

COMPRESSIVE VIDEO CLASSIFICATION IN A LOW-DIMENSIONAL MANIFOLD WITH LEARNED DISTANCE METRIC

George Tzagkarakis¹, Grigorios Tsagkatakis², Jean-Luc Starck¹ and Panagiotis Tsakalides²

¹ Commissariat à l'Énergie Atomique (CEA), Centre de Saclay, F-91191 Gif-Sur-Yvette cedex, France

² Foundation for Research & Technology - Hellas (FORTH-ICS), Crete, Greece

e-mail: {georgios.tzagkarakis, jstarck}@cea.fr, {greg, tsakalid}@ics.forth.gr

ABSTRACT

In this paper, we introduce an architecture for addressing the problem of video classification based on a set of compressed features, without the need of accessing the original full-resolution video data. In particular, the video frames are acquired directly in a compressed domain by means of random projections associated with a set of compressive measurements. This initial dimensionality reduction step is followed by distance metric learning for the construction of an informative distance matrix, which is then embedded in a manifold learning approach to increase the discriminative power of the random measurements in a lower-dimensional space. Classification results using a set of activity videos suggest that the proposed approach can be used effectively in cases when the acquisition and processing of full-resolution video data is characterized by increased consumption of the available power, memory and bandwidth, which may impede the operation of systems with limited resources.

Index Terms— Compressive video classification, distance metric learning, manifold learning

1. INTRODUCTION

Recent technological advances in the design of low-cost high-definition imaging devices and the availability of large digital databases with high-resolution video content, necessitates the efficient representation of an ever increasing volume of data in a precise and compact way to be exploited in carrying out tasks, such as detection and classification. In particular, automatic classification of video information is crucial for monitoring, indexing and retrieval purposes, while the advent of efficient computational models could be highly beneficial in cases when added constraints, such as power, storage and bandwidth, are imposed on the acquisition device. This is, for instance, the case in a remote sensing scenario, such as the use of unmanned aerial vehicles (UAVs) and terrestrial sensor networks in surveillance and reconnaissance applications,

with potentially limited resources of the acquisition hardware.

Besides, the increasing use of web-based video content during the last years motivated a growing interest on designing algorithms for video classification. For instance, in [1], motion and color features with Hidden Markov Models were used to classify sports videos, while in [2] audio features were extracted and used along with a multi-layer perceptron for real-time video classification. In both cases, the feature extraction step employs the original video content, which can be of very high dimension when the system operates at high resolution. Thus, dimensionality reduction arises naturally prior to further processing and decision making.

Techniques like *principal component analysis* (PCA), *independent component analysis* (ICA) and *linear discriminant analysis* (LDA) [3] have been employed widely in the framework of signal classification. The purpose of all these methods is the representation of the salient information in a low-dimensional space resulting in an improved classification performance. However, they suffer from certain limitations, such as data dependence and linearity restrictions. On the other hand, recently, there have been advances in the machine learning and pattern recognition communities for developing manifold learning algorithms to identify non-linear structures in low-dimensional manifolds from sample data points embedded in high-dimensional spaces, while preserving geometric distances and local neighborhood structures [4].

Operating at Shannon/Nyquist sampling rates may be excessive, or highly inefficient in case of limited-resource sensing systems, especially when signal processing tasks other than reconstruction, such as classification and detection are of interest. Recent works in the framework of *compressive sensing* (CS) [5] revealed that the task of classification can be performed efficiently using a highly reduced amount of data-independent linear incoherent random projections. These compressed measurements have been shown to preserve the meaningful information of the acquired signal [6], as well as the manifold structure of the acquired data [7].

Motivated by the success of random projections for face recognition [8], in a recent work [9] we addressed the video classification problem by working directly in the compressed

This work is supported by CS-ORION Marie Curie Industry - Academia Partnerships and Pathways (IAPP) project funded by the European Commission in FP7 (PIAP-GA-2009-251605).

domain. More specifically, a signature was generated for each video by projecting each frame on a random measurement matrix, and then the classification was performed by employing typical methods, such as, SVM and k NN, on these compressed signatures. Despite the satisfactory performance of the proposed computationally efficient approach, there were cases where the discriminative power of random measurements was not sufficient enough.

