

Video Trajectory Classification and Anomaly Detection Using Hybrid CNN-VAE

Kelathodi Kumaran Santhosh, *Student Member, IEEE*, Debi Prosad Dogra, *Member, IEEE*, Partha Pratim Roy, and Adway Mitra

Abstract—Classifying time series data using neural networks is a challenging problem when the length of the data varies. Video object trajectories, which are key to many of the visual surveillance applications, are often found to be of varying length. If such trajectories are used to understand the behavior (normal or anomalous) of moving objects, they need to be represented correctly. In this paper, we propose video object trajectory classification and anomaly detection using a hybrid Convolutional Neural Network (CNN) and Variational Autoencoder (VAE) architecture. First, we introduce a high level representation of object trajectories using color gradient form. In the next stage, a semi-supervised way to annotate moving object trajectories extracted using Temporal Unknown Incremental Clustering (TUIC), has been applied for trajectory class labeling. Anomalous trajectories are separated using t-Distributed Stochastic Neighbor Embedding (t-SNE). Finally, a hybrid CNN-VAE architecture has been used for trajectory classification and anomaly detection. The results obtained using publicly available surveillance video datasets reveal that the proposed method can successfully identify some of the important traffic anomalies such as vehicles not following lane driving, sudden speed variations, abrupt termination of vehicle movement, and vehicles moving in wrong directions. The proposed method is able to detect above anomalies at higher accuracy as compared to existing anomaly detection methods.

Index Terms—Convolutional Neural Network, Deep Learning, Variational Autoencoder, Dirichlet Process Mixture Model, Visual Surveillance, Trajectory Classification, Traffic Anomaly Detection.

I. INTRODUCTION

Timely detection of traffic anomaly is one of the prerequisites of an Intelligent Transportation Systems (ITS). If not done timely, anomalies may create cascading effects leading to chaos in traffic. Typical examples of traffic anomalies are, lane driving violation, over-speeding, collision, red-light violation, etc. Anomaly detection using video object trajectories with deep learning has not yet been explored much. In this paper, we propose a color gradient approach for representing vehicular trajectories extracted from videos. These trajectories are then used for classification and anomaly detection at traffic junctions using a hybrid CNN-VAE architecture.

Most commonly used features for video guided scene understanding are trajectories. A trajectory is a time series data with object locations indexed in temporal order. Classifying trajectories using neural networks is not trivial due to variation in the data length. Key to the success of a time series signal

classification lies in finding an effective representation of the data. Neural networks-based classifiers need fixed size inputs. CNN, Long Short Term Memory (LSTM) and Recurrent Neural Network (RNN) have been used for time series classification [10], [13], [30]. However, time series data can be of varying length. . Therefore, classification of varying length data can be applied after preprocessing, e.g. converting them into fixed length data either by padding or subsampling. If the trajectory length variance is large, preprocessing is mandatory.

Video anomaly detection at traffic junctions is highly challenging due to its contextual nature. For example, when a signal turns green at a traffic junction, only a few of the paths or directions are allowed for vehicle movement. Any motion that violates direction, is assumed to be anomaly though such motions can be normal in a different context.

A. Related Work

Traditional features such as basis transform coding using wavelet and Fourier coefficients [4], time series mean and covariance [4], and symbolic representation [17] have been used for classification of time-series data using neural networks. Also, other models such as Deep Belief Networks (DBN) [26] have been used for human activity detection [23]. On the other hand, CNNs are primarily used in image classification [16], [31], activity recognition in videos [12], speech recognition [6], etc.

Long Short Term Memory networks (LSTMs) [11] are a special kind of Recurrent Neural Network (RNN) that can be used for handling sequential/time series data. Authors of [8], [22] have proposed a recurrent network connecting LSTMs to CNNs to perform action recognition and video classification, respectively. Donahue et al. [8] have tested the learned models for activity recognition, image description and video description. The work proposed in [28] has achieved the state-of-the-art performance in video classification by connecting CNNs and LSTMs under a hybrid deep learning framework. Sequential Deep Trajectory Descriptor (DTD) has been used for action recognition [25] from the video sequences. Deep Neural Network (DNN)-based trajectory classification has been applied on Global Positioning System (GPS) trajectories [9]. Dense feature trajectories used have been utilized for action recognition in videos [27]. The LSTM-based work proposed in [13] uses fixed size features to classify trajectories of surrounding vehicles at four way intersections based on LIDAR (LIght Detection And Ranging), GPS, and inertial measurement unit (IMU) measurements.

K. K. Santhosh, D. P. Dogra and A. Mitra are with School of Electrical Sciences, Indian Institute of Technology Bhubaneswar, Odisha, India e-mail: (sk47@iitbbs.ac.in, dpdogra@iitbbs.ac.in, adway@iitbbs.ac.in).

P. P. Roy is with the Department of Computer Science and Engineering, Indian Institute of Technology, Roorkee, India. e-mail:(proy.fcs@iitr.ac.in).



Fig. 1. Depiction of temporal characteristics using color gradient forms. (a) Trajectories from QMUL [18] junction dataset using color gradient representation. Temporal characteristics are similar at similar locations, marked as P, Q, and R. (b) Illustration of trajectory of a vehicle (A) that does not stop and another vehicle (B) that stops and proceeds. (c) Temporal characteristics are significantly different for these two objects.

Dense trajectories extracted using neural networks have also been used for action recognition in videos including classifying a person when walking, running, jumping [3], [27], etc. These methods cannot handle multiple actions present in a scene. However, in real life scenario, multiple objects can interact resulting more than one action within the scene. Training neural networks for action recognition can be challenging in presence of multiple activities. However, object trajectories extracted using traditional methods [1], [2], [21] can be used for learning the motion patterns using DNNs as they can automatically extract features from trajectories. The trained/learned model can then be used in classification and action recognition applications.

In this work, we encode video trajectories using a high-level representation, named color gradient, that embeds spatio-temporal information of the objects-in-motion. The high-level representation is then used for trajectory classification and anomaly detection using a hybrid CNN-VAE architecture.

