

# EFFICIENT ANOMALY DETECTION IN THE REAL TIME SURVEILLANCE VIDEOS

<sup>1</sup>Dr., C. Gulzar Mtech, Ph.D <sup>2</sup>A.V Hari Chandra Reddy

<sup>1</sup>Associate Professor <sup>2</sup>P.G.Scholar

<sup>1,2</sup>DEPARTMENT OF CSE

<sup>1,2</sup>Dr. K.V SUBBA REDDY INSTITUTE OF TECHNOLOGY

Dupadu Railway Station, Lakshmipuram Post, Kurnool, Andhra Pradesh, India - 518218

## ABSTRACT

Surveillance videos are able to capture a variety of realistic anomalies. In this paper, we propose to learn anomalies by exploiting both normal and anomalous videos. To avoid annotating the anomalous segments or clips in training videos, which is very time consuming, we propose to learn anomaly through the deep multiple instance ranking framework by leveraging weakly labeled training videos,

i.e. the training labels (anomalous or normal) are at video-level instead of clip-level. In our approach, we consider normal and anomalous videos as bags and video segments as instances in multiple instance learning (MIL), and automatically learn a deep anomaly ranking model that predicts high anomaly scores for anomalous video segments. Furthermore, we introduce sparsity and temporal smoothness constraints in the ranking loss function to better localize anomaly during training.

We also introduce a new large-scale first of its kind dataset of 128 hours of videos. It consists of 1900 long and untrimmed real-world surveillance videos, with 13 realistic anomalies such as fighting, road accident, burglary, robbery, etc. as well as normal activities. This dataset can be used for two tasks. First, general anomaly detection considering all anomalies in one group and all normal activities in another group. Second, for recognizing each of 13 anomalous activities. Our experimen-

tal results show that our MIL method for anomaly detection achieves significant improvement on anomaly detection performance as compared to the state-of-the-art approaches. We provide the results of several recent deep learning baselines on anomalous activity recognition. The low recognition performance of these baselines reveals that our dataset is very challenging and opens more opportunities for future work. The dataset is available at:

<https://webpages.uncc.edu/cchen62/dataset.html>

**KEYWORDS:** Anomaly , cnn, opencv, Machine Learning

## I. INTRODUCTION

Surveillance cameras are increasingly being used in public places e.g. streets, intersections, banks, shopping malls, etc. to increase public safety. However, the monitoring capability of law enforcement agencies has not kept pace. The result is that the reimagining deficiency in the utilization of surveillance cameras and an unworkable ratio of cameras to human monitors. One critical task in video surveillance is detecting anomalous events such as traffic accidents, crimes or illegal activities. Generally, anomalous events rarely occur as compared to normal activities. Therefore, to alleviate the waste of labor and time, developing intelligent computer vision algorithms for automatic video anomaly detection is a pressing need. The goal of a practical anomaly detection system is to timely signal an activity that deviates normal patterns

and identify the time window of the occurring anomaly. Therefore, anomaly detection can be considered as coarse level video understanding, which filters out anomalies from normal patterns. Once an anomaly is detected, it can further be categorized into one of the specific activities using classification techniques.

A small step towards addressing anomaly detection is to develop algorithms to detect as specific anomalous event, for example violence detector[30]and traffic accident detector[23, 34]. However, it is obvious that such solutions cannot be generalized to detect other anomalous events, therefore they render a limited use in practice.

Real-world anomalous events are complicated and diverse. It is difficult to list all of the possible anomalous events. Therefore, it is desirable that the anomaly detection algorithm does not rely on any prior information about the events. In other words, anomaly detection should be done with minimum supervision. Sparse-coding based approaches [28,41] are considered as representative methods that achieve state-of-the-art anomaly detection results. These methods assume that only a small initial portion of a video contains normal events, and therefore the initial portion is used to build the normal event dictionary. Then, the main idea for anomaly detection is that anomalous events are not accurately reconstructable from the normal event

