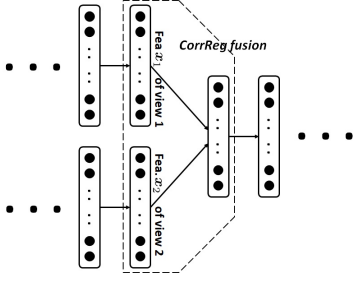


Fig. 3. Illustration of a *CorrReg fusion* layer (inside the dashed box).

performing a random two-way partition on feature dimensions of \mathbf{x} , producing two internal features $\mathbf{y}_1 = \mathbf{W}_1^\top \mathbf{x}_1$ and $\mathbf{y}_2 = \mathbf{W}_2^\top \mathbf{x}_2$, where the partition is fixed once determined. CorrReg is indeed to improve correlations between each dimension pair of \mathbf{y}_1 and \mathbf{y}_2 . It is straightforward to extend the above scheme of *single* two-way partition in CorrReg to its version of *multiple* two-way partitions. For a specified FC layer, one may simply perform multiple, random two-way partitions on feature dimensions of \mathbf{x} ; each of them produces their respective internal features, and also their respective gradients w.r.t. layer weights and layer inputs. The overall regularization imposed on this layer can be obtained by averaging the gradients from these multiple two-way partitions. Experiments investigating the efficacy of this scheme of multiple two-way partitions are reported in Section V-A.

For a conv layer, we perform single or multiple random two-way partitions on its input feature maps in the same way as for an FC layer. Each two-way partition produces two internal feature maps (corresponding to \mathbf{y}_1 and \mathbf{y}_2 in Figure 2) for each output feature map of the layer. Although features/observations at nearby locations of an image are generally correlated, we do not explicitly exploit such correlations. Instead, we independently apply regularization at each spatial location/pixel of each output feature map of the layer, so that correlations between the corresponding spatial locations in each pair of the internal feature maps are improved, where regularization is again applied before nonlinear activation. Applying our proposed CorrReg to modern architectures of DNNs (e.g., ConvNets [36], variants of ResNets [28], [29], or DenseNets [30]) is very simple: one simply replaces FC or conv layers with their CorrReg versions.

III. USE OF CORRREG FUSION IN DEEP REPRESENTATION LEARNING FROM MULTI-VIEW DATA

We present in this section how CorrReg can be readily applied to deep networks that take as inputs data of multiple views/modalities and learn deep representations from them. Different from generic DNNs, these networks by design have lower, parallel streams for data of individual views, and it is natural to apply CorrReg to the *fusion* layers where features of different views are concatenated for use in the subsequent layers. The fusion layer in such a network is usually based on FC or conv layers. To use CorrReg, one may simply treat output features of two parallel streams as the input two-view features of CorrReg (corresponding to \mathbf{x}_1 and \mathbf{x}_2 in

Section II), and replace the fusion layer with its CorrReg version, resulting in a *CorrReg fusion* layer. Figure 3 gives the illustration. In the following, we take network architectures used in our RGB-D recognition experiments to instantiate the use of CorrReg fusion.

Our networks for RGB-D object recognition and scene recognition are based on ConvNets [13] and ResNets [28]. Take an 8-layer ConvNet as the example. We modify it by enabling it to take inputs of both RGB and depth channels. The 8-layer ConvNet is composed of two lower, parallel streams followed by an upper, single stream. Outputs of the two lower streams are concatenated as inputs of the upper stream, and CorrReg is applied to the first layer of the upper stream, which thus becomes a CorrReg fusion layer. In this work, we also investigate the effects when the CorrReg fusion layer is at different “heights” (lower, middle, or upper layers) of the network. Figure 5 gives the illustration. Adaptation of ResNets is similar as above. We use such adapted network architectures for experiments of RGB-D object recognition and scene recognition in Section V-B.

IV. RELATIONS WITH EXISTING WORKS

Our proposed CorrReg method is related to four categories of existing research: regularization of generic DNNs, deep learning research that specially focuses on multi-view representation learning, RGB-D object/scene recognition, and multi-view recognition of 3D object shapes. We respectively discuss the relations as follows.

Network regularization In the literature of DNNs, various regularization techniques have been proposed to address the issue of overfitting in network training, including the traditional early stopping, weight decay, and data augmentation [37], and also the more recent Dropout [20], Dropconnect [38], batch normalization [18], and all conv layer based networks [39], [40]. Among them, Dropout and Dropconnect are most related to our proposed method. In the original proposal of Dropout [20], each hidden neuron is randomly dropped (usually with a probability of 0.5) at each training iteration, and the network is then updated on weights that are connected to the remaining neurons. During inference, all network weights are used after halving their values. Baldi and Sadowski [41] quantitatively analyze that random operations of dropout training and the associated inference can be understood as a good approximation to the expectation of outputs of a subnetwork ensemble, by introducing a bridging quantity Normalized Weighted Geometric Mean. They further show that the expectation of dropout gradients w.r.t. a network weight is approximately the gradient of the subnetwork ensemble regularized by *adaptive weight decay*. Wager *et al.* [42] present alternative interpretation of dropout training as adaptive weight decay by treating dropout as feature noising in generalized linear models. Analysis similar to [41] can be applied to Dropconnect [38] by randomly dropping weight connections rather than network neurons.

Dropout and Dropconnect achieve regularization by first sampling features/subnetworks (of shared weights), and then averaging over outputs of the subnetwork ensemble. Different from them, our CorrReg scheme explicitly increases the

correlations between (internal) features of different views, and regularization is achieved by suppressing view-specific noisy patterns. We empirically show in this work the usefulness of CorrReg in improving network training, and leave its theoretical connections with classical regularization to future research. Note that a few recent methods of network regularization explicitly reduce correlations across dimensions of the output features of a layer [43], or achieve similar effects by enforcing orthogonality of the layer weight matrix [44], [45]. These methods impose regularization complementary to our proposed CorrReg, and we are interested in the investigation of their combined use in future research.

Deep multi-view representation learning Recent deep multi-view representation learning methods include those based on CCA [2], [15], [46] and those based on auto-encoders (AE) [3]. AE based methods typically learn a shared bottleneck layer on top of lower view-specific layers, and the learned joint representation at the bottleneck layer is used for reconstruction of multiple views. Deep CCA [2] directly applies CCA to the output layers of two deep networks, so that the learned networks can produce maximally correlated features at the output layers. Wang *et al.* [15] extend deep CCA as Deep Canonically Correlated Auto-Encoders, by balancing the correlation objective between the two views with their respective reconstruction objectives.

Most of existing deep multi-view learning works take a two-stage strategy for the final tasks of interest: they first learn from data of multiple views/modalities deep features in a common space, and then use the learned features in the common space either to train classifiers for multi-view or across-view classification, or to reconstruct data of missing views. In contrast, our use of correlation based multi-view learning is to regularize training of end-to-end networks, where parameters that project multi-view data into a common space are exactly those of an intermediate network layer.

RGB-D object and scene recognition RGB-D object recognition [47], [48], [8], [49] has drawn research attention recently as a typical application of multi-view learning. Lai *et al.* [31] collect the first large-scale, hierarchical RGB-D object dataset using a Kinect camera; they show that depth information substantially helps object recognition, by concatenating hand-crafted depth features (e.g., spin images [50]) with those of RGB ones, and using the concatenated features for classification. In [16], a hierarchical learning model of Convolutional and Recursive Neural Networks (CNNs and RNNs) are proposed, where CNNs are used for learning low-level features and RNNs with random weights for efficiently extracting higher-order features; this combined model is applied to RGB and depth images separately, and the resulting features are concatenated for classification.

Eitel *et al.* [48] propose a multi-modal deep learning architecture for RGB-D object recognition, which fuses, via feature concatenation, outputs of two parallel streams of modality-specific subnetworks (composed of conv and FC layers), and uses two additional FC layers for feature transformation and softmax classification; the whole network is trained via standard back-propagation with no consideration of multi-

view learning/regularization criteria. Built on top of two parallel streams of ResNets (after removing their respective last layers of classifier) [28], a Correlated and Individual Multi-modal (CIM) learning layer is proposed in [51]; CIM aims to learn, in a discriminative and complementary manner, both correlated and modality-specific features from output vectors of the two ResNets, where “correlation” is measured by the Euclidean distance between projected features of the two ResNets’ outputs; parameters of the whole network in [51] are updated in an alternating manner: those of the two lower streams of ResNets are updated after updating of the CIM parameters (projection matrices) converges. In [52], a deep learning framework termed MDSI-CNN is proposed to learn highly discriminative and spatially invariant multi-modal feature representations at different hierarchical levels, which is technically achieved by introducing spatial transformer network [53] and Fisher encoding into CNN architectures. The problem of RGB-D image classification with limited training samples is addressed in [54]. It takes a domain adaptation approach and enforces the prediction consistency between two classifiers that are respectively learned either from the combined RGB and depth features or from the RGB features alone. Results on RGB-D object recognition show the efficacy of the proposed approach.