B. Motivation and Contributions

Since accurate classification is the key to detect anomalies, a classifier that can handle time series data with length variations, has been preferred. Typical neural networks-based methods need fixed input size. Therefore, varying length trajectories cannot directly be used in such classifiers. Conventional methods such as the one proposed in [30] convert the varying length time series data into fixed size by sampling. This is similar to quantization, which leads to information loss. The question is: Why can't a trajectory represented using an image be given as an input to a classifier? However, trajectories representing movement of more than one object in between two locations may look visually similar when projected in 2D space. Such representations fail to preserve temporal relations between successive points of a trajectory. Encoding of time information in the form of color gradient (red \rightarrow violet) reveals, similar patterns produce similar color gradient as depicted in Fig.1(a). Similarly, the trajectories with possible anomalies exhibit different spatio-temporal characteristics as depicted in Fig.1(b). This has motivated us to propose the following:

- (i) A high-level representation of object trajectories using color gradient that encodes spatio-temporal information of trajectories of varying length.
- (ii) A semi-supervised labeling technique based on modified Dirichlet Process Mixture Model (mDPMM) [24] clustering to identify the trajectory classes.

- (iii) A method using t-Distributed Stochastic Neighbor Embedding (t-SNE) [20] to eliminate anomalous trajectories in the training data.
- (iv) Detection of traffic anomalies using a hybrid CNN-VAE architecture.

Rest of the paper is organized as follows. In Section II, we present the proposed methodology. Section III presents experimental results and Section IV presents conclusion.

II. METHODOLOGY

First we discuss the background of the terms and concepts used in the work. A scene represents the view captured using static camera. We use observation or data to represent a trajectory. A cluster is a collection of trajectories of similar characteristics. A class is a set of trajectories having some selected common characteristics. Here, a class typically represents a unique path in a scene. A model is a representation of a real-world phenomenon. Here, model represents the weight parameters of the trained neural networks. We assume a model can represent a scene. Reconstruction loss (of CNN-VAE architecture) represents a measure of deviation from the input. A typical anomaly represents deviation from the normal path. Some anomalies are known a-priori. For example, when a signal turns green at a traffic junction, only a few of the paths are allowed for vehicle motion. Any motion that conflicts/intersects the allowed path, is considered as known anomaly. However, some anomalies may not be present in the training data. We refer them to as unknown anomalies.

Object trajectories are obtained using [1], [24]. A trajectory (τ_i) can be represented using (1), where (x_l, y_l) represents the position of moving object at time t_l and L_i be its length. A cluster is a collection of trajectories of similar characteristics. A class is a set of trajectories having some selected common characteristics. It can be trajectories in the same lanes, trajectories following same route, etc.

$$\tau_i = \langle (x_0, y_0, t_0), (x_1, y_1, t_1), \dots, (x_{L_i}, y_{L_i}, t_{L_i}) \rangle \quad (1)$$

Traffic anomalies can be classified into two types; *known* and *unknown*. *Known* anomalies correspond to trajectories that may be allowed in different contexts. On the contrary, *unknown* anomalies correspond to trajectories that are not present in the training data. In order to detect both types of anomalies, it is important to learn the normal trajectory patterns or classes. The overall anomaly detection framework is presented in Fig. 2.

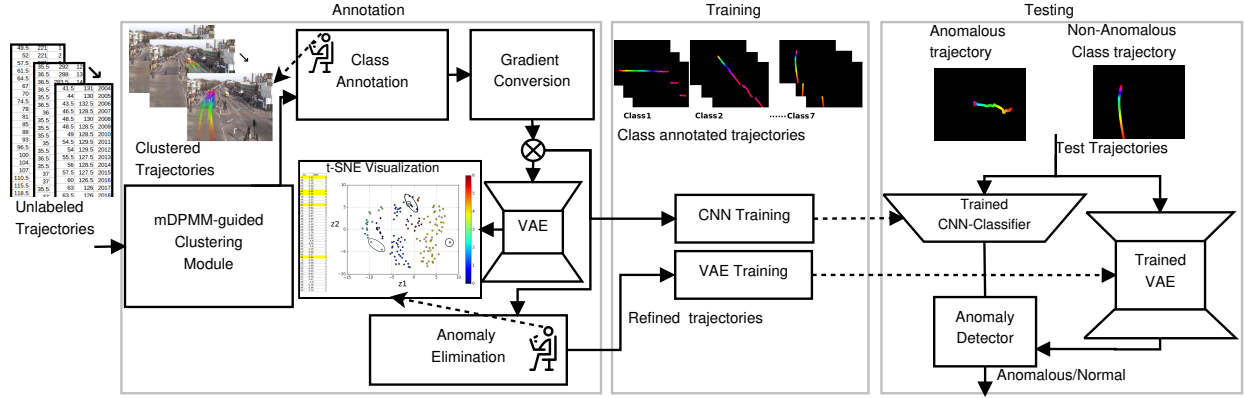


Fig. 2. Proposed anomaly detection framework. Unlabeled trajectories are grouped using the modified Dirichlet Process Mixture Model (mDPMM) [24]. Clusters are then mapped to different classes of trajectories using manual annotation and gradient representation is produced. These trajectories are then fed to train CNN and VAE to get t-SNE for eliminating anomalous trajectories (if any). The refined normal trajectories represented in color gradient form are fed to VAE to train the anomaly model of VAE. Anomaly detection is done once average reconstruction loss is known for the trained VAE with normal trajectories. Trained models are then used for classification and anomaly detection.

A. Background

1) *Modified DPMM Guided Clustering*: When raw trajectories are obtained from some tracking algorithms, they need to be clustered to identify different patterns. In [24], we have proposed a modified DPMM (mDPMM) to group pixels having similar characteristics. Here, we use mDPMM to group trajectories to learn the motion patterns. The model is expressed using (2 - 5).

$$z_i | \pi \sim \text{Discrete}(\pi) \quad (2)$$

$$\tau_i | z_i, \theta_k \sim F(\theta_{z_i}) \quad (3)$$

$$\pi | e^{-\beta} \sim \text{Dirichlet}(e^{-\beta}/K, \dots, e^{-\beta}/K) \quad (4)$$

$$\theta_k | H \sim H \quad (5)$$

Here, τ_i is a random variable representing the trajectory and z_i corresponds to the latent variable representing cluster labels. z_i takes one of the values from $k = 1 \dots K$, where K is the number of clusters. $\pi = (\pi_1, \dots, \pi_K)$, referred to as mixing proportion, is a vector of length K representing the probabilities of z_i to be k . θ_k is the parameter of cluster k and $F(\theta_{z_i})$ denotes the distribution defined by θ_{z_i} . $e^{-\beta}$ is the concentration parameter of Dirichlet distribution and its value decides the number of clusters formed. β is referred to as concentration radius. Trajectory clustering is to be done by taking τ_i as $\langle x_s, y_s, x_e, y_e, t_d \rangle$, where (x_s, y_s) represents the start position, (x_e, y_e) the end position and t_d is the duration/length of the trajectory.