### **Motivation and contributions.**

Although the above-mentioned approaches are appealing, they are based on the assumption that any pattern that deviates from the learned normal patterns would be considered as an anomaly. However, this assumption may not hold true because it is very difficult or impossible to define a normal event which takes all possible normal patterns/behaviors into account[9]. More importantly, the boundary between normal and anomalous behaviors is often ambiguous. In addition, un-

der realistic conditions, the same behavior could be a normal or an anomalous behavior under different conditions. Therefore, it is argued that the training data of normal and anomalous events can help an anomaly detection system learn better. In this paper, we propose an anomaly detection algorithm using weakly labeled training videos. That is we only know the video-level labels, i.e. a video is normal or contains anomaly somewhere, but we do not know where. This is intriguing because we can easily annotate a large number of videos by only assigning video-level labels. To formulate a weakly supervised learning approach, we sort to multiple instance learning(MIL)[12,4]. Specifically, we propose to learn anomaly through a deep MIL framework by treating normal and anomalous surveillance videos as bags and short segments/clips of each video as instances in a bag. Based on training videos, we automatically learn an anomaly ranking model that predicts high anomaly scores for anomalous segments in a video. During testing, a long-untrimmed video is divided into segments and fed into our deep network which assigns anomaly score for each video segment such that an anomaly can be detected. In summary, this paper makes the following contributions.

We propose a MIL solution to anomaly detection by leveraging only weakly labeled training videos. We propose a MIL ranking loss with sparsity and smoothness constraints for a deep learning network to learn anomaly scores for video segments. To the best of our knowledge, we are the first to formulate the video anomaly detection problem in the context of MIL.

We introduce a large-scale video anomaly detection dataset consisting of 1900 real-world surveillance videos of 13 different anomalous events and normal activities captured by surveillance cameras. It is by far the largest dataset with more than 15 times videos than

existing anomaly datasets and has a total of 128 hours of videos.

## II. RELATED WORK

Anomaly detection. Anomaly detection is one of the most challenging and long standing problems in computer vision [39, 38, 7, 10, 5, 20, 43, 27, 26, 28, 42, 18, 26].

For

video surveillance applications, there are several attempts to detect violence or aggression [15, 25, 11, 30] in videos. Datta et al. proposed to detect human violence by exploiting motion and limbs orientation of people. Kooij et al. [25] employed video and audio data to detect aggressive actions in surveillance videos. Gao et al. proposed violent flow descriptor to detect violence in crowd videos. More recently, Mohammadi et al. [30] proposed a new behavior heuristic based approach to classify violent and non-violent videos.

Beyond violent and non-violent patterns discrimination, authors in [38, 7] proposed to use tracking to model the normal motion of people and detect deviation from that normal motion as an anomaly. Due to difficulties in obtaining reliable tracks, several approaches avoid tracking and learn global motion patterns through histogram-based methods [10], topic modeling [20], motion patterns [31], social force models [29], mixtures of dynamic textures model [27], Hidden Markov Model (HMM) on local spatio-temporal volumes [26], and context-driven method [43]. Given the training videos of normal behaviors, these approaches learn distributions of normal motion patterns and detect low probability patterns as anomalies.

Following the success of sparse representation and dictionary learning approaches in several computer vision problems, researchers in [28, 42] used sparse representation to learn the dictionary of normal behaviors. During testing, the patterns which have large reconstruction errors are considered as

anomalous behaviors. Due to successful demonstration of deep learning for image classification, several approaches have been proposed for video action classification [24, 36]. However, obtaining annotations for training is difficult and laborious, specifically for videos.

Recently, [18, 39] used deep learning based autoencoders to learn the model of normal behaviors and employed reconstruction loss to detect anomalies. Our approach not only considers normal behaviors but also anomalous behaviors for anomaly detection, using only weakly labeled training data.

Ranking. Learning to rank is an active research area in machine learning. These approaches mainly focused on improving relative scores of the items instead of individual scores. Joachims et al. [22] presented rank-SVM to improve retrieval quality of search engines. Bergeron et al. [8] proposed an algorithm for solving multiple instance ranking problems using successive linear programming and demonstrated its application in hydrogen abstraction problem in computational chemistry. Recently, deep ranking networks have been used in several computer vision applications and have shown state-of-the-art performances. They have been used for feature learning [37], highlight detection [40], Graphics Interchange Format (GIF) generation [17], face detection and verification [32], person re-identification [13], place recognition [6], metric learning and image retrieval [16]. All deep ranking methods require a vast amount of annotations of positive and negative samples.

In contrast to the existing methods, we formulate anomaly detection as a regression problem in the ranking framework by utilizing normal and anomalous data. To alleviate the difficulty of obtaining precise segment-level labels, we rely on weakly labeled data (i.e. video-level labels – normal or abnormal which are much easier to obtain than temporal annotations) to learn the anomaly model and

detect video segment level anomaly during testing.