Methods for RGB-D scene recognition [55], [56], [57], [58] largely follow those of RGB-D object recognition. In particular, multi-modal deep architectures similar to that of [48] are still the main workhorse to get good recognition performance. We also use such a type of networks for RGB-D recognition, but with our proposed CorrReg fusion layers that have in-built neuron-wise correlation regularization. Training of CorrReg fusion based networks has no difference from standard back-propagation.

Multi-view recognition of 3D object shapes Recent research shows that multi-view images are of a promising representation for recognition of 3D object shapes. Given a 3D object model (mesh), multiple 2D images can be rendered by placing virtual cameras around the object, and recognition is based on the rendered 2D images of multiple views. Among recent methods, MVCNN [59] is a representative one that uses parallel streams of conv layers to extract features from individual views, and then aggregates these features as a global signature simply via max pooling across different views. Subsequent works improve over MVCNN by strengthening interaction among feature learning of individual views. For example, MHBN [60] uses harmonized bilinear pooling to aggregate local features, and GVCNN [61] proposes a group-view framework to model correlations among different views at a hierarchy of multiple levels. In this work, we adapt the architectural design of MVCNN by incorporating into it CorrReg fusion layers. We use such an adapted architecture to verify the usefulness of CorrReg for multi-view based 3D object recognition.

V. EXPERIMENTS

In this section, we first present control experiments of image classification to investigate the effectiveness of our proposed

TABLE I

ERROR RATES (%) ON CIFAR-10 [26] WHEN APPLYING CORRREG, WITH VARYING NUMBERS n_{Reg} OF RANDOM TWO-WAY PARTITIONS, TO DIFFERENT LAYERS OF A VARIANT OF LeNet [27]. SETTING n_{Reg} AS 0 INDICATES NO CORRREG IS APPLIED TO ANY NETWORK LAYER. EXPERIMENTS OF EACH SETTING ARE RUN FOR 5 TIMES, AND RESULTS ARE IN THE FORMAT OF MEAN (STANDARD DEVIATION).

	No CorrReg	CorrReg Conv2	CorrReg Conv3	CorrReg FC4	CorrReg FC5	CorrReg FC6
$n_{\text{Reg}} = 0$	17.42 (0.16)	-	-	-	-	-
$n_{\text{Reg}} = 1$	-	17.15 (0.11)	17.14 (0.25)	17.13 (0.27)	16.75 (0.20)	16.92 (0.30)
$n_{\text{Reg}} = 3$	-	17.17 (0.21)	17.19 (0.17)	17.25 (0.23)	16.81 (0.28)	17.08 (0.11)
$n_{\text{Reg}} = 5$	-	17.08 (0.13)	17.23 (0.23)	17.28 (0.15)	16.64 (0.20)	16.84 (0.18)
$n_{\text{Reg}} = 10$	-	17.15 (0.40)	17.12 (0.12)	17.16 (0.20)	16.87 (0.18)	17.14 (0.31)

CorrReg for regularization of network training. We use generic DNNs including ConvNet (LeNet) [27], and modern deep architectures of ResNet [28], Wide ResNet [29], DenseNet [30], and ResNeXt [62]. These experiments are conducted on the benchmark datasets of CIFAR-10, CIFAR-100 [26], and ImageNet [63]. We then present experiments of RGB-D object/scene recognition and multi-view 3D object recognition to evaluate the usefulness of CorrReg for practical multi-view learning problems. We use the benchmark datasets of RGB-D object [31], RGB-D scene [32], and ModelNet40 [64] for these experiments, and compare with the state-of-the-art results.

We use cross-entropy loss to train all these networks. Training is based on SGD with momentum. Without mentioning otherwise, we use mini-batches of size 128, momentum of 0.9, and weight decay of 0.0001; network parameters are initialized using Gaussian random weights; batch normalization is applied, before ReLU nonlinearity, in all networks to accelerate their training. In each experiment, the initial learning rate, the value of λ in (4), and also the dropping rate of Dropout (when using Dropout regularization) are determined by using 10% of training samples as the validation set. As (4) suggests, we use constant λ values for all neurons that are specified to apply CorrReg. Learning rates are decayed at the rate of 0.1 when learning curves plateau. Our implementation and experiments are based on the Torch library [65].

A. Control Experiments of Image Classification

We use the CIFAR-10 dataset for our controlled studies on a plain ConvNet (a variant of LeNet [27]). The CIFAR-10 dataset consists of 10 object categories of 60,000 color images of size 32×32 (50,000 training and 10,000 testing ones). We follow [66] and preprocess the data using global contrast normalization and ZCA whitening. Our used LeNet variant consists of 3 conv layers (the first one is the input layer), followed by 3 FC layers (the last one is the output layer). Max or average pooling layers are applied after each conv layer. More layer specifics are given in Appendix B.

We first investigate the regularization effects when applying CorrReg to different network layers. To this end, we replace each of the network layers (except the input one), namely Conv2, Conv3, FC4, FC5, and FC6, with their CorrReg versions respectively, and compare the recognition performance. As indicated in Section II-B, the scheme of multiple (random) two-way partitions can be used when applying CorrReg to any network layer. We also investigate how different numbers n_{Reg} of two-way partitions in CorrReg achieve regularization,

for which we set $n_{\text{Reg}} = 1, 3, 5$, or 10. Note that when $n_{\text{Reg}} = 1$, we use the first half neurons/feature maps of the layer as a subset, and use the other half as the second subset; when $n_{\text{Reg}} > 1$, regularization is achieved by averaging over those of the multiple two-way partitions. We run experiments of each setting (the layer/ n_{Reg} pair) for 5 times, and report results in the format of mean (standard deviation). Error rates reported in Table I tell that applying CorrReg, with any number n_{Reg} of two-way partitions, to these layers consistently achieves performance boost over the LeNet variant baseline. In general, CorrReg is more effective for (upper) FC layers; this is reasonable since a densely connected FC layer contains much more trainable parameters than a conv layer does, and is thus more susceptible to overfitting. Setting $n_{\text{Reg}} > 1$ sometimes helps in getting even better results, but at the cost of slightly increased computation. In the subsequent experiments, we simply set $n_{\text{Reg}} = 1$ for computational efficiency.

As a technique for regularization of network training, CorrReg is related to the methods [20], [38], [18] discussed in Section IV, in particular Dropout [20]. To compare CorrReg with Dropout, we apply them to the FC5 layer of LeNet variant. Since their working mechanisms are different, one might be also interested in using them together. Table II reports the comparative results. CorrReg achieves improvement comparable to that of Dropout, where the dropping rate is optimally set as 0.2 by tuning from the range of (0,1) on the validation set. Using CorrReg together with Dropout further improves the performance, showing the complementary regularization benefit of CorrReg to that of Dropout.

TABLE II

COMPARISON OF IMAGE CLASSIFICATION ON CIFAR-10 [26] WHEN APPLYING CORRREG AND/OR DROPOUT TO AN UPPER FC LAYER OF A VARIANT OF LeNet [27]. EXPERIMENTS ARE RUN FOR 5 TIMES, AND RESULTS ARE IN THE FORMAT OF MEAN (STANDARD DEVIATION).

Methods	Error rates (%)
Plain LeNet variant	17.42 (0.16)
Dropout [20]	16.72 (0.18)
CorrReg	16.75 (0.20)
CorrReg + Dropout	16.36 (0.25)

CorrReg achieves regularization via improving correlations between the internal features produced by two-way partition of a layer, which suggests a natural alternative that halves the number of layer neurons. Halving the number of layer neurons reduces the model capacity and creates “bottlenecking” of information flow, thus implicitly imposing regularization.

Oppositely, one may be also interested in alternatives that increase the number of layer neurons with varying factors. To investigate the efficacy of CorrReg for models with different capacities, we apply these alternatives again to the FC5 layer of LeNet variant. Results in Table III show that larger models perform better than smaller ones do, and applying our proposed CorrReg to larger models further improves the performance. We also compute in Table III the averaged (pair-wise) correlation among features of training samples learned at different layer neurons, in order to understand how network capacities relate to the behavior of feature correlations across layer neurons and how CorrReg plays a role here as a regularization. Results show that as the numbers of layer neurons increase, feature correlations between layer neurons increase, and applying CorrReg further enhances this effect.

