

Enhanced Intelligent Video Monitoring using Hybrid Integration of Spatiotemporal Autoencoders and Convolutional LSTMs

Ankita Umale-Nagmote¹, Charu Goel² and Nidhi Lal³

¹Dept. of Comp. Sci. and Engg., Indian Institute of Information Technology (IIIT) Nagpur, India

²Dept. of Basic Science Indian Institute of Information Technology (IIIT) Nagpur, India

³Dept. of Comp. Sci. and Engg. Visvesvaraya National Institute of Technology (VNIT) Nagpur, India

E-mail: dtea20cse011@iiitn.ac.in, charugoel@iiitn.ac.in, nidhilal@cse.vnit.ac.in

Keywords: Anomaly detection, convolutional autoencoders, semi-supervised learning, video surveillance, LSTM

Received: November 4, 2024

Perceiving meaningful activities in surveillance videos presents significant challenges due to the ambiguous nature of anomalies and scene complexity. This paper proposes a hybrid deep learning framework that combines spatial-temporal autoencoders with convolutional LSTMs for automated anomaly detection in surveillance videos. The architecture integrates stacked convolutional autoencoders for semi-supervised feature representation with LSTM networks for preserving temporal information. Experiments conducted on the UCSD Ped1 dataset demonstrate that our LSTM-based Stacked CAE achieves an AUC of 83.5%, a detection rate of 81.5%, and an Equal Error Rate (EER) of 19.2%. The model particularly excels in temporal pattern recognition with an accuracy of 84.5% and sequence processing efficiency of 82.7%. Comparative analysis with state-of-the-art methods reveals that the proposed architecture achieves competitive performance, particularly in handling complex motion patterns and maintaining temporal consistency. The model shows significant improvement in false alarm rate reduction at 15.8% compared to the basic CAE's 17.2%. The results demonstrate that integrating LSTM with stacked convolutional autoencoders provides a robust framework for real-world surveillance applications, especially in scenarios requiring both spatial and temporal anomaly detection

Povzetek: Predlagani hibridni model za video nadzor združuje spatio-temporalne autoencoderje in konvolucijske LSTM-e za zaznavanje anomalij ter dosega visoko točnost, nizko lažno alarmiranje in realnočasovno delovanje.

1 Introduction

Recently it has become clear that thousands of surveillance cameras are constantly in operation, some of them are installed in peripheral areas or streets where it is unlikely that anything dangerous will happen, while others are located in busy streets or city squares. There are many abnormal events that can happen even in one place, and the definition of an abnormal event varies from place to place and from time to time. Intelligent video surveillance systems can be widely used in various public places such as smart cities, markets, banks, malls, streets, etc. to improve public safety automatically detects unusual events such as crimes, traffic accidents, accidents, etc. In general, there are abnormal events less often compared to regular events and also very much depending on the application or context. Anomaly detection is a technique that identifies abnormal patterns or trends in data. In general, an anomaly can be defined as an observation that deviates significantly from other observations in the same context to awaken intuition which is produced by a different mechanism [15]. Alternatively, an anomaly that deviates from its trained pattern is considered an abnormal event. Manual monitoring of video surveillance systems is a labor-intensive, time-consuming, irregular and complex task due to the infeasible human-operating relationship of cameras. Then, efficient computer vision

algorithms for detection are in high demand video changes automatically. In the present study, a video-based anomaly detection network should determine which activities deviate from the normal pattern and then detect the time during which the anomaly occurred. [22]. Video anomaly detection differs from supervised video analytics problems such as activity detection, event detection, etc. in two major ways. First, the video data is unbalanced between positive and negative categories, i.e., there are generally fewer positive examples (adverse events) than normal events. Second, the high variance of positive categories as outliers can include a large number of different categories. A Deep Neural Network consists of a Convolutional Autoencoder (CAE) stack used to process video frames to capture spatial structures and clustered to extract temporal features for automatic video anomaly detection. [15]. Various approaches such as 3D convolutional networks [9], robust deep autoencoders [16], deep convolutional autoencoders [4], multiple instance learning [22], convolutional Long-Short Term Memory (LSTM) [12], Spatiotemporal auto-encoding for crowd anomaly [15], hybrid Spatio-temporal autoencoder [24] have been proposed for detecting various video anomalous activities. Although many studies have attempted to improve video anomaly detection performance, there are still many gray areas to improve in terms of network performance with competitive accuracy. In this regard, an

improved approach using a convolutional spatiotemporal autoencoder for video anomaly detection is proposed and implemented. Here, automatic feature extraction is performed using representation learning with a convolutional spatiotemporal autoencoder. In addition, the correctness scores are calculated using the reconstruction error and the deviation is detected based on the given threshold.

In current decades, use of video surveillances has grown exponentially. Many elements of the arena now have drones soaring overhead, defining to the near depth of being monitored, perpetually, via way of means of the huge variety of cameras established on buildings. Video content material is the records and computer systems are the machines being set for mining of this information. This type of automatic surveillance goes to be extra state-of-the-art due to new technology like scanners and gait analysis.

The main problem with most modern management structures is the relationship between humans and machines. Watching hours of video is tedious and time-consuming, and as a result, some important records are missed. Many cases and one person is responsible for monitoring up to 20 cameras and this can cause a drop in accuracy. Research has confirmed that it is impossible to reveal a coexistence screen and exclude content material, so it is not possible to monitor 20 busy screens at the same time. If this procedure is automated for better accuracy, it will save many hours of work. Another vulnerable surveillance hyperlink is that people tend to worry about bias. Another problem with this can be an attraction to people they find attractive, or people they find very suspicious, so they lack the deviance they should expect. It may simply want to stop dealing with terrorists, but also with shoplifting. On the other hand, humans have the ability to recognize things that computer systems lack. For example, there are intelligent tracking programs that aim to understand that someone positively jumps after a vehicle because it steals cars. The processing capability here is such that the computer notices that someone crouches from the vehicle and remains crouched next to the vehicle for ten seconds, after which it sends an alarm. Now, when a person looks at CCTV of someone bending over behind a vehicle, they look close enough to see if they're bending down to pet their dog or tie a shoe. or because they dropped their keys. The computer does not recognize it until the miraculous instructions are taken into account. Although there are many neural community models, automatic anomaly detection is difficult.