In this study, the video classification problem is also addressed by considering the futuristic scenario of a sensing system equipped with a *single-pixel camera* [10] for the acquisition of CS measurements. The discriminative capability of the compressed measurements is now enhanced by employing a distance metric learning step on the random measurements to estimate a distance metric, which will enhance the inter-class separability, thus improving the classification accuracy under a nearest-neighbor (NN) decision rule. Finally, the learned distance is embedded in a manifold learning procedure for the generation of an appropriate projection matrix. The classifier employs a simple k NN rule in the manifold space by utilizing the learned distance. The experimental evaluation reveals a clear improvement of the proposed approach when compared with our previous one. Moreover, we emphasize that the computational cost at the encoder is exactly the same, while it increases at the decoder due to the distance and manifold learning processes prior to classification. However, this does not affect the overall system’s performance, since increased processing resources are available at the decoder.

The paper is organized as follows: in Section 2, the representation of a given video sequence in a CS measurements domain is introduced. Section 3 describes in detail the proposed compressive video classification system, followed by an experimental evaluation of the classification accuracy in Section 4. Finally, conclusions and directions for further extensions are outlined in Section 5.

2. COMPRESSED VIDEO ACQUISITION

In the following, let $\mathbf{V} = \{\mathbf{x}_1, \dots, \mathbf{x}_F\}$ denote a video sequence consisting of F frames. We consider the case of acquiring grayscale video frames of size $N_r \times N_c$, where each frame is represented by the *luminance* component. Then, a set of compressed measurements is obtained by means of a single-pixel camera with a digital micromirror device (DMD) array generating the random basis patterns.

Let $\mathbf{x}_j \in \mathbb{R}^{N_r \times N_c}$, $j = 1, \dots, F$, denote the luminance component for the j -th frame. Then, a vector of compressed measurements is generated as follows:

$$\mathbf{g}_j = \Phi_j \mathbf{x}_j, \quad (1)$$

where the measurement matrix $\Phi_j \in \mathbb{R}^{M \times N}$ with $M < N$. We note that in the above equation the frame is viewed as a column vector with $N = N_r \times N_c$ elements, while also in general a different measurement matrix can be used for each frame. However, for simplicity, we consider that the same

measurement matrix $\Phi_j \equiv \Phi$ is used to capture all the frames of a given video sequence.

Common choices for Φ are random matrices with independent, identically distributed (i.i.d.) Gaussian or Bernoulli entries. However, their use in case of high-dimensional data can be memory and computationally intensive, thus increasing the burden of a system with limited resources. A family of matrices admitting a “hardware-friendly” implementation are the so-called *structurally random matrices* [11]. The block Walsh-Hadamard (BWHT) operator, which is a typical member of this family, is employed in the proposed architecture.

It is emphasized once again that the use of a single-pixel camera yields directly a set of measurements in the compressed domain without the need for accessing the original frames at full resolution, thus reducing significantly the processing and storage expenses of the sensing device.

In compressive video classification (CVC), we consider that the given video belongs to class $c \in \{1, \dots, C\}$. Working in a supervised learning framework, a set of training samples is obtained for each class, $\mathcal{T}_c = \{\mathbf{V}_1^c, \dots, \mathbf{V}_Q^c\}$. For simplicity, an equal number of training samples Q is considered for all the classes. Let also $\mathcal{T} = \{\mathcal{T}_1, \dots, \mathcal{T}_C\}$ be the overall set of training samples. Then, the CVC problem is stated as follows: *Given a low-dimensional signature of the acquired video, a training dictionary \mathbf{P} , and a measurement matrix Φ , estimate the correct class $c \in \{1, \dots, C\}$.*

3. PROPOSED CVC ARCHITECTURE

The generation of random measurements serves as a first dimensionality reduction step, which maps directly the original full-resolution video data into a lower-dimensional compressed domain. The advantage of this process, as opposed to other commonly used dimensionality reduction approaches, such as PCA, is that its linear character offers fast and efficient computations, while also, and most importantly, it is data independent. This allows for a better generalization, which is important, especially when we deal with very large databases.

Moreover, if the measurement matrix employed in (1) satisfies a restricted isometry property (RIP) [5], then, it is guaranteed with high probability that the relative distances between the samples (frames) in the original space are preserved in the lower-dimensional compressed space. This, in turn, affects the complexity of algorithms that depend on the dimensionality of the input data. Moreover, these key properties of the random projections were generalized for signal manifolds [7]. More specifically, it has been shown that random projections preserve the metric structure of a manifold, that is, the set of pairwise geodesic distances, as well as its curvature.