TABLE III

EXPERIMENTS ON CIFAR-10 [26] USING A VARIANT OF LENET [27]. CORRREG IS OPTIONALLY APPLIED TO AN UPPER FC LAYER WITH VARYING NUMBERS OF LAYER NEURONS. RESULTS ARE IN THE FORMAT OF ERROR RATE (%) / CORRELATION COEFFICIENT ($1e^{-2}$). REFER TO THE MAIN TEXT FOR HOW CORRELATION IS COMPUTED.

Neuron No.	128	256	512	1024
W/O CorrReg	17.85/0.53	17.42/0.70	17.17/0.76	17.01/0.79
With CorrReg	17.35/0.65	16.75/0.75	16.35/0.79	16.08/0.82

CorrReg is a neuron-wise scheme of CCA regularization. One might be interested in the performance when using CCA as the regularizer. To this end, we apply the CCA objective (2) as a regularizer to the FC5 layer of the LeNet variant, where regularization parameter is set as $1e^{-6}$ by optimally tuning on the validation set. Computational complexity of CCA regularization is significantly worse than that of CorrReg, and Table IV shows that it practically consumes more time per iteration of SGD training (measured on an M40 GPU and Intel Xeon CPU running at 2.2 GHz). Although CCA regularization improves performance over that of plain LeNet variant, its results with different sizes of mini-batches are worse than those of CorrReg; we hypothesize that this is because optimization of CCA objective (particularly the constraints in (2)) is incompatible with SGD based network training.

The penalty parameter λ in (4) controls the amount of regularization that CorrReg imposes on network training. To investigate how performance of CorrReg depends on λ values, we use 10% of training samples as validation, and apply CorrReg to different layers of the LeNet variant (i.e., the settings of Table I with $n_{\text{Reg}} = 1$) using a range of λ values. Results are plotted in Figure 4. Figure 4 shows that smaller values of λ are usually better when applying CorrReg to (lower) conv layers, and larger values of λ are usually better for (upper) FC layers. This is reasonable due to two compound reasons: (1) when applying CorrReg to an upper network layer, larger values of λ are needed in order to back-propagate the regularization for better learning of features/parameters of all the layers below; (2) conv layers already have intrinsic regularization via weight sharing. This inconsistency of optimal λ values across different network layers makes use of CorrReg less convenient. Fortunately, results in this section suggest that

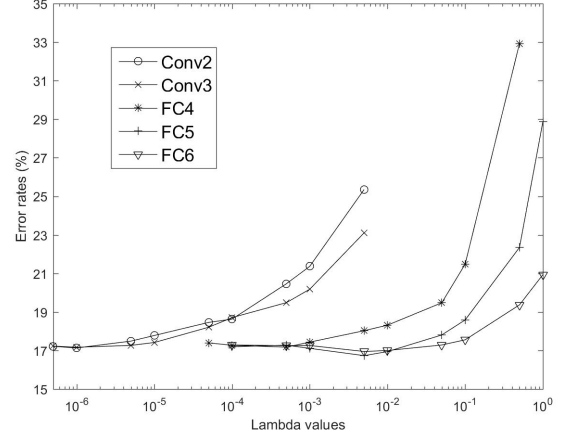


Fig. 4. Investigation of the penalty parameter λ on CorrReg’s performance. Each line represents the validation errors on CIFAR-10 [26] when applying CorrReg to a layer of a variant of LeNet [27] using a range of λ values. Each error rate is a mean out of 5 runs. These results are obtained by using 10% of CIFAR-10 training samples as validation.

to get the most effective regularization, one may simply apply CorrReg to an upper (FC) network layer, and set the optimal λ values accordingly. Setting $\lambda \in [1e^{-3}, 1e^{-1}]$ typically gives good results. Experiments in the subsequent sections follow this empirical rule.

1) Results on Modern Deep Architectures: We further investigate the regularization effects of CorrReg on modern deep architectures. For the datasets of CIFAR-10 and CIFAR-100, we use the representative architectures of ResNet [68], Wide ResNet [29], and DenseNet [30]. The CIFAR-100 dataset is an adaptation of CIFAR-10, consisting of 100 object categories of 60,000 color images. We use simple data augmentation following [67]: during training, we zero-pad 4 pixels along each image side, and sample a 32×32 region crop from the padded image or its horizontal flip; during testing, we simply use the original non-padded image. Our use of ResNet, Wide ResNet, and DenseNet for the CIFAR datasets is as follows: we use a pre-activation ResNet [68] of 68 weight layers, whose layer specifics are given in Appendix B; we use the exactly same top-performing architecture of “WRN-28-10” as in [29]; we also use the exactly same top-performing architecture of “DenseNet-BC” (with the growth rate $k = 40$) as in [30]. These architectures commonly aggregate features of lower layers via a top global average pooling layer, followed by a final FC layer of classification. For each network, we follow the empirical rule established in Section V-A and apply CorrReg to the final FC layer. We fix λ of CorrReg as $5e^{-3}$ for all the three networks. We train ResNet and Wide ResNet for a total of 160 epochs; learning rates are initialized as 0.1, and decay after 80 and 120 epoches of training. For DenseNet, we follow [30] and train for an extended duration of 300 epochs, using mini-batches of size 64. All other training hyperparameters are the same as described in the beginning of Section V (not necessarily the same as used in [68], [29], [30]). To compare with Dropout regularization, we use two schemes (denoted as Dropout-S1 and Dropout-S1S2

TABLE IV

COMPUTATION AND RECOGNITION PERFORMANCE ON CIFAR-10 [26] WHEN APPLYING CORRREG OR CCA REGULARIZATION, WITH DIFFERENT SIZES m OF MINI-BATCHES, TO AN UPPER FC LAYER OF A VARIANT OF LeNet [27]. WALL-CLOCK TIME IS MEASURED ON AN M40 GPU AND INTEL XEON CPU RUNNING AT 2.2 GHZ. EXPERIMENTS ARE RUN FOR 5 TIMES, AND ACCURACIES ARE IN THE FORMAT OF MEAN (STANDARD DEVIATION). $n_l = 256$ AND $n_{l+1} = 64$ DENOTE THE NUMBERS OF INPUT AND OUTPUT NEURONS OF THE LAYER RESPECTIVELY.

Methods	Plain LeNet variant			CorrReg			CCA regularization		
	$m = 128$	$m = 256$	$m = 512$	$m = 128$	$m = 256$	$m = 512$	$m = 128$	$m = 256$	$m = 512$
Computational Complexity	$\mathcal{O}(mn_l n_{l+1})$			$\mathcal{O}(mn_l n_{l+1})$			$\mathcal{O}(mn_l n_{l+1} + mn_{l+1}^2 + n_{l+1}^3)$		
wall-clock time per iter. (sec.)	0.009	0.016	0.030	0.013	0.020	0.035	0.025	0.035	0.048
Error rates (%)	17.42 (0.16)	17.64 (0.17)	18.33 (0.19)	16.75 (0.20)	17.13 (0.18)	17.43 (0.12)	17.24 (0.33)	17.39 (0.22)	17.65 (0.24)

TABLE V

RESULTS (ERROR RATES %) ON CIFAR-10 AND CIFAR-100 [26] OF SEVERAL DEEP ARCHITECTURES WITH OR WITHOUT REGULARIZATION. STANDARD DATA AUGMENTATION IS USED AS IN [67]. REFER TO THE MAIN TEXT FOR HOW DROPOUT-S1, DROPOUT-S1S2, CORRREG, AND CORRREG-DROPOUTS2 ARE APPLIED TO THESE ARCHITECTURES.