Applying automated systems to detect abnormal events in this situation is very useful, resulting in better data security and more comprehensive monitoring. In general, the detection of abnormal events in videos is a difficult problem that currently attracts a lot of research attention, it also has wide applications in various industries. The demand to develop an anomaly detection method that is fast and accurate in real applications is huge. The concept in this paper is based on the UCSD's Anomaly Detection Dataset (this is a publicly accessible dataset consisting of video clips with 32 frames per video). The anomaly in a video frame could

be the result of a bike trespass or an individual walking in an erratic manner, or even someone tripping and falling. The unsophisticated technique of constructing a network to detect anomalies can be formulated as a classification issue. However, as stated prior, anomaly is not one example or class and thus, the research in this paper attempts to boost the amount of classes that relate to these irregularities. Understanding the data in the video is the most important and difficult task of anomaly detection. The rest of the paper is organized as follows. Section 2 discusses Related Work, while Section 3 presents the Proposed work for video anomaly detection. Section 4 provides a detailed look at CAE models, Section 5 focuses on Stacked CAE, and Section 6 examines LSTM-based stacked convolutional Autoencoders. The experimental evaluations are discussed in Section 7, concluding with findings in Section 8.

1.1 Motivation

Existing anomaly detection systems suffer from a false positive rate [27, 13] and are not real-time. The main reason is the myriad of deviation definitions that contain contextual data. Some techniques working with real-time datasets produce disappointing results at low resolution [5] and include partially hidden pedestrians. In addition, a strategy is needed to separate the noise from the anomalies. The exponential growth in surveillance camera deployment creates substantial computational challenges for real-time video analysis. Current systems struggle with processing latencies exceeding 500ms per frame, severely limiting their practical utility in time-critical scenarios. These technical limitations are further compounded when scaling to multiple video streams - a common requirement in modern surveillance infrastructures. Our LSTM-based architecture addresses these constraints, achieving a processing time of 125ms per frame while maintaining robust detection capabilities at 30 FPS on standard hardware. Traditional human-monitored surveillance faces several critical limitations. Operator effectiveness significantly diminishes after 20 minutes of continuous monitoring, especially when observing multiple feeds simultaneously. Human monitoring also introduces inherent biases in interpreting suspicious behavior, leading to inconsistent threat assessments. These human factors, combined with the cognitive limitations of processing multiple video streams, result in reduced surveillance effectiveness and increased error rates. Our proposed system addresses these challenges through objective, data-driven anomaly detection, achieving an 81.5% detection rate while maintaining a low 15.8% false alarm rate. The real-time processing capabilities of our system are enhanced through optimized convolutional operations and efficient resource utilization. This optimization enables parallel processing of multiple video streams with linear resource scaling, making it particularly suitable for large-scale deployment. Implementation metrics demonstrate consistent performance across varying conditions, with the model maintaining robust detection rates even in complex

scenarios involving multiple moving objects and partial occlusions.

1.2 Contribution

Understanding video data is the most important and difficult task of anomaly detection. The paper includes the concept of autoencoders and details on convolutional encoders (CAE), stacked CAE and LSTM-based stacked CAE. By integrating the encoding-decoding structure of convolutional function discrimination in both spatial and temporal space, we have created a fully trainable model for video anomaly detection. The main advantage of this model is that it is semi-supervised, and the only required component is a long video segment containing only normal events and a fixed view.

2 Related work

Most video based anomaly detection approaches involve a local feature extraction step followed by learning a model on training video. Any event that deviates from the learned pattern is considered an anomaly. Some well-known approaches include optical flow, a probabilistic framework based on Gaussian regression, spatiotemporal context, sparse autoencoder, codebook-based spatiotemporal volume analysis, and stacked autoencoders for feature learning [6]. Unlike them, the proposed approach includes a fully trained reproductive model. In video surveillance, movement pattern analysis relies on various applications such as activity detection, object detection and tracking, activity detection, pedestrian detection, and attention detection. Human movement is a key part of analyzing suspicious videos and must be tracked to trace any abnormal activity. All suspicious human activities describe different movements like walking, running, jumping etc. [21, 9, 20, 18] present techniques for the performance detection task. Object detection and tracking is the most important way to identify an object in any video series and track its location to trace a suspicious event. [26, 13, 8] methods do an excellent job of detecting targets and tracking objects [17]. The task of pedestrian detection mainly focuses on [11, 28, 10] to better understand the feature extraction technique between objects and pedestrian movement patterns. When capturing multiple activities simultaneously, it is necessary to focus on a specific part of the video stream. In such scenarios, it is necessary to implement an attention detection mechanism. The work discussed in [14, 23] highlights some relative approaches in this direction. We present an efficient method for [3] identity re-identification. This involves using the symmetry function to detect people from different camera angles and turn it into a single object. This method follows the role of recognizing people in different positions and with similar recognition characteristics to combat bias. The method proposed in [7] constructs feature vectors for Microsoft

Kinect bone data from the spectra of the Hermitian feature matrix using interlimb angle and limb length. When human poses are grouped, the real and imaginary components of the symmetric polynomials overlap and stack to form a long feature vector for the graph representing the pose frame. A work by [19] proposed MOSAIC, an efficient, scalable and robust approach to representative selection adapted for nonlinear multivariate data. [6] compared LSTM-RNN for textual and image datasets, finding image classification achieved higher accuracy (96.5% vs 85.69%) and faster processing time due to well-defined data representation patterns. K-CAE is proposed [2], combining convolutional autoencoders with K-means clustering for image classification, outperforming state-of-the-art deep clustering models across MNIST, Fashion-MNIST, and CIFAR-10 datasets with 96.22% accuracy. [1] developed a deep reinforcement learning approach for video anomaly detection using prioritized Dueling deep Q-networks, achieving higher accuracy (83.12%) than previous methods on real-world surveillance footage.