Using the inference method given in [24], clustering of trajectories can be done. These clusters can be typically grouped into two types. First type contains large number of trajectories and they represent prominent patterns in the scene. The second type of clusters contain less number of trajectories. They can either correspond to less frequently occurring patterns or anomalies.

2) *Gradient Conversion of the Trajectories*: A trajectory in time series is mapped into a color gradient form by varying hue using $\text{hue}(x_l, y_l) = (t_l - t_0)/L_i * 180, 0 \leq l \leq L_i$ within an image frame. These gradient frames become inputs to the CNN and VAE.

3) *Anomaly Elimination in Training Data using t-SNE*: t-SNE [20] is a machine learning algorithm for visualizing high-dimensional data in a low-dimensional space. We use this for visualizing latent features of a trained VAE in two dimensions. Trajectories belonging to same class typically lie in close proximity in the visualization plane. However, trajectories that are far away from a class are inspected again for manual anomaly checking.

B. Trajectory Annotation

Suppose a set of trajectories captured from a traffic junction or road are given. These trajectories must belong to any one of the defined set of paths (classes). Applying mDPMM helps to identify prominent patterns from these trajectories. Like any unsupervised method, clustering algorithm can only identify different possible patterns from the trajectory data. Though prominent patterns can correspond to normal trajectories, clusters with less number of trajectories can represent a rare pattern or an anomaly. This necessitates to have an additional annotation process to identify allowed classes. Clustering reduces the load of the manual labeling process as an initial grouping is done through mDPMM. The annotator can identify these rare patterns through visual observation of the scene and separate the anomalous trajectories to finalize the allowed classes. This process is called class annotation.

More refinements are possible within a class. It is possible that two trajectories with similar endpoints and duration may follow different paths, out of which one may be normal. This may not always be detected through visual observation. Therefore, t-SNE has been used to visualize the distribution of trajectories within the classes. This helps to remove noises (anomalies) from the training set being prepared for VAE.

C. Training CNN and VAE Framework

A CNN classifier typically consists of repeated occurrences of cascaded convolution, activation, and pooling layers followed by fully connected layers. The architecture used in this work is depicted in Fig. 3(a). During the training stage,

a cost/loss function representing the cross-entropy between the expected and predicted class is minimized using Adam optimizer [14] with a learning rate of λ .

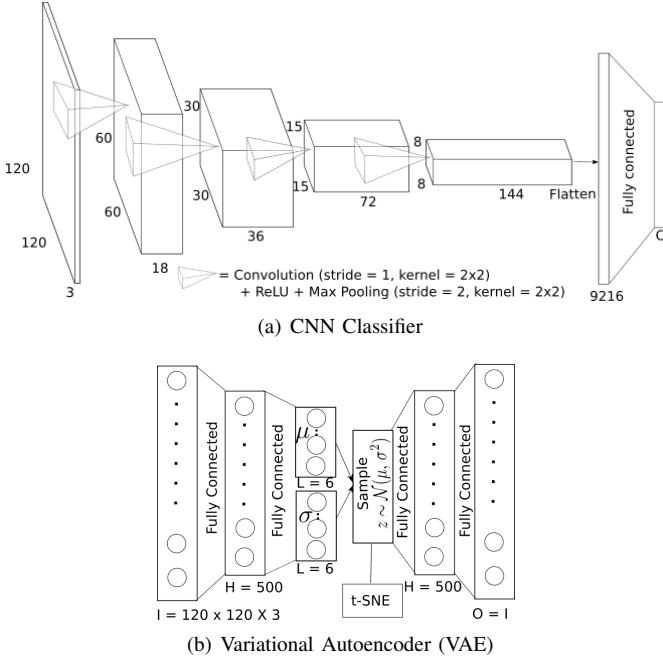


Fig. 3. CNN and VAE details. (a) Proposed CNN architecture for trajectory classification. (b) Autoencoder with the dimensions of each layer. I/O , H and L represent input/output, hidden and latent dimensions, respectively.

We use a variational autoencoder (VAE) similar to [15] to detect unknown anomalies. It typically consists of encoding and decoding stages. Input to the encoder $q_\theta(z|\tau)$ is τ . Output is a hidden/latent feature z , where θ represents weights and biases of encoder network. Decoder $p_\phi(\tau|z)$ takes latent feature z and regenerates τ , where ϕ represents weights and biases of decoder. Loss function (l_i) for a trajectory τ_i is given in (8) in terms of log likelihood (ll) as given in as given (6) and Kullback-Leibler Divergence (KLD) as given in as given (7). Adam optimizer minimizes the average loss function during training. Once trained, VAE can detect anomalies using the average reconstruction loss on the trained VAE.

$$ll = \mathbb{E}_{z \sim q_\theta(z|\tau_i)} [\log p_\phi(\tau_i|z)] \quad (6)$$

$$KLD = q_\theta(z|\tau_i) \| p(z) \quad (7)$$

$$l_i(\theta, \phi) = -ll + KLD \quad (8)$$

D. Anomaly Detection

Classification is performed on the trained CNN using test trajectories represented in gradient form to obtain class c . Let δ be the threshold of reconstruction loss value for normal classes on the trained VAE. δ is derived using the variance of loss values on the training trajectories. A trajectory can be considered anomalous when $c \notin A_s$ or $l_i(\theta, \phi) > \delta$, such that A_s is a set of allowed trajectory classes of a particular signal s .

However, a classifier is needed for anomaly detection to handle conflicting trajectories. In a typical traffic junction, a

set of flows may be allowed at a given time. For example, the QMUL dataset (Fig. 4) suggests, any two flows, e.g, south-to-north on left side and north-to-south on right side, are allowed at a given time. Any other movements can be termed anomalous though individually such movements may be allowed at a different time. VAE cannot detect such known anomalies. Therefore, CNN helps to detect such conflicting anomalies. It also helps to identify the anomalous path.

III. EXPERIMENTAL RESULTS

We have used tensorflow and openCV for developing the classification and anomaly detection framework. We have used three datasets, namely T15 [29], QMUL [18] and a junction video dataset (referred to as 4WAY). Context tracker [7] has been used for creating trajectories from QMUL dataset, and Temporal Unknown Incremental Clustering (TUIC) [24] has been used for obtaining 4WAY trajectories. Inputs to CNN-VAE are resized to $120 \times 120 \times 3$. CNN training has been completed with 50 epochs with a learning rate of $\lambda = 1e^{-3}$ on T15 and 4WAY dataset videos. A batch size of 20 has been used for the QMUL dataset. VAE for T15 has been trained with $\lambda = 5e^{-4}$ in 500 epochs using a batch size of 20. VAE for QMUL dataset has been trained with $\lambda = 1e^{-4}$ and batch size = 10 in 500 epochs.