However, the generation of compressed measurements for each frame of a given high-resolution video sequence may still result in a large amount of data even for relatively small sampling ratios, while also the selection of an appropriate distance metric, which is fundamental to any learning algorithm, is highly problem-dependent and determines the success or

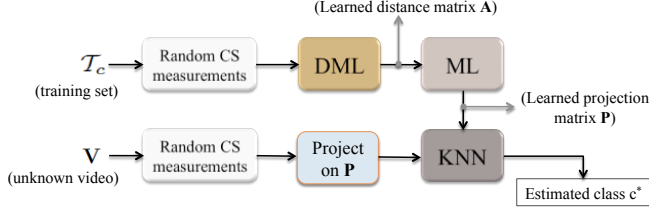


Fig. 1. Proposed CVC architecture.

failure of the classification process.

We overcome the above drawbacks by employing a *distance metric learning* (DML) combined with a *manifold learning* (ML) approach. In particular, DML is used in the domain of compressed measurements to estimate a distance (or equivalently a linear transformation) that satisfies the class constraints, while ML acts as a second dimensionality reduction step, where the final classification of a new video is performed by measuring the distance between the manifold embedded training samples and the new data. The overall architecture of the proposed CVC system is shown in Fig. 1.

A typical classification system consists of two main phases, namely, a *training phase*, where a more compact representation of the original information is generated in a low-dimensional space, with the goal of preserving a high discriminative power, and a *classification phase*, where the extracted feature vector of a new sample is compared with the corresponding features of the training samples by means of a suitable similarity criterion resulting in the estimated class.

3.1. Training phase

As Fig. 1 shows, the training phase consists of three distinct steps: i) initial dimensionality reduction through random CS measurements, ii) learning of a suitable distance metric, and iii) second dimensionality reduction via manifold learning.

3.1.1. Random CS measurements

During this step, a set of CS measurements is generated for the frames of each training video. More specifically, let $\mathbf{x}_{j,q}^c$ denote the j -th frame of the q -th video belonging in class c , with $j = 1, \dots, F$, $q = 1, \dots, Q$ and $c = 1, \dots, C$. Then, a low-dimensional (feature) measurement vector $\mathbf{g}_{j,q}^c \in \mathbb{R}^{M \times 1}$ is assigned to $\mathbf{x}_{j,q}^c$ as given by (1). The overall signature for the q -th video of class c is given by

$$\mathbf{V}_q^c \mapsto \mathcal{S}_q^c = \{\mathbf{g}_{1,q}^c, \dots, \mathbf{g}_{F,q}^c\}, \quad (2)$$

and by augmenting all the training signatures we get the overall signature for the training dataset, $\mathcal{S} = \{\mathcal{S}_q^c\}_{q=1, \dots, Q}^{c=1, \dots, C}$.

3.1.2. Distance metric learning

Depending on the problem setup, selecting an appropriate distance metric is crucial for the performance of many learning

algorithms, such as the k NN and the k -means. Typically, Euclidean distance may work empirically. On the other hand, the family of Mahalanobis distances have been shown to generalize well the standard Euclidean distance by admitting arbitrary linear scalings and rotations of the feature space, resulting in a more meaningful connection between the data. Furthermore, additional information, such as class labels, can be incorporated in the learned metric. For two data points $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^N$ the squared Mahalanobis distance is parameterized by a positive definite matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ as follows,

$$d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{A} (\mathbf{x}_i - \mathbf{x}_j). \quad (3)$$

The objective in DML in the supervised case is to learn a new distance that will satisfy the pairwise constraints imposed by class label information. Formally, for two vectors \mathbf{x} and \mathbf{y} , the learned distance should satisfy $d_{\mathbf{A}}(\mathbf{x}, \mathbf{y}) \leq l$ if $label(\mathbf{x}) = label(\mathbf{y})$ and $d_{\mathbf{A}}(\mathbf{x}, \mathbf{y}) \geq u$ if $label(\mathbf{x}) \neq label(\mathbf{y})$.

In the proposed CVC architecture, DML is carried out by means of a recently introduced information-theoretic metric learning (ITML) algorithm [12], which reduces the problem of estimating the distance matrix \mathbf{A} into a problem of minimizing the differential relative entropy between two multivariate Gaussians under similarity or dissimilarity constraints and pairwise relations of distances. Unlike most of the existing methods, no eigenvalue decompositions or semi-definite programming is required, which makes ITML fast and scalable, while it also generalizes well to unseen test data. ITML takes as inputs the signature \mathcal{S} with the training CS measurements, along with their corresponding class labels, and gives an estimate of the distance matrix \mathbf{A} as an output.