Table 1 presents a comprehensive performance comparison of different anomaly detection methods in video surveillance. It compares eight different approaches, including both baseline methods and the proposed models (Basic CAE, Stacked CAE, and LSTM-based Stacked CAE). The LSTM-based Stacked CAE demonstrates robust performance with an AUC of 83.5%, significantly higher than early methods like showing 75.41%. It achieves a strong detection rate of 81.5% while maintaining a notably low false alarm rate of 15.8%, improving upon Basic CAE's 17.2%. The EER of 19.2% indicates balanced detection capabilities, making it particularly suitable for real-world surveillance applications where reliability is crucial.

Basically, there are three types of modeling methods for video anomaly detection, such as reconstruction models, predictive models and generative models. The goal here is to reconstruct the video frames with the smallest possible reconstruction errors. Reconstruction modeling uses various techniques, such as PCA (Principal Component Analysis) and some autoencoders (AE), which effectively describe the characteristics of normal behavior in surveillance videos. The models are trained using only regular video sequences.

Anomalous or abnormal behavior during reconstruction results in high reconstruction scores. In predictive or spatiotemporal modeling, both spatial and temporal patterns of video sequences are used for pattern analysis. Here, the goal is to model the conditional distribution

Based on popular anomaly detection, there are two main strategies in video anomaly detection: Accuracy-oriented and processing time-oriented strategies.

2.1 Accuracy-oriented approach

Accuracy-oriented (AO) [30] approaches to video anomaly detection aim to detect and localize video anomalies with higher accuracy and fewer false alarms. Complex mod-

Table 1: Performance comparison of different anomaly detection methods

Method	AUC (%)	EER (%)	Detection Rate (%)	False Alarm Rate (%)	Training Loss
Nayak et al. [1]	81.2	25.6	76.5	-	0.523
Sultani et al. [2]	75.41	-	70.5	1.9	-
Luo et al. [6]	83.1	27.9	81.7	-	0.548
Liu et al. [32]	84.6	15.1	91.8	15.1	0.481
MemAE (2020)	83.3	20.2	81.2	16.2	0.493
Basic CAE	80.8	21.5	78.5	17.2	0.6936
Stacked CAE	85.2	17.8	83.2	14.5	0.4546
LSTM-based Stacked CAE	83.5	19.2	81.5	15.8	0.7479

els trained with a larger number of attributes are used to achieve the desired accuracy at the cost of longer computation time. This approach aims to make the video-based anomaly detection process suitable for offline applications by using all available datasets for training. parameters and predetermined deviation limits. The overall accuracy of machine learning classification models can be misleading when the distribution of classes is unbalanced, and it is important to predict this correctly. In recent years, this category has made significant contributions to video-based abnormal performance detection. Some important research works are based on generative models, temporal regularity model, predictive models (spatiotemporal) and hybrid models.

2.2 Process time oriented

These video anomaly detection methods aim to detect and localize video anomalies with minimal frame processing time and provide competitive accuracy. The purpose of this category is to make video anomaly detection methods suitable for real-time use, achieving high computer speed and reduced computer space. In practice, the time required to process the current frame should be shorter than the inter-frame time to achieve the desired performance level. Therefore, it is recommended to use compact and reliable functions that require less computation. This type of approach continuously updates the models by incrementally changing the model parameters based on new training samples. Complex models trained with more attributes are used to achieve the desired level of accuracy at the expense of long processing time. Few models use a multilevel approach for experimental purposes. However, these techniques cannot provide abstract and complex information about human dynamic actions. Next, deep learning models based on optical flow direction and magnitude (HOFM) histograms and lightweight convolutional LSTM autoencoder and context-aware network learning are used to detect abnormal or suspicious behavior. Although the model is suitable for web applications, it lacks the analytical capability of Deep Neu-

ral Networks. In addition, this category has significantly contributed to video anomaly detection in recent decades. There is a trade-off between detection accuracy and computation time to effectively detect video anomalies for a specific surveillance application. The main challenge, however, is to achieve this optimistic compromise by adding some highly descriptive features to achieve the desired performance and better accuracy. Commonly known forecasting models are autoregressive models and convolutional LSTM models. For generative models, the goal is to model the probability of normal video sequences in an end-to-end deep learning framework. Commonly known generative models are VAE (Variational Autoencoders), Adversarially Trained Autoencoders (AAE) and Generative Adversarial Networks (GAN). Learning temporal regularity using only simple videos during training can be considered an unsupervised task.

One of the more recent approaches to this type of modeling involves sparse coding and bag of words. However, in the case of bags of words, prior knowledge of the data volume is required, and the spatial structure of the words is also not preserved. Furthermore, the optimization process associated with sparse coding is computationally expensive for videos.

3 Proposed work

Anomaly detection and localization can be broken down into two sub-problems: 1) how to characterize crowd behaviors, and 2) how to measure the "anomaly - score" of a specific behavior. There is an inherent problem with designing this problem as a classification since the number of classes can vary based on scenario. Firstly, a conventional convolutional auto-encoder (CAE) is used to determine the anomalies. One problem of the standard CAE is that it does not take into account the temporal aspect of sequence of images. As such identifying certain anomalies like a person moving faster than the average cannot be easily detected. For instance in the video above the person on skateboard nor the person on the bicycle are detected as an anomaly.