## 1. DATASET

Due to the limitations of previous datasets, we construct an ewlarge-scale data set to evaluate our method. It consists of long untrimmed surveillance videos which cover 13 real-world anomalies, including Abuse, Arrest, Arson, Assault, Accident, Burglary, Explosion, Fighting, Robbery, Shoot-ing, Stealing, Shoplifting, and Vandalism. These anomalies are selected because they have a significant impact on public safety. We compare our dataset with previous anomaly detection datasets in Table 1.

**Video collection.** To ensure the quality of our dataset, we train ten annotators (having different levels of computer vision expertise) to collect the dataset. We search videos on YouTube [1] and Live Leak [2] using text search queries (with slight variations e.g. "car crash", "road accident") of each anomaly. In order to retrieve as many videos as possible, we also use text queries in different languages (e.g. French, Russian, Chinese, etc.) for each anomaly, thanks to Google translator. We remove videos which fall into any of the following conditions: manually edited, prank videos, not captured by CCTV cameras, taking from news, captured using a hand-held camera, and containing compilation. We also discard videos in which the anomaly is not clear. With the above video pruning constraints, 950 unedited real-world surveillance videos with clear anomalies are collected. Using the same constraints, 950 normal videos are gathered, leading to a total of 1900 videos in our dataset. In Figure 2, we show four frames of an example video from each anomaly.

**An notation.** For our anomaly detection method, only video-level labels are required for training. However, in order to evaluate its performance on testing videos, we need to know the temporal annotations, i.e. the start and ending frames of the anomalous event in

each testing anomalous video. To this end, we assign the same videos to multiple annotators to label the temporal extent of each anomaly. The final temporal annotations are obtained by averaging the annotations of different annotators. The complete data set is finalized after intense efforts of several months. **Training and testing sets.** We divide our dataset into two parts: the training set consisting of 800 normal and 810 anomalous videos (details shown in Table 2) and the testing set including the remaining 150 normal and 140 anomalous videos. Both training and testing sets contain all 13 anomalies at various temporal locations in the videos. Furthermore, some of the videos have multiple anomalies. The distribution of the training videos in terms of length (in minute) is shown in Figures 3. The number of frames and percentage of anomaly in each testing video are presented in Figures 4 and 5, respectively.

## 2. EXPERIMENTS

### 2.1. IMPLEMENTATION DETAILS

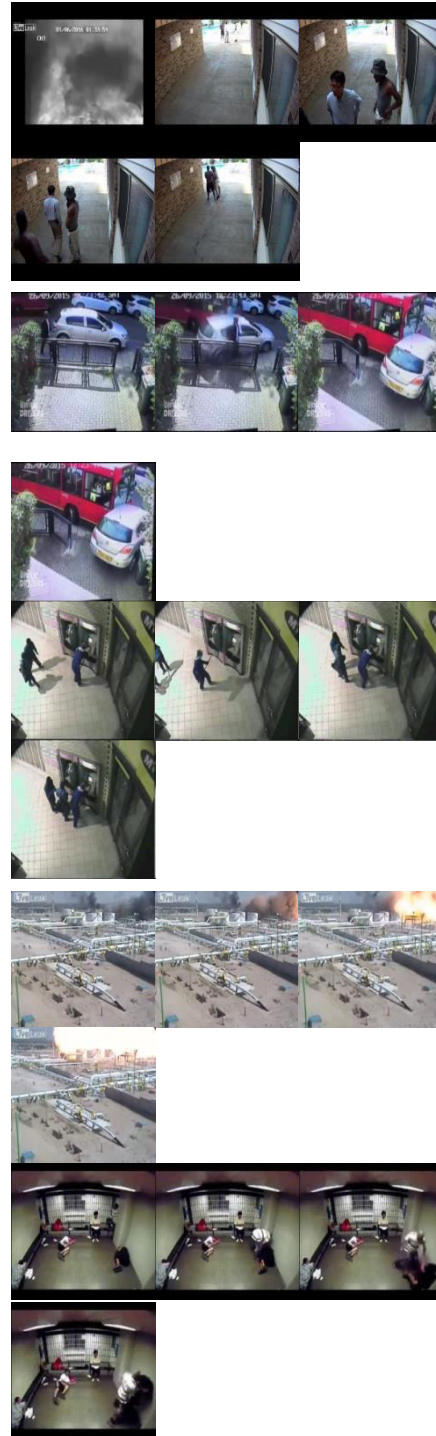
We extract visual features from the fully connected (FC) layer FC6 of the C3D network [36]. Before computing features, we re-size each video frame to 240x320 pixels and fix the frame rate to 30 fps. We compute C3D features for every 16-frame video clip followed by 12 normalization. To obtain features for a video segment, we take the average of all 16-frame clip features within that segment. We input these features (4096D) to a 3-layer FC neural network. The first FC layer has 512 units followed by 32 units and 1 unit FC layers. 60% dropout regularization [33] is used between FC layers. We experiment with deeper networks but do not observe better detection accuracy. We use ReLU [19] activation and Sigmoid activation for the first and the last FC layers respectively, and employ Adagrad [14] optimizer with the initial learning rate of 0.001. The parameters of sparsity and smoothness constraints in the MIL ranking