A. Experiments on Trajectory Clustering and Annotation

TABLE I
DATASET ANNOTATION RESULTS. T REPRESENTS THE VIDEO DURATION, N THE NUMBER OF TRAJECTORIES, K THE NUMBER OF CLUSTERS OBTAINED USING mDPMM, C THE NUMBER OF VALID CLASSES AND N_A THE NUMBER OF ANOMALOUS TRAJECTORIES. T15 IS A LABELED DATASET.

Dataset	β	T (min)	N	K	C	N_A
T15	-	-	1500	NA	15	31
QMUL	180	10	166	23	7	17
4WAY	100	28	3861	193	18	808

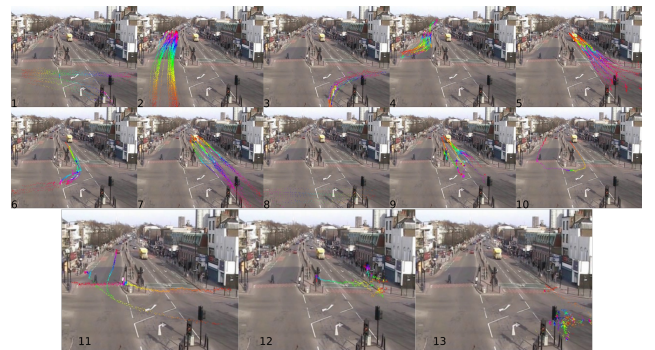


Fig. 4. Illustration of unsupervised clustering using mDPMM on QMUL trajectories. First two rows provide visual clue about the possible patterns in the scene. Last row images indicate rare patterns or possible outliers. The images labeled 5, 7, and 9 can be grouped together to form a single class indicating the downward traffic flow.

The annotation aspects of unlabeled trajectories using mDPMM are shown in Fig.4. Trajectory details are presented

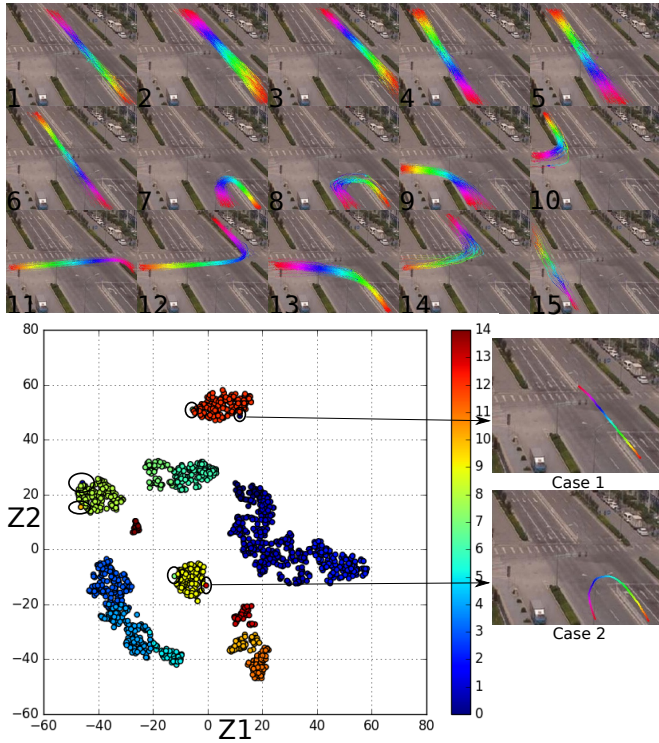


Fig. 5. Illustration of t-SNE on T15. z_1 and z_2 represent the dimensions after feature transformation. Top three rows represent trajectories of 15 classes. Most of the anomalies can be identified from the t-SNE visualization. For example, Case:1 the trajectory is anomalous due to the truncation and Case 2 is a normal trajectory. The latter is away from the respective class distribution due to the U-turn variation of the vehicle.

in Table I. Since T15 dataset readily comes with associated class annotation, unsupervised clustering has not been applied on this dataset. Fig.5 presents the t-SNE guided refinement.

B. Experiments on Classification and Comparisons

Trajectories of T15, QMUL and 4WAY datasets have been used for classification. Classification results are shown in Fig. 6 and summarized in Table II. It can be observed that the proposed method performs accurate classification across all datasets. We have randomly selected 75% of the trajectories for training and the rest for testing. Our proposed classification method has been compared with other state-of-the-art classification methods such as HAR-CNN [30], LSTM [13], LSTM+CNN [28] typically used for time series data classification. We have converted the input trajectories to 128 samples by downsampling or upsampling depending on their size. The comparative results are shown in Table II. The results reveal, our proposed method performs better than the existing approaches across all datasets. However, classification without color gradient degrades slightly, even though it performs better than most of the existing work.

C. Experiments on Anomaly Detection

T15 dataset has been used to evaluate the anomaly detection framework. Reconstructions using VAE are depicted

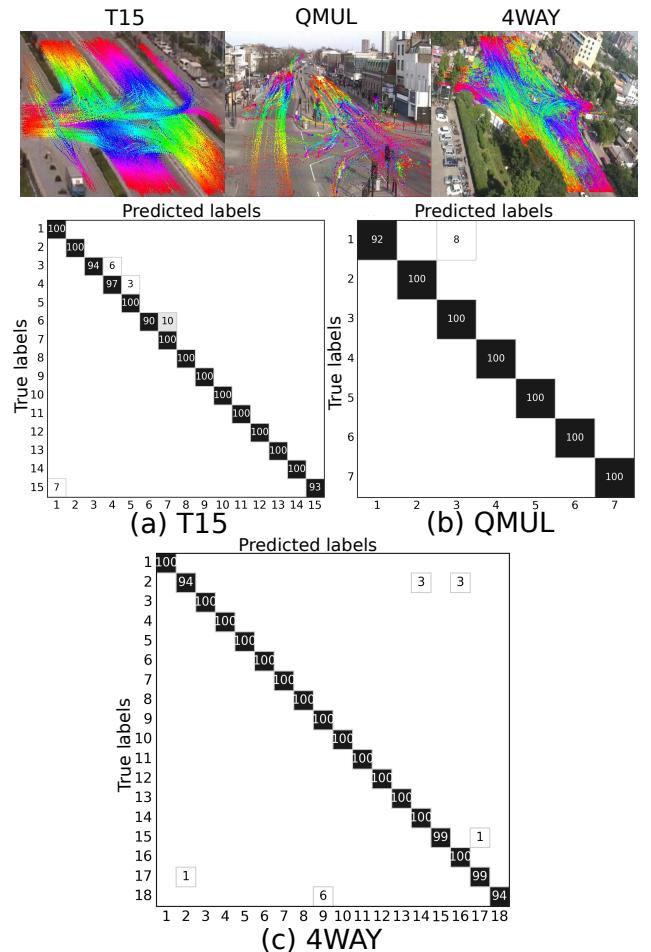


Fig. 6. Illustration of test results using normalized confusion matrices (in %) for T15, QMUL and 4WAY datasets.