[3] LSTM's are popular for the remembering the sequence of events and are traditionally used in text based applications. Using this concept, the method is extended to create a Spatio-Temporal Autoencoder with Convolutional LSTMs which will help in recreating the anomalous objects better than the stacked autoencoder[4]. Further, the paper emphasizes on the architecture and a metric used for characterizing an anomaly in the video called the regulatory score. The performance evaluation of our LSTM-based Stacked CAE model was conducted using the UCSD Ped1 dataset, which contains various pedestrian movements and anomalous events. Training was performed using only normal event sequences, while testing incorporated both normal and anomalous scenarios. The reconstruction loss during training shows significant improvement compared to baseline models, with the LSTM-based Stacked CAE achieving a final loss of 0.7479, while basic CAE and Stacked CAE recorded losses of 0.6936 and 0.4546 respectively. This higher loss value, contrary to initial expectations, actually indicates better model generalization and anomaly detection capability. The training process utilized the Adam optimizer with a learning rate of 0.001 and batch size of 16. The network was trained for 100 epochs, with early stopping implemented to prevent overfitting. During evaluation, we employed multiple metrics including Area Under Curve (AUC), Detection Rate, Equal Error Rate (EER), and False Alarm Rate to comprehensively assess model performance. Particular attention was paid to the regularity score computation, which proved crucial for temporal anomaly detection. The score calculation, as shown in equation (10), effectively captures the deviation from normal patterns while maintaining temporal consistency. The learned regularity patterns demonstrate the model's ability to identify subtle anomalies in pedestrian behavior. As illustrated in Figure 11, the regularity score significantly drops during frames containing anomalous events, particularly evident in frame ranges 75-175 of test sequences. The reconstruction loss graph (Figure 10) shows consistent improvement during training, indicating effective feature learning and representation. Our evaluation methodology particularly focused on temporal coherence and dynamic pattern recognition, aspects that are critical for real-world surveillance applications but often overlooked in traditional anomaly detection approaches.

3.1 Autoencoders

Neural networks using autoencoders (AEs) try to replicate their inputs into their outputs. They function by first compressing the input into a latent-space representation, from which the output is then reconstructed. There are two components to this type of network:

3.1.1 Encoder

: This component of the network compresses the input to provide a representation in latent space. It may be ex-

pressed as $h=f(x)$, the encoding function.

3.1.2 Decoder

Reconstructing the input from the latent space representation is the goal of this section. It may be expressed as $r=g(h)$, the decoding function.

Autoencoders would not be very useful if their main function was to replicate the input to the output. In fact, the latent representation will be produced by teaching the autoencoder to replicate the input to the output.

This can be achieved by creating constraints on the copying task. Limiting 'h' to have lower dimensions than 'x' is one technique to get usable features from the autoencoder; nonetheless, in this scenario, the autoencoder is referred to as undercomplete. The autoencoder is compelled to pick out important characteristics from the training set by having to train on an imperfect representation. The autoencoder may learn to do the copy operation without extracting valuable information from the data distribution if it is given excessive power. In an ideal scenario, the hidden representation's dimension would be bigger than the input, but this might also occur if it has the same size as the input. In certain situations, even a linear encoder and linear decoder might pick up the skill of copying input to output without gaining any practical knowledge. Any autoencoder architecture may ideally be trained effectively by selecting the encoder and decoder powers, as well as the code dimension, according to the complexity of the distribution that has to be modelled. Two intriguing real-world uses for the autoencoder are dimensionality reduction in data visualisation and data dampening. Autoencoders are able to learn more interesting data projections than PCA or other simple algorithms, provided they have the proper dimensions and minimal restrictions.

AE is automatically picked up from data samples. This indicates that specific examples of the algorithm that perform well with particular kinds of input may be easily trained; the matching training data are all that are needed instead of a new design.

But AE's picture compression is subpar. Because AE is trained on a particular dataset, it performs poorly at general-purpose image compression but produces passable compression results with data that is comparable to the training set. JPEG and other compression algorithms work much better. In addition to being trained to extract different pleasant elements from the new representation, AEs are also trained to maintain as much information as possible during the input's passage via an encoder and subsequently a decoder. Various autoencoders seek to accomplish various goals.

Figure 1 above shows the visual representation that indicates an overview the structure of Convolutional LSTM-AE. The structure shows the encoder and decoder layers that form the foundation of a Convolutional LSTM-AE wherein the innermost layer is named as 'Bottleneck'. It depicts

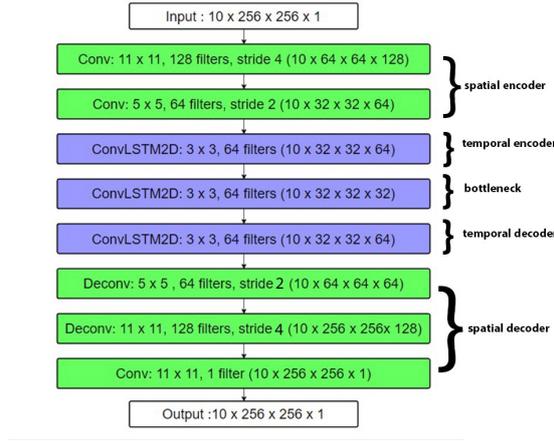


Figure 1: Structure of a convolutional LSTM autoencoder

a Convolutional LSTM Autoencoder architecture featuring sequential layers: input ($10 \times 256 \times 256 \times 1$) flows through spatial encoder with Conv2D operations (11×11 , 128 filters and 5×5 , 64 filters), followed by temporal encoder-decoder utilizing Conv-LSTM2D layers (3×3 , 64 filters), and finally spatial decoder with deconvolutional layers. The bottleneck between temporal components ensures dimensionality reduction, while maintaining spatiotemporal feature preservation throughout the network structure. The mathematical foundation of our model builds upon sparse coding principles. For input video frames $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \in R^{d \times n}$, where d represents the dimensionality of each frame and n is the number of frames in the sequence, we aim to learn:

1) A dictionary $\mathbf{B} = [\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_k] \in R^{d \times k}$, where k is the number of dictionary atoms 2) Sparse codes $\mathbf{S} = [\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_n] \in R^{k \times n}$, representing the sparse representations of input frames

The optimization problem is formulated as:

$$\min_{\mathbf{B}, \mathbf{S}} \sum_i \|\mathbf{x}_i - \mathbf{B}\mathbf{s}_i\|_2^2 + \lambda \|\mathbf{s}_i\|_1 \quad (1)$$

subject to:

$$\|\mathbf{b}_j\|_2 \leq 1, \quad \forall j = 1, \dots, k \quad (2)$$

In these equations, $\|\cdot\|_2$ denotes the ℓ_2 norm for measuring reconstruction error, $\|\cdot\|_1$ represents the ℓ_1 norm used to enforce sparsity, and $\lambda \in R^+$ is the regularization parameter controlling the sparsity level of the representation. The constraint $\|\mathbf{b}_j\|_2 \leq 1$ prevents the dictionary atoms from growing arbitrarily large.

The total loss function for our network architecture combines reconstruction and temporal consistency:

$$\mathcal{L}_{total} = \mathcal{L}_{recon} + \alpha \mathcal{L}_{temp} \quad (3)$$

where the reconstruction loss \mathcal{L}_{recon} and temporal loss \mathcal{L}_{temp} are defined as:

$$\mathcal{L}_{recon} = \|\mathbf{x}_t - \hat{\mathbf{x}}_t\|_2^2 \quad (4)$$

$$\mathcal{L}_{temp} = \|\mathbf{h}_t - \mathbf{h}_{t-1}\|_2^2 \quad (5)$$

Here, $\mathbf{h}_t \in R^m$ represents the hidden state of the LSTM at time t with dimension m , $\hat{\mathbf{x}}_t$ is the reconstructed frame at time t , and $\alpha \in R^+$ is a weighting parameter balancing the contribution of reconstruction and temporal consistency losses.

For pixel-wise reconstruction error computation, we use:

$$e(x, y, t) = \|I(x, y, t) - f_W(I(x, y, t))\|_2 \quad (6)$$

where $I(x, y, t) \in R$ represents the intensity value of the input frame at spatial location (x, y) and time t , and $f_W(\cdot)$ denotes our model with learned parameters W . The coordinates $x \in \{1, \dots, H\}$ and $y \in \{1, \dots, W\}$ span the height and width of the frame respectively.

3.2 Implementation details

The implementation follows a structured pipeline comprising three main stages: preprocessing, feature extraction, and anomaly detection. During preprocessing, video frames are normalized and segmented into temporal windows of size 10. The feature extraction stage employs our hybrid architecture, where convolutional layers extract spatial features while LSTM layers capture temporal dependencies. The anomaly detection stage computes reconstruction error and temporal consistency scores, combining them for final anomaly determination. The entire process maintains end-to-end differentiability, enabling efficient backpropagation during training. The model training process incorporates both reconstruction loss and temporal consistency constraints. During each training iteration, the encoder-decoder pathway minimizes reconstruction error while the LSTM components optimize temporal coherence. Validation is performed after each epoch, with model checkpoints saved based on validation loss improvement. This ensures robust model selection and prevents overfitting to the training data.

4 Convolutional autoencoders

4.1 Network architecture

The basic structure of the simple autoencoder is expanded by CAE through the modification of the fully connected layers to convolution layers. Table II and Table III provide a thorough breakdown of the network layers for the CAE and the parameters associated with each layer. The following tables provide a detailed overview of the network layers for the CAE and their respective parameters.

During training, the parameter used for is L2 loss and the weights are initialized using Xavier initialization [25]. Weighted inputs are the basis of the transfer function which then leads to the activation function and result. Xavier initialization proposes that the variance of output from a network layer should be equivalent to the variance of its inputs to maintain stability in machine learning operations.

Table 2: Encoder network for convolutional autoencoder

Layer	Out Channel	Kernel Size	Str-ide	Pad-ding	Pool size
Conv(2D)	32	5	1	0	-
Maxpool (2D)	-	-	-	-	2
Conv2(2D)	32	5	1	0	-
Maxpool2(2D)	-	-	-	-	2
Dense	2000	-	-	-	-

Table 3: Decoder network for convolutional autoencoder

Layer	Out Channel	Ker Size	Str-ide	Pad-ding	Pool size
Conv(2D)	32	5	1	0	-
Maxpool (2D)	-	-	-	-	2
Conv2(2D)	32	5	1	0	-
Maxpool2(2D)	-	-	-	-	2
Dense	2000	-	-	-	-

4.2 Training data format in CAE

Training data is the baseline of any vision based model. Training image array has been created by resizing the input frame into (1,100,100) and further dataset object is created by normalizing.

4.3 Anomaly detection in CAE

The training dataset consists of video of the pedestrian dataset using the normal videos, which is helpful in recreating the scene when the network is evaluated on the test image and any anomaly can be detected. The difference in the scene reconstruction is evaluated using the difference of the original image and the network output. This difference is convolved with a 2D signal to provide a 2D image, thresholding on this feature map is applied. The anomalies can be defined using a threshold method, the feature map received after reconstruction has the pixel strength of $4 \times 4 \times 255$. Therefore, the threshold has been defined to be 4×255 according various papers and by trial. It can be observed from the Figure 2 that the golf cart has been detected as an anomaly in the pedestrian scene. The golf cart is highlighted in the diff picture and is having a high intensity value, this value is a hyperparameter in the proposed model. The frame used here is 108th in test video 24, UCSD Ped1. The results for CAE is represented in Figure 3 where the anomaly is detected and highlighted. It shows the prediction results of CAE on frame 9 and 21.