loss are set to  $\lambda_1=\lambda_2=8 \times 10^{-5}$  for the best performance. We divide each video into 32 non-overlapping segments and consider each video segment as an instance of the bag. The number of segments (32) is empirically set. We also experimented with multi-scale overlapping temporal segments but it does not affect detection accuracy. We randomly select 30 positive and 30 negative bags as a mini-batch. We compute gradients by reverse mode automatic differentiation on computation graph using Theano[35]. Specifically, we identify set of variables on which loss depends, compute gradient for each variable and obtain final gradient through chain rule on the computation graph. Each video passes through the network and we get the score for each of its temporal segments. Then we compute loss as shown in Eq. 6 and Eq. 7 and back-propagate the loss for the whole batch. Evaluation Metric. Following previous works on anomaly detection[27], we use frame based receiver operating characteristic (ROC) curve and corresponding area under the curve (AUC) to evaluate the performance of our method. We do not use equal error rate (EER)[27] as it does not measure anomaly correctly, specifically if only a small portion of along video contains anomalous behavior

	#of videos	Average frame s	Dataset length	Example anomalies
UCSDPed1[27]	70	201	5min	Bikers, small carts, walking across walkways
UCSDPed2[27]	28	163	5min	Bikers, small carts, walking across walkways
SubwayEntrance[3]	1	121,749	1.5hours	Wrong direction, No payment
SubwayExit[3]	1	64,901	1.5hours	Wrong direction, No payment
Avenue[28]	37	839	30min	Ram, throw, no object
UMIN[2]	5	1290	5min	Ram
BOSS[1]	12	4092	29min	Harm, Disease, Panic
Ours	1900	7247	128hours	Abuse, arrest, arson, assault, accident, burglary, fighting, robbery

Table 1. A comparison of anomaly datasets. Our dataset contains large number of longer surveillance videos with more realistic anomalies.