TABLE II
COMPARISON OF CLASSIFICATION ACCURACIES ON THREE DATASETS

Method	T15	QMUL	4WAY
Proposed	99.0%	97.3%	99.5%
HAR-CNN [30]	94.9%	97.3%	98.7%
LSTM [13]	93.4%	88.6%	93.0%
LSTM+CNN [28]	93.4%	91.3%	94.1%
Proposed (no gradient)	98.0%	94.6%	99.1%

in Fig. 7. Four kinds of anomalous trajectories are used in our experiments: (i) Trajectories terminating abruptly. (ii) Speed variation as compared to normal trajectories of the same class. (iii) Trajectories of objects traveling in opposite direction of the normal traffic. (iv) Trajectories corresponding to vehicles violating lane driving. Since T15 dataset does not contain type three anomalous trajectories, we have created a few such trajectories by gradient conversion in reverse order. We have used two times the converged loss value as a threshold for detecting anomaly based on the empirical study on anomalous and normal trajectories as shown in Fig 8(a). Anomaly detection results are shown in Fig. 8(b). We have used 69 randomly selected normal trajectories that are not used in the training and 31 identified anomalous trajectories

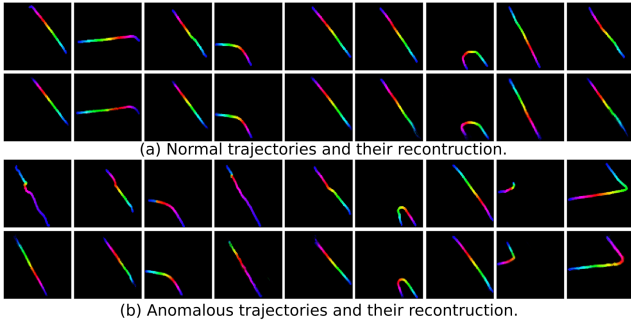


Fig. 7. Illustration of anomaly detection on T15 dataset. First row presents test trajectories and second row presents corresponding reconstructed patterns. (a) Correct reconstruction happens for the non-anomalous trajectories. (b) Reconstruction fails on anomalous trajectories. Columns 2 and 5 represent lane change anomaly. Columns 1 and 3 represent speed variations. Column 4 represents vehicle stopping then moving and column 8 represents terminated trajectory. Columns 6, 7 and 9 represent vehicle moving in opposite direction of normal traffic.

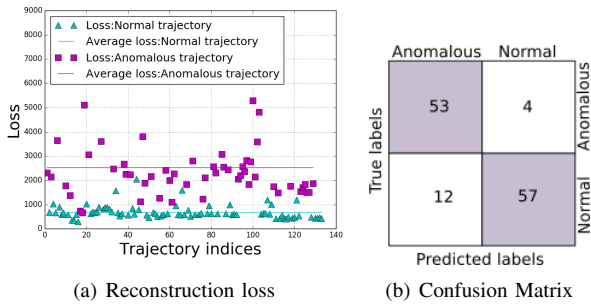


Fig. 8. Depiction of anomaly detection. (a) Reconstruction loss (l_i) for normal and anomalous trajectories. The loss depends on the amount of deviation from the normal path. (b) Confusion matrix for anomaly detection experiments on T15 dataset.

TABLE III
COMPARISONS OF ANOMALY DETECTION

Method	Accuracy	Precision	Recall
Without Gradient	49.2%	46.7%	85.9%
Without t-SNE	86.5%	83.1%	87.5%
With t-SNE	87.3%	81.5%	93.0%

along with synthetically created ones. We have created 11 synthetic trajectories for lane change and 15 corresponding to each class for opposite direction driving anomalies. The comparisons of trajectory projection on image plane using VAE under different conditions are presented in Table III. We are able to detect anomalies with an accuracy of 87.3% when t-SNE is used. This reveals, without gradient representation, anomaly detection accuracy drops significantly (49.2%).

D. Comparison of Anomaly Detections

Since video trajectory-based anomaly detection method using DNNs proposed in this paper is of the first kind, we could not find benchmark datasets that can be used in comparison with neural network-based anomaly detection. Hence we have performed high-level comparison with the state-of-the-art anomaly detection techniques presented in [19] and [5] using the input reconstruction property. The work proposed in [19] uses sparse combination learning for learning normal

behavior, while [5] learns the model from the spatio-temporal video segments using Autoencoder. Several experiments have been conducted on QMUL dataset. Training videos have been created by splitting the original traffic video into 42 segments starting from the frame number 8610 by eliminating anomalous segments from the scene. Testing has been conducted using the video segment prior to the frame number 8610. We have trained our proposed architecture using trajectories obtained with the help of the method proposed in [32] with $\delta = 836$ for the testing. Training for the method proposed in [19] has been done using the same configuration as reported in their work, while testing has been conducted with an error threshold of 0.4. For training the model proposed in [5], we have used a sequence length (T) = 10, batch size = 4 and number of epochs = 200.

The test results are depicted through Figs. 9-12. It can be observed that both the methods proposed in [19] and [5] report several false positives on the QMUL dataset. Moreover, these methods cannot detect contextual anomalies. A deeper analysis reveals that the false positives are mainly due to the unseen characteristics present in the scene with heterogeneous data, making it difficult to learn all spatio-temporal features. Such methods can work only when the video duration is long enough that can learn all types of object motions possible within a scene. However, it may be difficult to train as separating normal video segments from the anomalous can be very challenging when anomalies are present throughout the video. As our method is trajectory-based, individual trajectories can be characterized as normal or abnormal rather than declaring a video segment normal or anomalous. Moreover, training a deep neural network using video frames can be time consuming. On the contrary, a trajectory has been condensed into a single video frame as done in our method. In a nutshell, we are using the advantages of conventional trajectory extraction methods as well as the feature extraction capabilities of deep neural network to achieve classification which is fast. Table IV summarizes the comparative results.