Architecture	CIFAR-10					CIFAR-100				
	W/O Regu.	Dropout-S1	Dropout-S1S2	CorrReg	CorrReg-DropoutS2	W/O Regu.	Dropout-S1	Dropout-S1S2	CorrReg	CorrReg-DropoutS2
ResNet [68]	6.69	6.25	6.13	6.15	6.02	27.68	28.08	26.58	27.32	26.33
Wide ResNet [29]	4.39	4.50	4.22	4.11	4.05	21.40	22.07	20.75	21.32	20.38
DenseNet [30]	3.90	3.81	3.75	3.45	3.51	18.99	18.19	18.15	18.06	17.99

respectively): scheme 1 applies Dropout to (inputs of) the final FC layer of each network, which is the same as our use of CorrReg; scheme 2 follows the way in [29] and *additionally* applies Dropout to (inputs of) the second one of the two conv layers in each residual block of these networks. Dropping rates are tuned on the validation set with the optimal one set as 0.2. We also try CorrReg together with the above scheme 2 (denoted as CorrReg-DropoutS2), to compare fairly with Dropout regularization. Results in Table V show that Dropout-S1S2 improves over Dropout-S1 by providing additional regularization, especially for the CIFAR-100 dataset that contains much fewer training samples per category than CIFAR-10 does. When applying to the top FC layer alone, our proposed CorrReg consistently outperforms Dropout-S1. Moreover, the best results are obtained by CorrReg-DropoutS2 that has the combined benefit of CorrReg and Dropout regularization.

For the ImageNet dataset, we use the representative architectures of ResNet [28], Wide ResNet [29], and ResNeXt [62] (more specifically, the “ResNet-101”, and the top-performing “WRN-50-2-bottleneck” and “ResNeXt-101 (64×4d)” of these architectures). For data augmentation, we adopt the same scheme as in [28]: during training, we randomly sample a 224×224 region crop from an image or its horizontal flip; during testing, we use a single crop of size 224×224 . Learning rates are initialized as 0.1, and decay by a factor of 0.1 at 50% and 75% of the total 90 training epochs, using mini-batches of size 256. For each network, we again apply CorrReg to the top FC layer, using a fixed λ value of $1e^{-3}$. We also apply Dropout to (inputs of) the same FC layers of these networks, where dropping rate is again set as 0.2. Table VI shows the comparative results. While Dropout regularization may not have effect on these architectures, our proposed CorrReg steadily achieves performance improvement.

TABLE VI

RESULTS OF SINGLE-CROP TESTING ON THE IMAGENET VALIDATION SET [63] OF SEVERAL DEEP ARCHITECTURES WITH OR WITHOUT REGULARIZATION. RESULTS ARE IN THE FORMAT OF TOP-1/TOP-5 ERROR RATES (%).

Architecture	W/O Regu.	With Dropout	With CorrReg
ResNet [28]	22.39/6.25	22.59/6.30	22.14/6.09
Wide ResNet [29]	22.22/6.35	22.25/6.25	22.10/6.17
ResNeXt [62]	20.90/5.46	20.82/5.70	20.58/5.40

Remarks We note that experiments in this section are not intended to compare with the best results on the benchmark image classification datasets. They are to show the efficacy of our proposed CorrReg for regularization of network training: even though input data are from the same source, intermediate features may be learned to be overfitting to view-specific patterns. CorrReg effectively regularizes network training so that the final task of classification can collaboratively benefit from all feature views.

B. RGB-D Object and Scene Recognition

In this section, we use RGB-D object dataset [31] and SUN RGB-D scene dataset [32] to investigate the efficacy of CorrReg for practical problems of multi-view learning. The RGB-D object dataset contains 207,920 RGB-D video frames of 51 classes of 300 object instances captured from different views, with roughly 600 frames per object instance. We sample these videos by every 5th frame of each video. We use the 10 random dataset splits provided by [31], with each split containing different object instances of all the 51 classes. For these splits, there are on average about 34,000 images for training and 6,900 images for testing. This dataset is intensively used for comparative studies of alternative baselines,

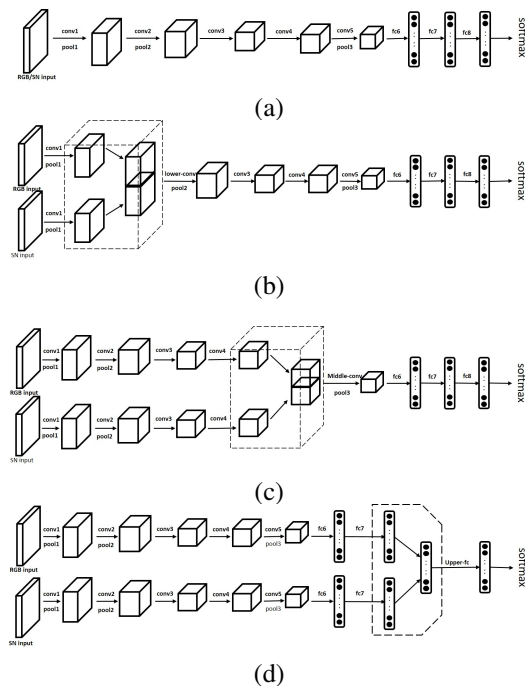


Fig. 5. Network architectures used for experiments of RGB-D object recognition. The corresponding layer parameters (numbers of feature maps, filter sizes, etc) are given in Appendix B. (a) The plain (RGB or depth) ConvNets. (b)-(d) The RGB-Depth ConvNets concatenating features of the two lower streams respectively at lower, middle, or upper “heights” of network layers.

investigation of our proposed method under different settings, and also for robustness test against contamination of input data. SUN RGB-D scene dataset is a benchmark suite for indoor scene understanding, including 10,355 RGB-D images. For scene recognition, we follow [32] and select the 19 most common categories, each of which has at least 80 RGB-D images in the dataset; we then divide the data into training and test sets, giving 4,845 images for training and 4,659 ones for testing. For both RGB-D object and SUN RGB-D datasets, we compute surface normal (SN) [69] from each depth image as input depth features.

1) *Comparative Studies of Alternative Baselines:* Our comparative studies of alternative baselines for RGB-D object recognition are based on adaptations of an 8-layer ConvNet. The adaptations consist of two lower, parallel streams followed by one upper, single stream, whose model architectures are shown in Figure 5 and whose layer parameters (numbers of feature maps, filter sizes, etc) are given in Appendix B. RGB and depth/SN images are respectively taken as inputs of the two lower streams, whose outputs are concatenated as inputs of the upper stream. We term such adapted networks as RGB-Depth ConvNets. In this work, we investigate the effects of concatenating features of the two lower streams at different “heights” (lower, middle, or upper layers) of RGB-Depth ConvNets, as illustrated in Figure 5. Alternative to RGB-Depth ConvNets are plain ConvNets that consist of one of the two lower streams of RGB-Depth ConvNets and the upper stream, as shown in Figure 5. Such plain ConvNets can be used for both RGB and depth images, and we term them

as RGB ConvNet and Depth ConvNet respectively.

The above networks suggest several baseline methods for RGB-D object recognition. In particular, given existence of both RGB and depth images during training and test phases, one may separately train RGB ConvNet or Depth ConvNet using single-modal images, and use the trained models for single-modal inference. Alternatively, one may use the above RGB-Depth ConvNets that concatenate features of individual modalities at different heights for multi-modal training and inference. To train these baseline networks, we use common data augmentation practices on the RGB-D object dataset: we first re-scale each training image to the size of 150×150 , from which or the horizontal flip of which we randomly crop a region of the size 143×143 . The learning rate is initialized as 0.01 and decays by a factor of 0.1 when learning curves plateau. Table VII shows that while concatenating features at lower or middle layers of RGB-Depth ConvNets is not effective, feature concatenation at an upper layer of RGB-Depth ConvNet achieves improved performance over single-modal networks.

TABLE VII
RECOGNITION ACCURACIES(%) ON RGB-D OBJECT DATASET [31] USING ARCHITECTURES IN FIGURE 5 AND VARIOUS REGULARIZATION METHODS. RESULTS ARE IN THE FORMAT OF MEAN \pm STANDARD DEVIATION.

Methods	Accuracy
RGB ConvNet	79.83 \pm 2.06
Depth Convnet	83.49 \pm 2.00
RGB-Depth ConvNet (concat. at a lower height)	79.48 \pm 2.64
RGB-Depth ConvNet (concat. at a middle height)	79.56 \pm 3.66
RGB-Depth ConvNet (concat. at an upper height)	87.99 \pm 1.51
RGB-Depth ConvNet with Dropout	88.33 \pm 1.69
RGB-Depth ConvNet with L2Regu	88.65 \pm 1.37
RGB-Depth ConvNet with CorrReg	89.16 \pm 1.18
RGB-Depth ConvNet with L2Regu & Dropout	88.58 \pm 1.85
RGB-Depth ConvNet with CorrReg & Dropout	89.65 \pm 1.40

TABLE VIII
ACCURACIES(%) OF DIFFERENT λ VALUES WHEN APPLYING CORRREG TO A RGB-DEPTH CONVNET (FIGURE 5-(D)) FOR RGB-D OBJECT RECOGNITION [31].