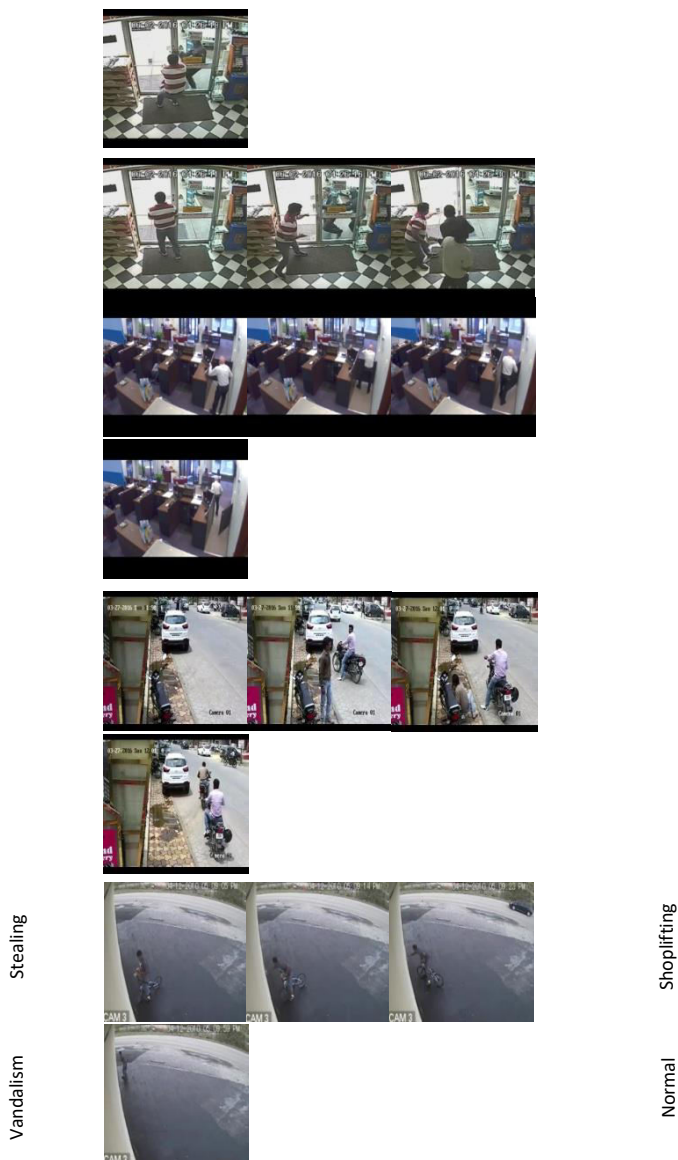


Figure 2. Examples of different anomalies from the training and testing videos in our dataset.

### III. CONCLUSIONS

We propose a deep learning approach to detect real-world anomalies in surveillance videos. Due to the complexity of these realistic anomalies, using only normal data alone may not be optimal for anomaly detection. We attempt to exploit both normal and anomalous surveillance videos. To avoid labor-intensive temporal annotations of anomalous segments in training videos, we learn a general model of anomaly detection using deep multiple instance ranking framework with weakly labeled data. To validate the proposed approach, a new large-scale anomaly dataset consisting of a variety of real-world anomalies is introduced. The experimental results on this dataset

show that our proposed anomaly detection approach performs significantly better than baseline methods. Furthermore, we demonstrate the usefulness of our dataset for the second task of anomalous activity recognition.

### IV. References

- [1] <http://www.multitel.be/image/research-development/research-projects/boss.php>.
- [2] Unusual crowd activity dataset of university of minnesota. In <http://mha.cs.umn.edu/movies/crowdactivity-all.avi>.
- [3] A. Adam, E. Rivlin, I. Shimshoni, and D. Reinitz. Robust real-time unusual event detection using multiple fixed-location monitors. *TPAMI*, 2008.
- [4] S. Andrews, I. Tsochanaridis, and T. Hofmann. Support vector machines for multiple-instance learning. In *NIPS*, pages 577–584, Cambridge, MA, USA, 2002. MIT Press.
- [5] B. Antani and B. Ommer. Video parsing for abnormality detection. In *CCV*, 2011.
- [6] R. Arandjelović, P. Gronat, A. Torii, T. Pajdla, and J. Sivic. NetVLAD: CNN architecture for weakly supervised place recognition. In *CVPR*, 2016.
- [7] A. Basharat, A. Gritai, and M. Shah. Learning object motion patterns for anomaly detection and improved object detection. In *CVPR*, 2008.
- [8] C. Bergeron, J. Zaretski, C. Breneman, and K. P. Bennett. Multiple instance ranking. In *ICML*, 2008.
- [9] V. Chandola, A. Banerjee, and V. Kumar. Anomaly detection: A survey. *ACM Comput. Surv.*, 2009.
- [10] X. Cui, Q. Liu, M. Gao, and D. N. Metaxas. Abnormal detection using interaction energy potentials. In *CVPR*, 2011.
- [11] A. Datta, M. Shah, and N. Da Vitoria Lobo. Person-on-person violence detection in video data. In *ICPR*, 2002.
- [12] T. G. Dietterich, R. H. Lathrop, and T. Lozano-Pérez. Solving the multiple instance problem with axis-parallel rectangles. *Artificial Intelligence*, 89(1):31–71, 1997.
- [13] S. Ding, L. Lin, G. Wang, and H. Chao. Deep feature learning with relative distance comparison for person re-identification. *Pattern Recognition*, 48(10):2993–3003, 2015.
- [14] J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *J. Mach. Learn. Res.*, 2011.
- [15] Y. Gao, H. Liu, X. Sun, C. Wang, and Y. Liu. Violence detection using oriented violent flows. *Image and Vision Computing*, 2016.
- [16] A. Gordo, J. Almazán, J. Revaud, and D. Larlus. Deep image retrieval: Learning global representations

- for image search. In *ECCV*, 2016.
- [17] M. Gygli, Y. Song, and L. Cao. Video2gif: Automatic generation of animated gifs from video. In *CVPR*, June 2016.
- [18] M. Hasan, J. Choi, J. Neumann, A. K. Roy-Chowdhury, and L. S. Davis. Learning temporal regularity in video sequences. In *CVPR*, June 2016.