TABLE IV
COMPARISONS OF ANOMALY DETECTION WITH STATE-OF-ART

Parameters	Proposed method	Sparse reconstruction [19]	Spatio-temporal autoencoder [5]
False alarm rate	Low	High	High
Unknown anomaly detection	Yes	Yes	Yes
Contextual anomaly detection	Yes	No	No
Training difficulty	Low	High	High
Anomaly localization	Yes	Yes	No
Detection time	Once trajectory is available	Per frame	Per sequence length

E. Discussions and Limitations

Key to accurate anomaly detection lies in training the model with normal trajectories. Apart from mDPMM-based

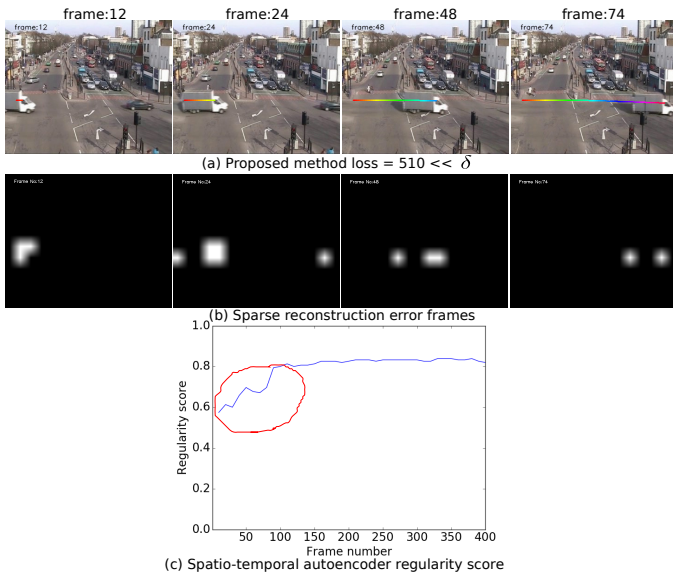


Fig. 9. Illustration of false alarms in the sparse reconstruction technique [19] and spatio-temporal Autoencoder-based method [5]. Though the traffic is open for vehicles at the junction for east and south bound traffic from top-peft lane, anomalies are reported for the scene for these methods. (a) A trajectory corresponding to a truck in color gradient form in different frames. (b) The respective sparse reconstruction error frames using [19]. White patches represent the anomalies detected by the method proposed in [19]. It can be seen that several false positives are present throughout the sequence. (c) The regularity score for the scene using [5] shown for the frame sequence. The highlighted portion indicates some anomaly.

clustering, t-SNE visualization plays an important role in eliminating anomalous trajectories. The need for a classifier is to detect known anomalies such as traffic rule violations by vehicles. While unknown anomalies are detected using VAE, the CNN classifier helps to identify known anomalies and to localize the path of unknown anomalies. The loss values in terms of KLD and likelihood are justified as they represent the distance of the trajectories from the allowed class distributions. A small offset from the converged loss can be a good estimate of the threshold. CNN classifier performs with higher accuracy as compared to other methods.

Some of the limitations of the proposed method are: (i) The method is tracking dependent. However, with improved tracking, we can overcome this issue. (ii) A large number of training samples need to be available to learn the allowed paths in a traffic junction.

IV. CONCLUSIONS

The key idea behind this work is to represent time varying visual data using color gradient form in order to train DNN-based systems for encoding temporal features. This method combines traditional object tracking-based results to be combined with neural network-based methods to use the advantages of both systems. It has been observed through experiments that the proposed color gradient feature using CNN performs better than existing classifiers. We are also able to detect a few types of trajectory anomalies using the proposed architecture. It performs better than some of the existing reconstruction-based anomaly detection methods. We plan to extend this work to develop a real-time anomaly detection

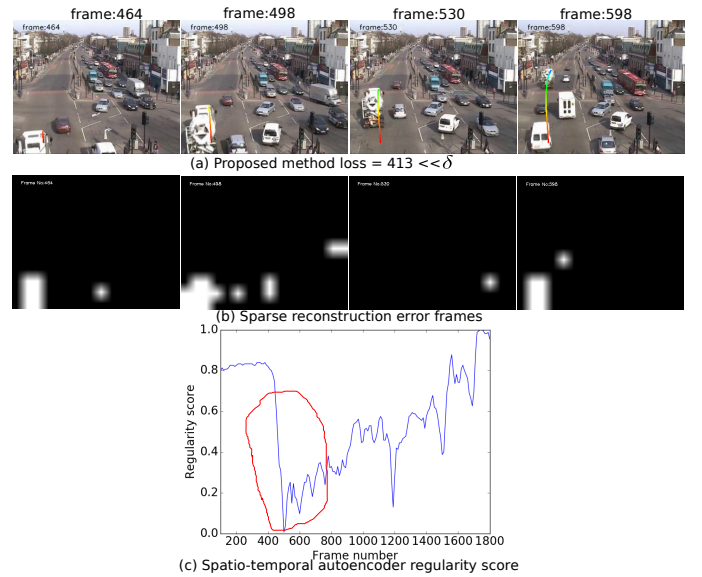


Fig. 10. Depiction of anomaly detection for a scene, where the signal opens for north-bound and south-bound traffics. Though the traffic movement seems to be normal, false positives are reported using [19] and [5]. (a) A trajectory corresponding to a truck in color gradient form in different frames. (b) Corresponding sparse reconstruction error frames using [19]. False positives are present even for other vehicles. (c) The regularity score for the scene using [5] during the frame sequence. False positives can be seen for longer duration when heavy traffic flow is underway during green signal.

system for traffic intersections using online trajectories which will be able to detect discussed anomalies as well as other anomalies such as over-speeding. We also plan to explore this method for time series data analysis in other domains.

ACKNOWLEDGMENT

We gratefully acknowledge the support of NVIDIA Corporation with the donation of the Quadro P5000 GPU used for this research.