λ	10^{-4}	10^{-3}	10^{-2}	10^{-1}	10^0	10^1
mean	88.18	87.86	88.08	89.16	89.15	86.08
\pm std	± 1.22	± 1.85	± 1.26	± 1.18	± 1.93	± 1.39

The above baselines fuse multi-modal features via direct concatenation. Discussions in Section I suggest that one may apply *network regularization at fusion* to help collaboratively learn from multi-view features, which includes existing methods such as Dropout [20], and also our proposed CorrReg that explicitly leverages multi-view learning criteria. More specifically, we apply CorrReg to the first layer of the upper stream of RGB-Depth ConvNet that concatenates multi-view features at upper “height”, making it become a CorrReg fusion layer, or alternatively apply Dropout to inputs of the first layer of the upper stream. As noted in Section IV, an

TABLE IX

ROBUSTNESS TEST BY ADDING RANDOM OCCLUSION BLOCKS OF VARYING SIZES TO TEST RGB AND DEPTH (SN) IMAGES OF THE RGB-D OBJECT DATASET [31]. TRAINED NETWORKS IN SECTION V-B1 ARE USED FOR THESE EXPERIMENTS. RESULTS ARE IN TERMS OF RECOGNITION ACCURACY(%).

Occlusion size	10×10	20×20	30×30	40×40	50×50
RGB ConvNet	79.25	75.29	63.52	45.58	28.35
Depth ConvNet	83.89	82.71	80.07	71.96	54.67
RGB-Depth ConvNet	88.05	86.81	80.60	66.32	43.53
RGB-Depth ConvNet with Dropout	88.34	86.98	80.71	66.10	44.66
RGB-Depth ConvNet with CorrReg	89.27	88.66	85.33	75.73	56.75
RGB-Depth ConvNet with CorrReg & Dropout	89.70	88.93	84.54	73.28	53.46

approximate measure of correlation is used in [51] that simply computes the (squared) Euclidean distance between features of individual views. We also consider this simple correlation measure as a baseline regularizer of (1), and term such a method as L2Regu. We compare with L2Regu by applying it to the same first layer of the upper stream of RGB-Depth ConvNet. We set the λ value of CorrReg as $1e^{-1}$, the same penalty parameter of L2Regu as $1e^{-1}$, and the dropping rate of Dropout as 0.5, which are determined by tuning on the validation set. Results in Table VII show that either Dropout, L2Regu, or CorrReg improves performance over that of direct concatenation, and CorrReg outperforms Dropout and L2Regu, showing the advantage of CorrReg in practical multi-view learning problems. Table VIII also gives results of CorrReg when using different λ values. To investigate whether the effect of CorrReg (or L2Regu) is complementary to that of Dropout, we also use both CorrReg (or L2Regu) and Dropout in RGB-Depth ConvNet. Using L2Regu together with Dropout may not improve over L2Regu itself; instead, using CorrReg together with Dropout further improves the performance, showing the advantage of CorrReg for complementary regularization.

2) *Robustness against Contamination of Input Data:* An important property of multi-view learning is that inference should be less influenced when input data are contaminated [70]. In this section, we simulate such testing scenarios by adding random occlusion blocks to input RGB and depth images. Occlusion blocks are obtained by setting pixel values of the occluded regions as 0. We use the trained networks in Section V-B1 for these investigations. Table IX reports comparative results under different sizes of random occlusion. Compared with RGB ConvNet and Depth ConvNet, direct feature concatenation using RGB-Depth ConvNet may not provide better robustness against contamination of input data. RGB-Depth ConvNet with our proposed CorrReg improves the robustness, and performs consistently better than plain RGB-Depth ConvNet and the one with Dropout do.

3) *Comparisons with the state of the art:* State-of-the-art results on RGB-D object recognition are obtained by using either advanced base models (e.g., ResNets [28]) with parameters pre-trained on ImageNet [51], or advanced feature encoding scheme [52]. We follow [51] and use two ResNet-50 [28] (after removing its last FC layer of classifier) as the lower, parallel streams, whose 2048-dimensional output feature vectors are concatenated as the input vector of a FC based CorrReg fusion layer, followed by the last layer of 51-way classifier. The two lower streams respectively take RGB

and SN images as inputs. We term such a constructed network as RGB-Depth ResNet. For data augmentation, we follow [49] by first re-scaling training images to the size of 256×256 , and then randomly cropping regions of the size 224×224 from them or their horizontal flips. RGB-Depth ResNet is fine-tuned with the initial learning rates of 0.0001 for the lower streams and 0.01 for the upper stream. We use mini-batches of size 64, and set the λ value of CorrReg as $1e^{-1}$ and the dropping rate of Dropout as 0.8. Other training hyperparameters are the same as described in the beginning of Section V.

TABLE X
RECOGNITION ACCURACIES (%) OF DIFFERENT METHODS ON THE RGB-D OBJECT DATASET [31]. RESULTS ARE IN THE FORMAT OF MEAN \pm STANDARD DEVIATION.

Methods	Accuracy
Nonlinear SVM[31]	83.9 ± 3.5
CKM [71]	86.4 ± 2.3
CNN-RNN [16]	86.8 ± 3.3
upgraded HMP [69]	87.5 ± 2.9
MMSS [49]	88.5 ± 2.2
Fus-CNN [48]	91.3 ± 1.4
CIMDL-ResNet [51]	92.4 ± 1.8
MDSI-CNN [52]	92.8 ± 1.2
RGB ResNet	90.5 ± 1.6
Depth ResNet	85.5 ± 2.4
RGB-Depth ResNet	92.5 ± 1.2
RGB-Depth ResNet with Dropout	93.1 ± 1.4
RGB-Depth ResNet with CorrReg	93.4 ± 1.6
RGB-Depth ResNet with CorrReg & Dropout	93.6 ± 1.6

Note that the method [51] uses the same ImageNet pre-trained base models (i.e., ResNet-50) as we do. It is interesting to observe in Table X that our result of RGB-Depth ResNet that concatenates RGB and depth features directly is better than those from most of existing methods. Regularizing RGB-Depth ResNet with either Dropout or CorrReg further improves the result, with CorrReg achieving better improvement. Using CorrReg together with Dropout achieves the new state of the art of 93.6%.

4) *Results on RGB-D Scene Recognition:* In this section, we report experiments of RGB-D scene recognition on the SUN RGB-D dataset [32]. We use the same network architectures and training manners as in Section V-B3, with the only difference that replaces the 51-way softmax classifiers with the 19-way ones. Results in Table XI tell that RGB-Depth ResNet using simple feature concatenation outperforms existing methods that have complicated feature fusion schemes and/or training criteria. Regularizing RGB-Depth ResNet with

CorrReg and/or Dropout further improves the result to the new state of the art.

TABLE XI
RECOGNITION ACCURACIES (%) OF DIFFERENT METHODS ON THE SUN
RGB-D SCENE DATASET [32].

Methods	Accuracy
GIST + RBF Kernel SVM [32]	23.0
Place-CNN + Linear SVM [32]	37.2
Place-CNN + RBF Kernel SVM [32]	39.0
SSCNN [55]	41.3
DMFF [72]	41.5
MDSI-CNN [52]	45.2
FV-CNN [56]	48.1
RGB-D-CNN(wSVM) [57]	52.4
BilinearCNN [58]	55.5
RGB ResNet	57.8
Depth ResNet	47.2
RGB-Depth ResNet	59.6
RGB-Depth ResNet with Dropout	60.0
RGB-Depth ResNet with CorrReg	60.7
RGB-Depth ResNet with CorrReg & Dropout	61.0

C. Multi-view Recognition of 3D Object Shapes

We conduct experiments of multi-view 3D object recognition on the ModelNet40 dataset [64] to investigate the efficacy of CorrReg for practical problems with data of more than two views. The ModelNet40 dataset contains 12,311 CAD models (meshes) of 40 object categories, with 9,843 models for training and 2,468 ones for testing. To prepare images of multiple views from each object model, we follow the 1st camera set-up in [59] and assume that each model is upright oriented; 12 virtual cameras, pointing towards the model centroid, are evenly distributed (with intervals of 30 degrees) around a horizontal circle that is elevated 30 degrees from the ground plane; 2D images are rendered from these 12 camera views.