3.1.3. Manifold learning

Manifold learning (or non-linear dimensionality reduction), which is the counterpart to PCA, consists of finding a low-dimensional representation of the input data, while attempting to preserve the local metric structure. Graph-based learning methods, such as Isomap [13] and local linear embedding (LLE) [14], employ Euclidean distances to describe local neighborhoods of the data. However, both methods fail in cases of data sets with high curvature, self-intersections, or non-convex sampling within the manifold space. Furthermore, manifold learning methods such as Isomap and LLE lack a straightforward extension of the learned mapping to new data, a problem called *out-of-sample extension*, which limits the generalization applicability of these methods. Moreover, the unsupervised nature of both methods leads in a limited generalization to new data, due to the direct mapping of input data to the manifold space, instead of a mapping from the input space to the manifold space.

Locality preserving projections (LPP) [15] was among the methods introduced to overcome the above drawbacks by performing linear approximations of a non-linear manifold in small neighborhoods. Although LPP was designed originally

for unsupervised dimensionality reduction, it was extended to the supervised case where the class label information is used to construct the weight matrix for the embedding.

Regarding classification, this task is independent of the preservation of manifold structure, which is the purpose of manifold learning, since data points belonging to distinct classes may still be close in the embedding manifold. This observation necessitates the use of a supervised approach, along with an appropriately selected distance metric for the extraction of the local neighborhood structure. In the proposed CVC setup, this information is exploited by employing the distance matrix \mathbf{A} learned from the training CS measurements, as described in the previous section. The final output of LPP is a projection matrix \mathbf{P} with the eigenvectors of the adjacency graph Laplacian.

3.2. Classification phase

When a new video is given as input to the proposed CVC system, it is represented directly by its corresponding set of CS measurements, due to the assumption that our acquisition system is equipped with a single-pixel camera. Let $\mathcal{S}_u = \{\mathbf{g}_{1,u}, \dots, \mathbf{g}_{F,u}\}$ be the matrix of CS measurements for the F frames of the unknown video \mathbf{V}_u . As a second step, the compressed measurements are projected on the learned manifold by multiplying with the learned projection matrix, that is, $\mathcal{S}_u \mapsto \mathbf{P}^T \mathcal{S}_u$. Then, a k NN classifier is applied for each projected vector $\mathbf{g}_{j,u}$ of the signature of \mathbf{V}_u in the trained manifold. Finally, the estimated class is the one with the highest frequency of appearance among the individually classified projected vectors.

4. EXPERIMENTAL EVALUATION

The classification performance of the proposed CVC system is evaluated on a subset of the UCF-50 dataset¹ consisting of videos categorized in 8 classes according to different activities, namely, “Basketball”, “Clean & jerk”, “Guitar”, “Piano”, “Rock climbing indoor”, “Rowing”, “Skiing” and “Tennis”. This dataset is considered to be particularly challenging due to large variations in camera motion and illumination conditions, object appearance and pose. In the present experimental setup, each class contains 30 videos of 50 frames per video, rescaled at 128×128 pixels for computational simplicity.

A distinct BWHT measurement matrix Φ is used for the acquisition of compressed measurements in each one of 50 Monte-Carlo runs. The sampling ratio r varies in $[0.005, 0.019]$ (or equivalently the number of CS measurements per frame, $M = \text{round}(rN)$, varies in $[82, 311]$). Regarding the manifold learning step, the dimension of the manifold varies as a percentage p of the number of CS measurements with $p \in [0.01, 0.15]$, which is equivalent to keeping the L most significant eigenvectors of the projection

¹<http://www.computervisiononline.com/dataset/ucf50-action-recognition-realistic-videos>

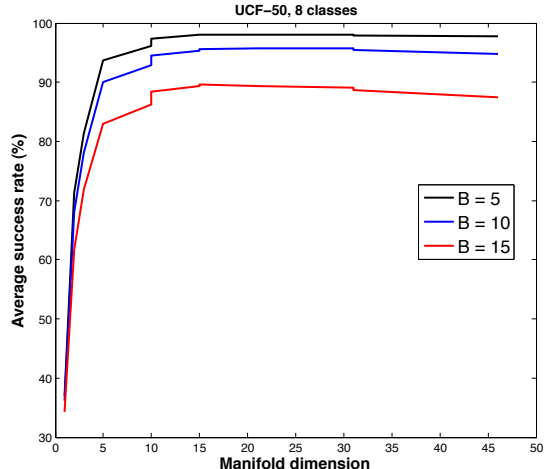


Fig. 2. Average classification rates as a function of the manifold dimension, for a varying number of testing samples B .