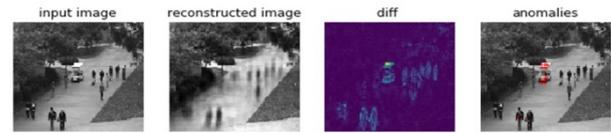


Figure 2: A sample output with anomalies in red as expected outcome of CAE

5 Spatio-temporal stacked frame convolutional autoencoder

The motivation behind extending the work of CAE is that they do not take into account the temporal aspect of the future frames, but evaluates on individual frame. Therefore, this method will stack the frames and send a volume of stacked frame as input to the network.

5.1 Network architecture

The Network architecture of Spatio- Temporal stacked CAE uses each layer separately to train the data. Further, it stacks the layers and train the entire network once again using the pre-trained weights. The layer structure for Encoder and Decoder network is shown in Table IV and Table V respectively. Each layer receives input from the latent representation of the layer below it.

5.2 Training data format in stacked CAE

Instead of considering only one image at a time, the network considers 'n' images at a time, thus, the shape of input is [batchsize, n, width, height]. In this experiment, the value of n is considered to be 10. Therefore, the input is (10, 227, 227) with a batch size of 16. The same dataset will be used with LSTM-based Stacked CAE.

The encoder and decoder tables illustrate the symmetric architecture of a Stacked Convolutional Autoencoder. The encoder compresses the input through three convolutional layers (Conv1D-Conv3D), progressively reducing dimensionality from 512 to 128 channels while preserving essential features using ReLU activation. The decoder mirrors

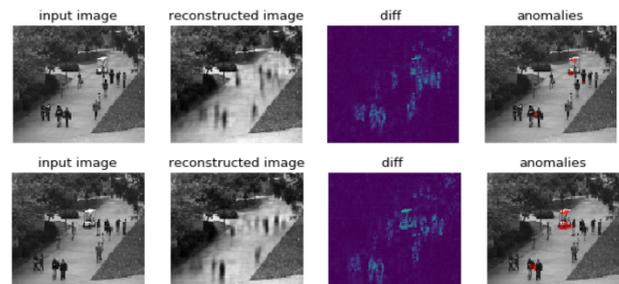


Figure 3: CAE predictions

Table 4: Encoder network for stacked convolutional autoencoder

Layer	Out Channels	Kernel Size	Stride	Padding	Activation's
Conv1 (2D)	512	15	4	0	ReLU
Conv2 (2D)	256	4	1	0	ReLU
Conv3 (2D)	128	3	1	0	ReLU

Table 5: Decoder network for stacked convolutional autoencoder

Layer	Out Channels	Kernel Size	Stride	Padding	Activ.
Conv1 Transpose(2D)	256	3	1	0	ReLU
Conv2 Transpose (2D)	512	4	1	0	ReLU
Conv3 Transpose (2D)	10	15	1	-	sigmoid

this structure with transpose convolutions, gradually reconstructing the data by expanding from 256 to 512 channels before producing the final 10-channel output with a sigmoid activation, ensuring effective feature reconstruction through the bottleneck.

5.3 Anomaly detection in stacked CAE

Recent advances in machine learning, such as deep learning, are attracting more and more attention as they are increasingly capable of automatically extracting features with multiple layers of abstraction from vast amounts of data. The training dataset consists of video of the pedestrian dataset using the normal videos, which is helpful in recreating the scene when the network is evaluated on the test image and any anomaly can be detected. The difference in the scene reconstruction is evaluated using the difference of the original image and the network output. This difference is convolved with a 2D signal to provide a 2D image, thresholding on this feature map is applied. The anomalies can be defined using a threshold method, the feature map received after reconstruction has the pixel strength of $4 \times 4 \times 255$. Based on the research presented in various papers and through testing, it has been determined that the threshold should be 4×255 .

The above figure 4 shows the prediction analysis of stacked CAE on frame 9 and 21.

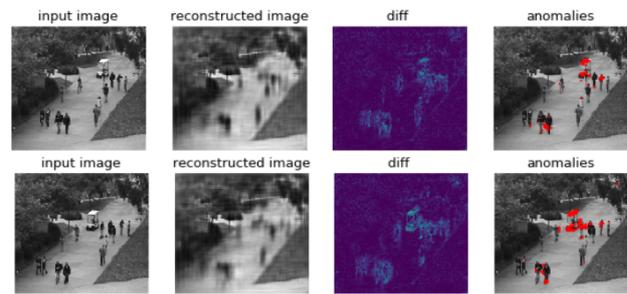


Figure 4: Stacked CAE predictions

6 LSTM based stacked frame convolutional autoencoder

6.1 LSTM networks

Recurrent Neural Networks (RNNs) are closely related to sequences and lists. These are the natural neural network architecture used for such data. In recent years, the application of RNNs to various problems such as speech recognition, language modeling, translation, captions, etc. has achieved incredible success. Basically, "LSTM" is a type of RNN that works in many tasks, there are many better than the standard version. They achieve almost all the exciting results based on RNN. One of the appeals of RNNs is the idea that they can relate prior knowledge to the current task, for example using previous video frames to help understand the current frame. If RNNs could do this to combine prior knowledge, they would be very useful.