Based on a very simple architecture of MVCNN [59] for multi-view based 3D object recognition, where features of individual views extracted from lower, parallel layer streams are aggregated via feature-wise max pooling, we design a Multi-view Fusion Network (MvFusionNet) as shown in Figure 6. By pairing neighboring views, MvFusionNet re-organizes feature vectors of individual views from lower streams into an equal number of pairs of feature vectors. Each of such pairs is then fed into a fusion layer where CorrReg can also be applied to form a CorrReg fusion layer. Feature-wise max pooling is subsequently applied to outputs of these fusion layers, and MvFusionNet ends with a FC layer of classifier. We investigate here whether CorrReg is helpful for feature aggregation of different views by regularizing such a constructed MvFusionNet.

Lower streams of MvFusionNet are adapted from ResNet-101 [28] that is pre-trained on ImageNet[51]. To train MvFusionNet, we use a mini-batch of 16 (i.e., $16 \times 12 = 192$ images); the learning rates start at 0.001 and decay at the rate of 0.1 when learning curves plateau. The penalty λ of CorrReg is set as $5e^{-4}$. We report in Table XII results of MvFusionNet without or with CorrReg regularization, where we also

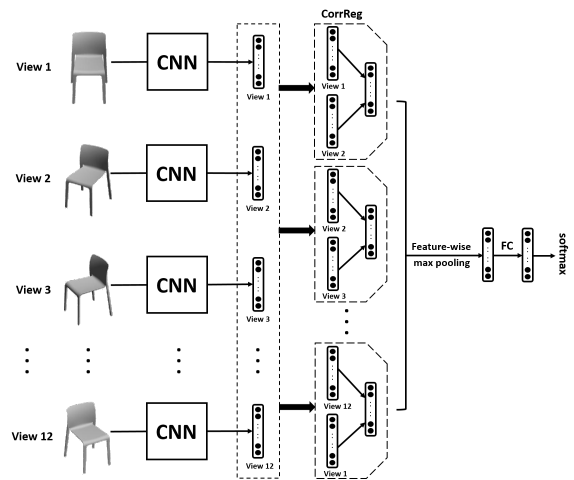


Fig. 6. Illustration of the multi-view fusion network used for experiments of multi-view 3D object recognition. Thick right arrows represent data re-organization by pairing feature vectors of individual views.

compare with recent state-of-the-art results [73], [60], [61] on ModelNet40 whose multi-view images are prepared following the same style of 1st camera set-up in [59] (i.e., 12 camera views pointing towards upright orientation of object models). Due to varying architectural designs, network optimizers, and feature aggregation schemes, results of different methods in Table XII may not be directly comparable; nevertheless, it confirms the efficacy of CorrReg for better feature learning and aggregation from multiple views of 3D object shapes. We note that results of multi-view based methods on ModelNet40 depend heavily on how multi-view images are prepared by positioning virtual cameras on a sphere enclosing the object model. For example, the current best result on ModelNet40 is obtained in [74] by selecting camera set-ups from a much richer set of camera positioning and viewpoints. We expect our results can also be boosted by using multi-view images rendered from these optimal camera set-ups.

TABLE XII
ACCURACIES (%) OF MULTI-VIEW 3D OBJECT RECOGNITION ON THE
MODELNET40 DATASET [64]. ALL METHODS USE THE 1st CAMERA
SET-UP IN [59] (12 CAMERA VIEWS POINTING TOWARDS UPRIGHT
ORIENTED MODEL). BEST RESULTS OF SOME METHODS AND THE
CORRESPONDING VIEW NUMBERS ARE ALSO QUOTED IN PARENTHESES.

Methods	Accuracy
MVCNN [59]	89.9 (90.1, 80 views)
Pairwise [75]	90.7
Dominant Set Clustering [73]	92.2 (93.0, 24 views)
MHBN [60]	93.4 (94.1, 6 views)
GVCNN [61]	92.6 (93.1, 8 views)
MvFusionNet	93.9
MvFusionNet with CorrReg	95.4

VI. CONCLUSION

We study in this paper deep multi-view learning in the context of regularized network training. We take a regularization approach via multi-view learning criteria, and propose a novel, effective, and efficient neuron-wise correlation-maximizing

regularizer. We also implement such regularizers collectively as a correlation-regularized network layer (CorrReg). CorrReg can be applied to either FC or conv based fusion layers that concatenate intermediate features of individual views. Controlled experiments of benchmark image classification show that CorrReg consistently improves performance of various modern deep architectures. Applying CorrReg to multi-modal deep networks achieves the new state of the art on the benchmark RGB-D object and scene recognition datasets. In future research, we are interested in applying CorrReg to other multi-view learning problems of practical interest.