matrix \mathbf{P} with $L = \text{round}(pM)$. The combination of all possible values of r , M , and p results in a manifold dimension which varies in $[1, 47]$. The ITML algorithm is applied first over the whole training set to estimate the distance matrix \mathbf{A} .

For the classification, a k NN classifier with $k = 1$ is used, while the classification accuracy is computed for a varying training and testing partition. More specifically, in each Monte-Carlo run the overall dataset is divided randomly in two subsets, a training and a testing one, consisting of $(30 - B)$ and B videos, respectively, with $B \in \{5, 10, 15\}$. In the following, the classification accuracy is expressed in terms of the average success rate, which is defined as the ratio of the number of correctly classified frames over the total number of query frames, where the average is taken over the 50 Monte-Carlo runs and all testing frames.

Fig. 2 shows the average success classification rates for the three partitions of the video database, as a function of manifold dimension. As we expected, the successful classification rate increases as the number of training samples increases (smaller B), while the system is able to achieve its optimum performance enough at a small number of dimensions. These results are also consistent with the theoretical perspective, which states that random projections using a measurement matrix satisfying a RIP condition, along with an appropriate embedding in a low-dimensional manifold, preserve the local metric structure, which is expressed by the ability of a nearest-neighbor classifier to discriminate the embedded features between distinct classes. Moreover, we can see that in the very low-dimensional regime the low classification accuracy is not improved by increasing the number of training samples, since the corresponding learned manifold is unable to represent the original information content of the higher-dimensional samples.

Furthermore, Table 1 shows the average confusion matrix between the 8 classes for a sampling rate $r = 0.01$ and $B = 5$

testing samples per class, for the proposed method, along with our previous one introduced in [9] (average success rates are shown in parentheses). In contrast to our previous CVC system based solely on the CS measurements, the classification accuracy of the one proposed here is significantly higher for all classes, while the larger misclassification errors observed in some cases, such as the (Skiing, Basketball) and (Skiing, Rowing) pairs, are again mainly due to the use of the luminance component only. For instance, in the first case there is a similarity between the relatively smooth surfaces of the snowy terrain and the courts, which is hardly distinguishable when we exclude the color information. This also highlights the importance of extracting features, which are as representative as possible, as we also suggest in the last section.

	Basketball	Clean & jerk	Guitar	Piano	Rock climbing	Rowing	Skiing	Tennis
Basketball	87.15 (64.69)	0.03 (3.61)	0.66 (2.44)	0.25 (4.21)	0.90 (2.18)	1.48 (6.23)	0.68 (5.95)	2.38 (8.92)
Clean & jerk	1.23 (4.79)	95.98 (69.09)	0.14 (5.37)	0 (2.76)	0.02 (3.37)	0.47 (5.41)	0.24 (4.78)	0 (1.31)
Guitar	0.70 (1.30)	0 (0.84)	98.67 (70.55)	0 (2.82)	0 (1.97)	0.02 (3.89)	0 (2.37)	0 (1.92)
Piano	0.77 (2.54)	0.09 (1.08)	0 (1.60)	97.90 (70.97)	1 (3.73)	0.01 (2.62)	0 (1.82)	0.12 (3.55)
Rock climbing	3.15 (5.32)	0.70 (6.51)	0.11 (4.03)	0.09 (3.32)	94.98 (68.53)	0.09 (3.15)	2.25 (6.99)	0 (2.59)
Rowing	2.57 (6.34)	2.1 (8.64)	0.20 (6.02)	0.33 (5.75)	0.30 (3.87)	93.96 (66.50)	3.55 (8.60)	0.59 (5.13)
Skiing	3.75 (13.52)	0.90 (9.87)	0.22 (8.41)	1.43 (6.55)	2.63 (14.05)	3.95 (8.42)	92.73 (64.74)	2.08 (6.78)
Tennis	0.68 (1.50)	0.20 (0.36)	0 (1.58)	0 (3.62)	0.17 (2.30)	0.02 (3.78)	0.55 (4.75)	94.83 (69.80)

Table 1. Average confusion matrix for $r = 0.01$, $B = 5$ for the proposed CVC system, and our previous one based solely on CS measurements (success rates shown in parentheses).