REFERENCES

- [1] S. H. Bae and K. J. Yoon. Confidence-based data association and discriminative deep appearance learning for robust online multi-object tracking. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PP(99):1–1, 2017.
- [2] Ben Benfold and Ian Reid. Stable multi-target tracking in real-time surveillance video. In *CVPR*, 2011.
- [3] A. F. Bobick and J. W. Davis. The recognition of human movement using temporal templates. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 23(3):257–267, March 2001.
- [4] A. Bulling, U. Blanke, and B. Schiele. A tutorial on human activity recognition using body-worn inertial sensors. *ACM Computing Surveys (CSUR)*, 46(3):33, 2014.
- [5] Y. S. Chong and Y. H. Tay. Abnormal event detection in videos using spatiotemporal autoencoder. In *ISNN*, 2017.
- [6] L. Deng, J. Li, J. T. Huang, K. Yao, D. Yu, F. Seide, M. Seltzer, G. Zweig, X. He, J. Williams, et al. Recent advances in deep learning for speech research at microsoft. In *ICASSP*, 2013.
- [7] T. B. Dinh, N. Vo, and G. Medioni. Context tracker: Exploring supporters and distracters in unconstrained environments. In *CVPR*, 2011.
- [8] J. Donahue, L. A. Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell. Long-term recurrent convolutional networks for visual recognition and description. In *CVPR*, 2015.
- [9] Y. Endo, H. Toda, K. Nishida, and J. Ikedo. Classifying spatial trajectories using representation learning. *International Journal of Data Science and Analytics*, 2(3):107–117, Dec 2016.

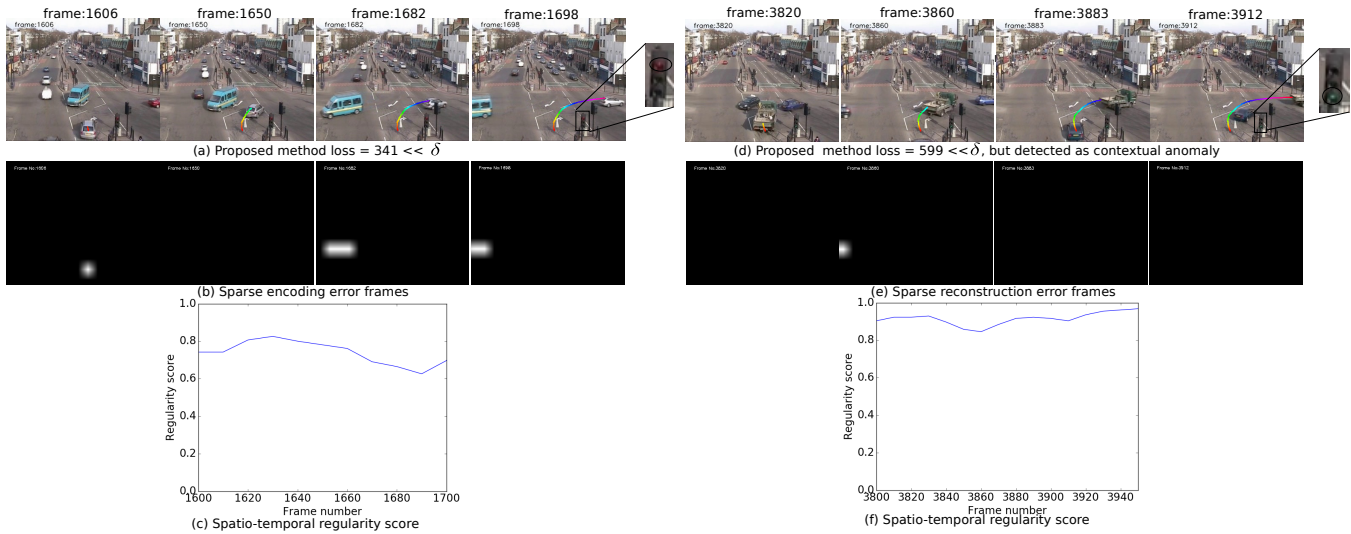


Fig. 11. Depiction of contextual anomaly detection. When the traffic signal is green, only north-bound and south-bound traffics are allowed. Going right or left is allowed only after the end of the signal. Traffic scenario depicted through (a-c) is a normal condition, while (d-f) depict contextual anomaly. (a) A trajectory corresponding to a car in color gradient during normal traffic flow. When the signal is red, the waiting vehicles are allowed to go to the right or left, i.e. the pattern is an allowed one. The loss value is well below the anomaly threshold, indicating that it is a normal flow. (b) The respective sparse reconstruction error frames using [19]. Though false positives are not observed throughout for the tracked vehicle during this period, it is present in some frames. False positives can also be observed for the vehicle turning left. (c) The regularity score for the scene using [5]. The regularity score does not indicate any anomaly. (d) A trajectory corresponding to a truck in color gradient form during an anomalous traffic flow. As the signal is green, even though the no vehicles can be seen heading south ward, vehicles are not supposed to cross towards east side. Though the loss is less than the threshold, this is categorized as an unknown anomaly using our method. (e) Corresponding sparse reconstruction error using [19] with no anomalies detected. (f) The regularity score for the scene using [5]. The regularity score does not indicate any anomaly, though there is a contextual anomaly.

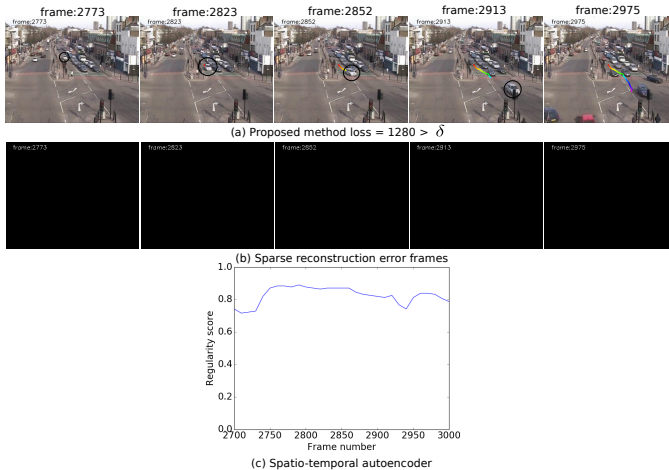


Fig. 12. Illustration of lane change anomaly with truncated trajectory using three different methods. (a) The highlighted vehicle gets tracked very late and the tracking fails and wrong trajectories are created. However, our method can detect it as an anomaly. (b) Corresponding sparse reconstruction error using [19] with no possible anomaly. (c) The regularity score of the scene using [5]. The regularity score does not indicate any anomaly.

[10] N. Y. Hammerla, S. Halloran, and T. Ploetz. Deep, convolutional, and recurrent models for human activity recognition using wearables. *arXiv preprint arXiv:1604.08880*, 2016.