REFERENCES

- [1] H. Hotelling, "Relations between two sets of variates," *Biometrika*, vol. 28, pp. 321–377, 1936. [1](#)
- [2] G. Andrew, R. Arora, K. Livescu, and J. Bilmes, "Deep canonical correlation analysis," in *Proceedings of the International Conference on Machine Learning*, 2013. [1](#), [3](#), [4](#), [6](#)
- [3] J. Ngiam, A. Khosla, M. Kim, J. Nam, H. Lee, and A. Y. Ng, "Multimodal deep learning," in *Proceedings of the International Conference on Machine Learning*, June 2011. [1](#), [6](#)
- [4] A. Frome, G. S. Corrado, J. Shlens, S. Bengio, J. Dean, M. Ranzato, and T. Mikolov, "Devise: A deep visual-semantic embedding model," in *Advances in Neural Information Processing Systems* 26, 2013, pp. 2121–2129. [1](#)
- [5] F. R. Bach and M. I. Jordan, "A probabilistic interpretation of canonical correlation analysis," Department of Statistics, University of California, Berkeley, Tech. Rep, Tech. Rep., 2005. [1](#)
- [6] D. A. Cohn and T. Hofmann, "The missing link - a probabilistic model of document content and hypertext connectivity," in *Advances in Neural Information Processing Systems* 13, 2001, pp. 430–436. [1](#)
- [7] K. Barnard, P. Duygulu, D. Forsyth, N. de Freitas, D. M. Blei, and M. I. Jordan, "Matching words and pictures," *J. Mach. Learn. Res.*, vol. 3, pp. 1107–1135, 2003. [1](#)
- [8] N. Srivastava and R. R. Salakhutdinov, "Multimodal learning with deep boltzmann machines," in *Advances in Neural Information Processing Systems* 25, 2012, pp. 2222–2230. [1](#), [6](#)
- [9] R. Arora, A. Cotter, K. Livescu, and N. Srebro, "Stochastic optimization for PCA and PLS," in *50th Annual Allerton Conference on Communication, Control, and Computing, Allerton, Allerton Park & Retreat Center, Monticello, IL, USA, October 1-5, 2012*, 2012, pp. 861–868. [1](#), [3](#)
- [10] D. R. Hardoon, S. R. Szedmak, and J. R. Shawe-taylor, "Canonical correlation analysis: An overview with application to learning methods," *Neural Comput.*, vol. 16, no. 12, pp. 2639–2664, Dec. 2004. [1](#)
- [11] P. L. Lai and C. Fyfe, "Kernel and nonlinear canonical correlation analysis," *Int. J. Neural Syst.*, vol. 10, pp. 365–377, 2000. [1](#)
- [12] D. R. Hardoon and J. Shawe-Taylor, "Sparse canonical correlation analysis," *Machine Learning*, vol. 83, no. 3, pp. 331–353, 2011. [1](#)
- [13] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems* 25, 2012, pp. 1097–1105. [1](#), [5](#)
- [14] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Region-based convolutional networks for accurate object detection and segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 1, pp. 142–158, Jan. 2016. [1](#)
- [15] W. Wang, R. Arora, K. Livescu, and J. Bilmes, "On deep multi-view representation learning," in *Proceedings of the International Conference on Machine Learning*, 2015. [1](#), [2](#), [3](#), [6](#)
- [16] R. Socher, B. Huval, B. Bath, C. D. Manning, and A. Y. Ng, "Convolutional-recursive deep learning for 3d object classification," in *Advances in Neural Information Processing Systems* 25, 2012, pp. 656–664. [1](#), [6](#), [11](#)
- [17] S. Akaho, "A kernel method for canonical correlation analysis," in *Proceedings of the International Meeting of the Psychometric Society*, 2001. [1](#)
- [18] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *Proceedings of the 32nd International Conference on Machine Learning*, 2015, pp. 448–456. [2](#), [3](#), [5](#), [7](#)
- [19] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *Proceedings of the 13th European Conference on Computer Vision*, 2014, pp. 818–833. [2](#)
- [20] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Improving neural networks by preventing co-adaptation of feature detectors," *CoRR*, vol. abs/1207.0580, 2012. [2](#), [5](#), [7](#), [10](#)
- [21] M. Mohri, A. Rostamizadeh, and A. Talwalkar, *Foundations of Machine Learning*. The MIT Press, 2012. [2](#)
- [22] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, "Understanding deep learning requires rethinking generalization," in *International Conference on Learning Representations (ICLR)*, 2017. [2](#)
- [23] M. Hardt, B. Recht, and Y. Singer, "Train faster, generalize better: Stability of stochastic gradient descent," in *Proceedings of The 33rd International Conference on Machine Learning*, vol. 48, 2016, pp. 1225–1234. [2](#)
- [24] I. Kuzborskij and C. H. Lampert, "Data-dependent stability of stochastic gradient descent," *CoRR*, vol. abs/1703.01678, 2017. [2](#)
- [25] C. Zhang, Q. Liao, A. Rakhlin, B. Miranda, N. Golowich, and T. A. Poggio, "Theory of deep learning ii: Optimization properties of SGD," *CoRR*, vol. abs/1801.02254, 2018. [2](#)
- [26] A. Krizhevsky, "Learning multiple layers of features from tiny images," *Tech. Report*, 2009. [2](#), [7](#), [8](#), [9](#)
- [27] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998. [2](#), [6](#), [7](#), [8](#), [9](#)
- [28] 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*, 2016, pp. 770–778. [2](#), [5](#), [6](#), [7](#), [9](#), [11](#), [12](#), [15](#)
- [29] S. Zagoruyko and N. Komodakis, "Wide residual networks," in *Proceedings of the British Machine Vision Conference*, 2016. [2](#), [5](#), [7](#), [8](#), [9](#)
- [30] G. Huang, Z. Liu, and K. Q. Weinberger, "Densely connected convolutional networks," *CoRR*, vol. abs/1608.06993, 2016. [2](#), [5](#), [7](#), [8](#), [9](#)
- [31] K. Lai, L. Bo, X. Ren, and D. Fox, "A large-scale hierarchical multi-view rgb-d object dataset," in *Proceedings of the IEEE International Conference on Robotics and Automation*, 2011. [2](#), [6](#), [7](#), [9](#), [10](#), [11](#)
- [32] S. Song, S. P. Lichtenberg, and J. Xiao, "Sun rgb-d: A rgb-d scene understanding benchmark suite," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2015, pp. 567–576. [2](#), [7](#), [9](#), [10](#), [11](#), [12](#)
- [33] I. Sutskever, J. Martens, G. E. Dahl, and G. E. Hinton, "On the importance of initialization and momentum in deep learning," in *Proceedings of the 30th International Conference on Machine Learning*, vol. 28, no. 3, May 2013, pp. 1139–1147. [3](#), [4](#)
- [34] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *Proceedings of the 27th International Conference on Machine Learning*, 2010, pp. 807–814. [3](#)
- [35] T. Salimans and D. P. Kingma, "Weight normalization: A simple reparameterization to accelerate training of deep neural networks," *CoRR*, vol. abs/1602.07868, 2016. [4](#)
- [36] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *CoRR*, vol. abs/1409.1556, 2014. [5](#)
- [37] D. Ciresan, U. Meier, and J. Schmidhuber, "Multi-column deep neural networks for image classification," in *Proceedings of the 2012 IEEE Conference on Computer Vision and Pattern Recognition*, ser. CVPR '12, 2012, pp. 3642–3649. [5](#)
- [38] L. Wan, M. D. Zeiler, S. Zhang, Y. LeCun, and R. Fergus, "Regularization of neural networks using dropconnect," in *Proceedings of the 30th International Conference on Machine Learning*, 2013, pp. 1058–1066. [5](#), [7](#)
- [39] M. Lin, Q. Chen, and S. Yan, "Network in network," in *Proceedings of the International Conference on Learning Representations*, 2013. [5](#)
- [40] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. A. Riedmiller, "Striving for simplicity: The all convolutional net," *CoRR*, vol. abs/1412.6806, 2014. [5](#)
- [41] P. Baldi and P. J. Sadowski, "Understanding dropout," in *Advances in Neural Information Processing Systems* 26, 2013, pp. 2814–2822. [5](#)
- [42] S. Wager, S. Wang, and P. S. Liang, "Dropout training as adaptive regularization," in *Advances in Neural Information Processing Systems* 26, 2013, pp. 351–359. [5](#)
- [43] M. Cogswell, F. Ahmed, R. Girshick, L. Zitnick, and D. Batra, "Reducing overfitting in deep networks by decorrelating representations," in *International Conference on Learning Representations*, 2016. [6](#)
- [44] P. Rodriguez, J. Gonzalez, G. Cucurull, J. M. Gonfau, and X. Roca, "Regularizing cnns with locally constrained decorrelations," in *International Conference on Learning Representations*, 2017. [6](#)
- [45] K. Jia, D. Tao, S. Gao, and X. Xu, "Improving training of deep neural networks via singular value bounding," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. [6](#)

- [46] S. Chandar, M. M. Khapra, H. Larochelle, and B. Ravindran, "Correlational neural networks," *CoRR*, vol. abs/1504.07225, 2015. 6
- [47] L. Bo, X. Ren, and D. Fox, "Hierarchical matching pursuit for image classification: Architecture and fast algorithms," in *Advances in Neural Information Processing Systems* 24, 2011, pp. 2115–2123. 6
- [48] A. Eitel, J. T. Springenberg, L. Spinello, M. Riedmiller, and W. Burgard, "Multimodal deep learning for robust rgb-d object recognition," in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2015, pp. 681–687. 6, 11
- [49] A. Wang, J. Cai, J. Lu, and T.-J. Cham, "Mmss: Multi-modal sharable and specific feature learning for rgb-d object recognition," in *Proceedings of the IEEE International Conference on Computer Vision*, December 2015. 6, 11
- [50] A. E. Johnson and M. Hebert, "Using spin images for efficient object recognition in cluttered 3d scenes," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 21, no. 5, pp. 433–449, May 1999. 6
- [51] Z. Wang, J. Lu, R. Lin, J. Feng, and J. Zhou, "Correlated and individual multi-modal deep learning for rgb-d object recognition," *CoRR*, vol. arXiv:1604.01655, 2016. 6, 11, 12
- [52] U. Asif, M. Bennamoun, and F. A. Sohel, "A multi-modal, discriminative and spatially invariant cnn for rgb-d object labeling," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 9, pp. 2051–2065, Sept. 2018. 6, 11, 12
- [53] M. Jaderberg, K. Simonyan, A. Zisserman, and k. kavukcuoglu, "Spatial transformer networks," in *Advances in Neural Information Processing Systems* 28, 2015, pp. 2017–2025. 6
- [54] X. Li, M. Fang, J.-J. Zhang, and J. Wu, "Learning coupled classifiers with rgb images for rgb-d object recognition," *Pattern Recognition*, vol. 61, no. C, pp. 433–446, Jan. 2017. 6
- [55] Y. Liao, S. Kodagoda, Y. Wang, L. Shi, and Y. Liu, "Understand scene categories by objects: A semantic regularized scene classifier using convolutional neural networks," in *Proceedings of the IEEE International Conference on Robotics and Automation*, 2016, pp. 2318–2325. 6, 12
- [56] A. Wang, J. Cai, J. Lu, and T.-J. Cham, "Modality and component aware feature fusion for rgb-d scene classification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 5995–6004. 6, 12
- [57] X. Song, L. Herranz, and S. Jiang, "Depth cnns for rgb-d scene recognition: Learning from scratch better than transferring from rgb-cnns," in *AAAI*, 2017, pp. 4271–4277. 6, 12
- [58] H. F. Zaki, F. Shafait, and A. Mian, "Learning a deeply supervised multi-modal rgb-d embedding for semantic scene and object category recognition," *Robotics and Autonomous Systems*, vol. 92, pp. 41–52, 2017. 6, 12
- [59] H. Su, S. Maji, E. Kalogerakis, and E. Learned-Miller, "Multi-view convolutional neural networks for 3d shape recognition," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 945–953. 6, 12
- [60] T. Yu, J. Meng, and J. Yuan, "Multi-view harmonized bilinear network for 3d object recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 186–194. 6, 12
- [61] Y. Feng, Z. Zhang, X. Zhao, R. Ji, and Y. Gao, "Gvcnn: Group-view convolutional neural networks for 3d shape recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 264–272. 6, 12
- [62] S. Xie, R. B. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2017, pp. 5987–5995. 7, 9
- [63] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "ImageNet Large Scale Visual Recognition Challenge," *International Journal of Computer Vision (IJCV)*, vol. 115, no. 3, pp. 211–252, 2015. 7, 9
- [64] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao, "3d shapenets: A deep representation for volumetric shapes," in *CVPR*, 2015, pp. 1912–1920. 7, 12
- [65] R. Collobert, S. Bengio, and J. Mariethoz, "Torch: a modular machine learning software library," Tech. Rep., 2002. 7
- [66] I. J. Goodfellow, D. Warde-Farley, M. Mirza, A. C. Courville, and Y. Bengio, "Maxout networks," in *Proceedings of the 30th International Conference on Machine Learning*, 2013, pp. 1319–1327. 7
- [67] C. Lee, S. Xie, P. W. Gallagher, Z. Zhang, and Z. Tu, "Deeply-supervised nets," in *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*, 2015. 8, 9
- [68] K. He, X. Zhang, S. Ren, and J. Sun, "Identity mappings in deep residual networks," in *Proceedings of the European Conference on Computer Vision*, 2016. 8, 9, 15
- [69] L. Bo, X. Ren, and D. Fox, "Unsupervised feature learning for rgb-d based object recognition," in *Experimental Robotics*. Springer, 2013, pp. 387–402. 10, 11
- [70] C. Xu, D. Tao, and C. Xu, "Multi-view intact space learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 12, pp. 2531–2544, 2015. 11
- [71] M. Blum, J. T. Springenberg, J. Wülfing, and M. Riedmiller, "A learned feature descriptor for object recognition in rgb-d data," in *Proceedings of the IEEE International Conference on Robotics and Automation*. IEEE, 2012, pp. 1298–1303. 11
- [72] H. Zhu, J.-B. Weibel, and S. Lu, "Discriminative multi-modal feature fusion for rgb-d indoor scene recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2969–2976. 12
- [73] C. Wang, M. Pelillo, and K. Siddiqi, "Dominant set clustering and pooling for multi-view 3d object recognition," in *British Machine Vision Conference 2017, BMVC 2017, London, UK, September 4-7, 2017*, 2017. 12
- [74] A. Kanezaki, Y. Matsushita, and Y. Nishida, "Rotationnet: Joint object categorization and pose estimation using multiviews from unsupervised viewpoints," in *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, 2018, pp. 5010–5019. 12
- [75] E. Johns, S. Leutenegger, and A. J. Davison, "Pairwise decomposition of image sequences for active multi-view recognition," in *CVPR*. IEEE Computer Society, 2016, pp. 3813–3822. 12