5. CONCLUSIONS

In this study, a compressive video classification method was introduced. The design of the proposed CVC system was based on the assumption of an imaging system with limited resources, without having access to the full-resolution frames, where the video data are captured directly in the CS domain using a single-pixel camera. A distance metric learning approach followed the initial dimensionality reduction in the compressed domain, to learn the relative distances between the training samples in a supervised way by exploiting the prior class label information. As a further dimensionality reduction step, while preserving the local metric structure, a manifold learning method was used by exploiting the distance matrix learned in the previous step. Finally, the estimated class for a query video was obtained using a simple

nearest-neighbor classifier in the low-dimensional manifold. The experimental results revealed a high classification accuracy, even at very low sampling rates and manifold dimension, which are significantly smaller than the rates and dimensions required for solving the problem of sparse reconstruction.

In the proposed CVC framework we do not exploit the sparsity of the original video data in an appropriate transform domain (*e.g.*, DCT, DWT), which is at the core of CS theory. As a direct extension, we expect that the sparsification of the original data before their embedding in the low-dimensional manifold, by learning a sparsifying dictionary, could enhance the discriminative power of the generated features, and consequently the classification success rate. In addition, the above could be combined with the generation of CS features by considering the color information, which can also increase the classification margin among the several classes.

6. REFERENCES

- [1] X. Gibert, H. Li, and D. Doermann, "Sports video classification using HMMs," in *Proc. IEEE Int. Conf. Multimedia & Expo (ICME)*, Vol. 2, Baltimore, July 2003.
- [2] M. Rouvier, G. Linares, and D. Matrouf, "On-the-fly video genre classification by combination of audio features," in *Proc. IEEE Int. Conf. on Ac., Speech & Signal Proc. (ICASSP)* Mar. 14–19, Dallas, TX, 2010.
- [3] C. Bishop, "Pattern recognition and machine learning." Springer, 2006.
- [4] R. Pless and R. Souvenir, "A survey of manifold learning for images," *IPSIJ Trans. on Comp. Vision and Applications*, Vol. 1, pp. 83–94, 2009.
- [5] E. Candès, J. Romberg, and T. Tao, "Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information," *IEEE Trans. on Inf. Th.*, Vol. 52 (2) pp. 489–509, Feb. 2006.
- [6] C. Hegde *et al.*, "Efficient machine learning using random projections," *NIPS Work. on Eff. Machine Learn.*, Whistler, Canada, Dec. 2007.
- [7] C. Hegde, M. Wakin, and R. Baraniuk, "Random projections for manifold learning," *NIPS Work. on Efficient Machine Learning*, Whistler, Canada, Dec. 2007.
- [8] G. Tsagkatakis and A. Savakis, "Face recognition using sparse representations and manifold learning," in *Proc. Int. Symp. on Visual Comp. (ISVC)*, Nevada, NV, 2010.
- [9] G. Tzagkarakis *et al.*, "Compressive video classification for decision systems with limited resources," in *Proc. Picture Coding Symp. (PCS'12)*, Krakow, Poland, May 7–9, 2012.
- [10] <http://dsp.rice.edu/cscamera>
- [11] T. Do, T. Tran, and L. Gan, "Fast compressive sampling with structurally random matrices," in *Proc. IEEE Int. Conf. on Acoustics, Speech & Signal Proc. (ICASSP)*, Mar. 30–April 4, Las Vegas, NV, 2008.
- [12] J. Davis *et al.*, "Information-theoretic metric learning," in *Proc. 24th Int. Conf. on Machine Learning*, Corvallis, OR, 2007.
- [13] J. Tenenbaum, V. Silva, and J. Langford, "A global geometric framework for nonlinear dimensionality reduction," *Science*, Vol. 290 (5500), pp. 2319–2323, Dec. 2000.
- [14] S. Roweis and L. Saul, "Nonlinear dimensionality reduction by locally linear embedding," *Science*, Vol. 290 (5500), pp. 2323–2326, Dec. 2000.
- [15] X. He and P. Niyogi, "Locality preserving projections," in *Proc. 17th Ann. Conf. Neural Inf. Proc. Systems (NIPS)*, Dec. 11–13, British Columbia, CA, 2003.