[11] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

[12] S. Ji, W. Xu, M. Yang, and K. Yu. 3d convolutional neural networks for human action recognition. *IEEE transactions on pattern analysis and machine intelligence*, 35(1):221–231, 2013.

[13] A. Khosroshahi, E. Ohn-Bar, and M. M. Trivedi. Surround vehicles trajectory analysis with recurrent neural networks. In *ITSC*, 2016.

[14] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.

[15] D. P. Kingma and M. Welling. Auto-encoding variational bayes. In *ICLR*, 2014.

[16] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.

[17] J. Lin, E. Keogh, S. Lonardi, and B. Chiu. A symbolic representation of time series, with implications for streaming algorithms. In *DMKD*, 2003.

[18] C. C. Loy, T. Xiang, and S. Gong. From local temporal correlation to global anomaly detection. In *ECCV*, 2008.

[19] C. Lu, J. Shi, and J. Jia. Abnormal event detection at 150 fps in matlab. In *ICCV*, 2013.

[20] L. V. D. Maaten and G. Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605, 2008.

[21] A. Milan, L. Leal-Taix, K. Schindler, and I. Reid. Joint tracking and segmentation of multiple targets. In *CVPR*, 2015.

[22] J. Y.-H. Ng, M. Hausknecht, S. Vijayanarasimhan, O. Vinyals, R. Monga, and G. Toderici. Beyond short snippets: Deep networks for video classification. In *CVPR*, 2015.

[23] T. Plötz, N. Y. Hammerla, and P. Olivier. Feature learning for activity recognition in ubiquitous computing. In *IJCAI*, 2011.

[24] K. K. Santhosh, D. P. Dogra, and P. P. Roy. Temporal unknown incremental clustering model for analysis of traffic surveillance videos. *IEEE Transactions on Intelligent Transportation Systems*, pages 1–12, 2018.

[25] Y. Shi, Y. Tian, Y. Wang, and T. Huang. Sequential deep trajectory descriptor for action recognition with three-stream cnn. *IEEE Transactions on Multimedia*, 19(7):1510–1520, 2017.

[26] Tijmen Tieleman. Training restricted boltzmann machines using approximations to the likelihood gradient. In *ICML*, 2008.

[27] H. Wang, A. Klser, C. Schmid, and C. L. Liu. Action recognition by dense trajectories. In *CVPR*, 2011.

[28] Z. Wu, X. Wang, Y. G. Jiang, H. Ye, and X. Xue. Modeling spatial-temporal clues in a hybrid deep learning framework for video classification. In *ACM MULTIMEDIA*, 2015.

[29] H. Xu, Y. Zhou, W. Lin, and H. Zha. Unsupervised trajectory clustering via adaptive multi-kernel-based shrinkage. In *ICCV*, 2015.

[30] J. Yang, M. N. Nguyen, P. P. San, X. Li, and S. Krishnaswamy. Deep convolutional neural networks on multichannel time series for human activity recognition. In *IJCAI*, 2015.

[31] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *ECCV*, 2014.

- [32] B. Zhou, X. Tang, and X. Wang. Measuring crowd collectiveness. In *CVPR*, 2013.



Kelathodi Kumaran Santhosh is a research scholar in the School of Electrical Sciences, IIT Bhubaneswar, India. He joined a Ph.D. program for resuming his research work that can help humanity. His interests are in the development of vision based applications that can replace human factor. He is a member of IEEE. Prior to joining IIT Bhubaneswar, he worked for Huawei Technologies India Pvt. Ltd. for 10 years (2005-2015) and in Defence Research Development Organization (DRDO) as a Scientist for around 2 years (2003-2004). During his tenure

with Huawei, he has worked in many signalling protocols such as Diameter, Radius, SIP etc. in the role of a developer, technical leader, project manager and also served the product lines HSS, CSCF etc. in Huawei China as a support engineer for closer to 1.5 years. In DRDO, he worked in the field of object tracking algorithms based on the data received from radars. More information on Santhosh can be found at <https://sites.google.com/site/santhoshkelathodi>.



Dr. Debi Prosad Dogra is an Assistant Professor in the School of Electrical Sciences, IIT Bhubaneswar, India. He received his M.Tech degree from IIT Kanpur in 2003 after completing his B.Tech. (2001) from HIT Haldia, India. After finishing his masters, he joined Haldia Institute of Technology as a faculty members in the Department of Computer Sc. & Engineering (2003-2006). He has worked with ETRI, South Korea during 2006-2007 as a researcher. Dr. Dogra has published more than 45 international journal and conference papers in the

areas of computer vision, image segmentation, and healthcare analysis. He is a member of IEEE. More information on Dr. Dogra can be found at <http://www.iitbbs.ac.in/profile.php/dpdogra>.



Dr. Partha Pratim Roy has obtained his M.S. and Ph. D. degrees in the year of 2006 and 2010, respectively at Autonomous University of Barcelona, Spainis. Presently he is an Assistant Professor in the Department of Computer Science and Engineering, IIT Roorkee, India in 2014. Prior to joining, IIT Roorkee, Dr. Roy was with Advanced Technology Group, Samsung Research Institute Noida, India during 2013-2014. Dr. Roy was with Synchromedia Lab, Canada in 2013 and RFAI Lab, France in 2012 as postdoctoral research fellow. His research

interests are Pattern Recognition, Multilingual Text Recognition, Biometrics, Computer Vision, Image Segmentation, Machine Learning, and Sequence Classification. He has published more than 65 papers in international journals and conferences.



Dr. Adway Mitra is a researcher, interested in Machine Learning and Data Mining, and especially in the application of these techniques to solve problems affecting the world. More specifically, he is interested in data-driven modeling and simulation of complex spatio-temporal processes. His background is in Computer Science and Engineering, and PhD thesis was related to semantic Video Analytics, using Bayesian modeling techniques. Many of the techniques and concepts he developed during PhD may be extended to spatio-temporal processes in other

domains. By doing so, he intend to build a career in interdisciplinary research. He is currently focusing on Climate Informatics - application of Computer Science (especially Data Science) concepts to solve problems in Climate Science. He is particularly interested in the following questions in this domain: 1) Realistic simulation of climatic processes, through stochastic processes 2) Understanding and Modeling the dynamics of Indian Monsoon and its various vagaries like onset, withdrawal and active/break spells 3) Extreme events - their various statistical properties, and links between different extreme events 4) Identification of widespread and long-lasting anomalies such as droughts and heat waves in huge volumes of climatic data 5) Causal relationships between different events, aimed at attribution.