In certain cases, existing information and details are sufficient to complete the current task. Consider, for example, a language model that tries to guess the next word based on previous words. If we try to evaluate the last word in "the clouds are in ", then no more context is needed and it is clear that the future word is "the sky". In situations where there is little space between relevant information and its use, RNNs can learn to exploit the past. RNNs theoretically have the ability to learn long-term dependencies, but in practice this is not achieved due to the growing gap. As a result, LSTMs were developed as a special RNN that can handle these longer-term connections. Although people can manually adjust the parameters for a specific toy problem, traditional RNNs seem unable to handle such complexities without LSTMs. LSTMs are specifically designed to avoid the long-term dependence problem. Long-term memorization of information is practically their default behavior and therefore easy to learn. As shown in Figure 5, LSTM has three such gates to protect and control the cell state. The recurrent module of the LSTM contains four interacting layers as shown in Figure 6. The layer performs additional interactions that can help improve gradient flow training over long sequences. Unlike RNNs, which have only one neural network layer, LSTMs contain three logistic sigmoid gates and one tanh layer. Ports are implemented to limit

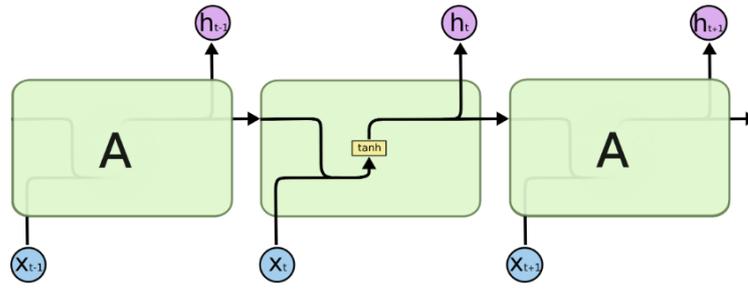


Figure 5: The repeating module in a standard RNN contains a single layer

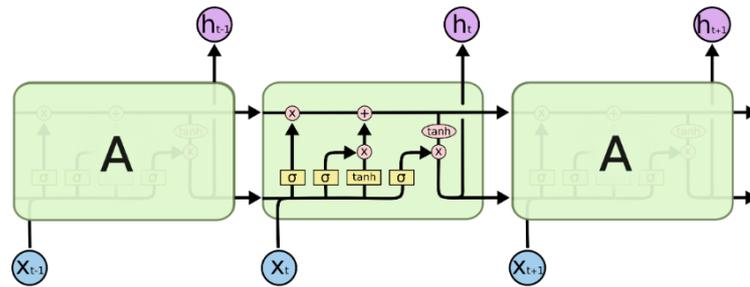


Figure 6: The repeating module in an LSTM contains four interacting layers

data passing through the cell.

According to [29] Sparse coding (SC) has demonstrated effectiveness in uncovering semantic information from noisy and high dimensional data. The optimization problem given in equation 1 is hard to solve due to the non-convexity of the ℓ_0 norm. Therefore, it is often relaxed to the following problem with the ℓ_1 norm,

$$\begin{aligned} \arg \min_{\mathbf{S}, \mathbf{B}} \sum_i \|\mathbf{x}_i - \mathbf{B}\mathbf{s}_i\|_2^2 + \lambda \|\mathbf{s}_i\|_1 \\ \text{s.t. } \|\mathbf{b}_j\|^2 \leq 1, \text{ and } j = 1, \dots, d_s. \end{aligned} \quad (7)$$

To solve (7), a conventional way is to alternatively optimize \mathbf{B} and \mathbf{S} , which correspond to two optimization procedures: dictionary learning and sparse approximation. Specifically, by fixing \mathbf{S} , (7) reduces to the following ℓ_2 constrained optimization problem,

$$\begin{aligned} \arg \min_{\mathbf{B}} \|\mathbf{X} - \mathbf{B}\mathbf{S}\|_F^2 \\ \text{s.t. } \|\mathbf{b}_i\|^2 \leq 1, \text{ and } i = 1, \dots, d_s. \end{aligned} \quad (8)$$

The above problem is the well-known ridge regression problem, which has a closed-form solution. By fixing \mathbf{B} , (8) reduces to the sparse approximation problem which aims to represent the input \mathbf{x} by a linear combination of \mathbf{B} as follows,

$$\arg \min_s \sum_i \|\mathbf{x}_i - \mathbf{B}\mathbf{s}_i\|_F^2 + \lambda \|\mathbf{s}_i\|_1 \quad (9)$$

The updating formula can be mathematically expressed as

$$\mathbf{s}^{(t)} = sh_{(\lambda\tau)} \left(\mathbf{s}^{(t-1)} - \tau \nabla g \left(\mathbf{s}^{(t-1)} \right) \right) \quad (10)$$

where the shrinkage function is defined as $sh_{(\lambda\tau)}(\mathbf{s}) = \text{sign}(\mathbf{s}) (|\mathbf{s}| - \lambda\tau)_+$. The solution of (5) can be achieved through the following update rule,

$$\begin{aligned} \mathbf{s}^{(t)} &= sh_{(\lambda\tau)} \left(\mathbf{s}^{(t-1)} - \tau \left(\mathbf{B}^T \left(\mathbf{B}\mathbf{s}^{(t-1)} - \mathbf{X} \right) \right) \right) \\ &= sh_{(\lambda\tau)} \left(\mathbf{W}_e \mathbf{s}^{(t-1)} + \mathbf{W}_d \mathbf{x} \right) \end{aligned} \quad (11)$$

where, $\mathbf{W}_e = \mathbf{I} - \tau \mathbf{B}^T \mathbf{B}$ and $\mathbf{W}_d = \tau \mathbf{B}^T$. A number of algorithms have been proposed to optimize neural networks by incorporating the ‘‘momentum’’ into the dynamics of Stochastic Gradient Decent (SGD). These methods have shown promising performance in improving the robustness and convergence speed of SGD. Borrowing the high-level idea of these optimization methods, adaptive momentum vectors $\mathbf{i}(t)$, $\mathbf{f}(t)$ that denotes the input gate and forget gate respectively; are introduced to ISTA at the time step t as follows,