APPENDIX A
GRADIENTS OF THE PROPOSED NEURON-WISE
CORRELATION-MAXIMIZING REGULARIZER

We use multi-variable chain rule to derive the gradients of the neuron-wise regularizer w.r.t. the weight vectors $\frac{\partial Corr}{\partial \mathbf{w}_1}$ and $\frac{\partial Corr}{\partial \mathbf{w}_2}$, and also w.r.t. the input features $\frac{\partial Corr}{\partial \mathbf{x}_1^i}$, $\frac{\partial Corr}{\partial \mathbf{x}_2^i}$, $i = 1, \dots, m$. Their explicit forms are presented as follows (with no simplification).

$$\frac{\partial Corr}{\partial \sigma_1^2} = -\frac{1}{2}(\sigma_1^2 \sigma_2^2 + \epsilon)^{-\frac{3}{2}} \cdot \sigma_2^2 \cdot \sum_{i=1}^m (y_1^i - \mu_1)(y_2^i - \mu_2)$$

$$\frac{\partial Corr}{\partial \sigma_2^2} = -\frac{1}{2}(\sigma_1^2 \sigma_2^2 + \epsilon)^{-\frac{3}{2}} \cdot \sigma_1^2 \cdot \sum_{i=1}^m (y_1^i - \mu_1)(y_2^i - \mu_2)$$

$$\frac{\partial Corr}{\partial \mu_1} = \frac{-\sum_{i=1}^m (y_2^i - \mu_2)}{\sqrt{\sigma_1^2 \sigma_2^2 + \epsilon}} + \frac{\partial Corr}{\partial \sigma_1^2} \cdot (-2) \sum_{i=1}^m (y_1^i - \mu_1)$$

$$\frac{\partial Corr}{\partial \mu_2} = \frac{-\sum_{i=1}^m (y_1^i - \mu_1)}{\sqrt{\sigma_1^2 \sigma_2^2 + \epsilon}} + \frac{\partial Corr}{\partial \sigma_2^2} \cdot (-2) \sum_{i=1}^m (y_2^i - \mu_2)$$

$$\frac{\partial Corr}{\partial y_1^i} = \frac{(y_2^i - \mu_2)}{\sqrt{\sigma_1^2 \sigma_2^2 + \epsilon}} + \frac{\partial Corr}{\partial \mu_1} \cdot \frac{1}{m} + \frac{\partial Corr}{\partial \sigma_1^2} \cdot 2(y_1^i - \mu_1)$$

$$\frac{\partial Corr}{\partial y_2^i} = \frac{(y_1^i - \mu_1)}{\sqrt{\sigma_1^2 \sigma_2^2 + \epsilon}} + \frac{\partial Corr}{\partial \mu_2} \cdot \frac{1}{m} + \frac{\partial Corr}{\partial \sigma_2^2} \cdot 2(y_2^i - \mu_2)$$

$$\frac{\partial Corr}{\partial \mathbf{x}_1^i} = \frac{\partial Corr}{\partial y_1^i} \cdot \mathbf{w}_1 \quad \frac{\partial Corr}{\partial \mathbf{x}_2^i} = \frac{\partial Corr}{\partial y_2^i} \cdot \mathbf{w}_2$$

$$\frac{\partial Corr}{\partial \mathbf{w}_1} = \sum_{i=1}^m \frac{\partial Corr}{\partial y_1^i} \cdot \mathbf{x}_1^i \quad \frac{\partial Corr}{\partial \mathbf{w}_2} = \sum_{i=1}^m \frac{\partial Corr}{\partial y_2^i} \cdot \mathbf{x}_2^i$$

APPENDIX B
NETWORK ARCHITECTURES

Our LeNet variant in Section V-A consists of 3 conv layers (the first one is the input layer), followed by 3 FC layers (the last one is the output layer). The first two conv layers have 32 filters, and the third one has 64 filters. All the three conv layers have filters of size 5×5 and stride 1. Max pooling of size 3×3 is applied after the first conv layer, and average pooling of size 3×3 is applied after both the second and third conv layers. The numbers of neurons for the three FC layers are respectively 256, 64, and 10.

Our used ResNet in Section V-A1 follows [28], [68]. In particular, we first build a ConvNet that starts with a conv layer of $16 \ 3 \times 3$ filters, and then sequentially stacks three types of $2X$ conv layers of 3×3 filters, each of which has the feature map sizes of 32, 16, and 8, and filter numbers 16, 32, and 64, respectively. Spatial sub-sampling of feature maps is achieved by conv layers of stride 2. The ConvNet ends with a global average pooling and FC layers, with $6X + 2$ weight layers in total. Based on this ConvNet, we do (1) using an “identify

TABLE XIII
LAYER SPECIFICS OF CONVNETS IN FIGURE 5.

Layer	Filter size/Filter no./ Stride/Padding
conv1	$7 \times 7/96/2/3$
conv2	$5 \times 5/96/1/2$
conv3	$3 \times 3/112/1/1$
conv4	$3 \times 3/128/1/1$
conv5	$3 \times 3/128/1/1$
fc6	-/1024/-/-
fc7	-/512/-/-
fc8	-/51/-/-
(max) pool1, pool2, pool3	$2 \times 2/-/2/1$
lower-conv	$5 \times 5/192/1/2$
middle-conv	$3 \times 3/256/1/1$
upper-fc	-/1024/-/-

shortcut” to connect every two conv layers of 3×3 filters, and a “projection shortcut” when sub-sampling of feature maps is needed; (2) changing it to the pre-activation version according to [68]. We set $X = 11$ that gives 68 weight layers.

Layer specifics of the RGB ConvNet, Depth ConvNet, and RGB-Depth ConvNets used in Section V-B1 (Figure 5) are presented in Table XIII.