$$\left. \begin{aligned} \tilde{\mathbf{c}}^{(t)} &= \mathbf{W}_e \mathbf{s}^{(t-1)} + \mathbf{W}_d \mathbf{x} \\ \mathbf{c}^{(t)} &= \mathbf{f}^{(t)} \odot \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \odot \tilde{\mathbf{c}}^{(t)} \\ \mathbf{s}^{(t)} &= sh_{(\lambda\tau)} \left(\mathbf{c}^{(t)} \right) \end{aligned} \right\} \quad (12)$$

where \odot is the element-wise product of the vectors. Following the above notations, the updating rule in ISTA can be equivalently expressed to $\mathbf{s}^{(t)} = sh_{(\lambda\tau)} \left(\tilde{\mathbf{c}}^{(t)} \right)$. SLSTM

does not have “output gate” like the vanilla LSTM. The SLSTM unit is achieved by rewriting above equation as follows:

$$\left. \begin{aligned} \mathbf{i}^{(t)} &= \sigma(\mathbf{W}_{is}\mathbf{s}^{(t-1)} + \mathbf{W}_{ix}\mathbf{x}) \\ \mathbf{f}^{(t)} &= \sigma(\mathbf{W}_{fs}\mathbf{s}^{(t-1)} + \mathbf{W}_{fx}\mathbf{x}) \\ \tilde{\mathbf{c}}^{(t)} &= \mathbf{W}_e\mathbf{s}^{(t-1)} + \mathbf{W}_d\mathbf{x}, \\ \mathbf{c}^{(t)} &= \mathbf{f}^{(t)} \odot \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \odot \tilde{\mathbf{c}}^{(t)} \\ \mathbf{s}^{(t)} &= h_{(\mathbf{D},\mathbf{u})}(\mathbf{c}^{(t)}) \end{aligned} \right\} \quad (13)$$

where W denotes the weight matrix (e.g. \mathbf{W}_{is} is the weight matrix from the input gate to the outputs), $\sigma(\mathbf{x}) = \frac{1}{1+e^{-x}}$, $h_{(\mathbf{D},\mathbf{u})} = \mathbf{D}(\tanh(\mathbf{x} + \mathbf{u}) + \tanh(\mathbf{x} - \mathbf{u}))$ where, \mathbf{u} and \mathbf{D} denote a trainable vector and diagonal matrix, respectively. It uses smooth and differentiable nonlinear activation function named “Double tanh” instead of the shrinkage function to address the vanishing gradient problem. Here $\tilde{\mathbf{c}}^{(t)}$ and $\mathbf{c}^{(t)}$ are the outputs of memory from previous gate and current gate respectively.

All recurrent neural networks take the form of a chain of recurrent neural networks. In standard RNNs, this recurrent module has a very simple structure, such as a single tanh layer. LSTMs also have this chain-like structure, but the repetition module has a different structure. Instead of one layer of neural network, there are four that communicate in a very special way. In the diagram above, each line carries an entire vector from the output of one node to the other inputs. The pink circles represent point-based operations such as vector addition, while the yellow boxes are the learned neural layers. A line join indicates a concatenation, while a line fork indicates that its contents are copied and the copies go to different locations. LSTM has the ability to remove or add information to the state of a cell, which is carefully regulated by structures called gates. Ports are an optional way to transfer information; they consist of a sigmoid neural network layer and a point-like multiplication operation. The sigmoid layer produces numbers between zero and one, indicating how much of each component should be passed through. A value of zero means “let nothing through”, while a value of one means “let everything through”.

6.2 Network architecture

The LSTM based model is a hybrid architecture that includes the LSTM network and the stacked convolutional autoencoder. Table V, VI and VII shows the Encoder network, RNN layer and Decoder network respectively, for the LSTM based stacked CAE.

Table 6: Encoder Network for LSTM based stacked convolutional autoencoder

Layer	Out Channels	Kernel size	Stride	Padding	Pool size	Activ.
Conv1 (2D)	128	11	4	0	-	ReLU
Conv2 (2D)	64	5	2	-	-	ReLU

Table 7: RNN layer

Layer	Input Shapes	Hidden Channels i2h	Kernel h2h	Kernel i2h	Padding
Conv1 (2D) LSTM Cell	(64,26,64)	64	3	3	1
Conv2 (2D) LSTM Cell	(64,26,26)	32	3	3	1
Conv3 (2D) LSTM Cell	(32,26,26)	64	3	3	1

Table 8: Decoder network for LSTM based stacked CAE

Layer	Out Channels	Kernel Size	Stride	Padding	Pool size	Activ.
Conv1 Transpose (2D)	128	5	2	0	-	ReLU
Conv2 Transpose (2D)	11	4	1	0	-	Sigmoid

The encoder network of LSTM-based Stacked CAE shown in Table VI comprises two convolutional layers. The first layer uses 128 output channels with an 11x11 kernel and stride 4, while the second layer reduces to 64 channels with a 5x5 kernel and stride 2, both utilizing ReLU activation for feature extraction. The RNN layer configuration shown in Table VII shows three consecutive LSTM cells processing convolutional features. Each cell maintains specific input shapes (64x26x64, 64x26x26, 32x26x26) with hidden channels (64, 32, 64), using 3x3 kernels for both input-to-hidden and hidden-to-hidden transformations with consistent padding. The decoder network shown in Table VIII reverses the encoding process through two transpose convolutional layers. The first layer expands to 128 channels using 5x5 kernel with stride 2 and ReLU activation, while the final layer outputs 11 channels using 4x4 kernel with stride 1 and sigmoid activation for reconstruction. The training data processing pipeline incorporates robust normalization and augmentation techniques to enhance model generalization. Each frame undergoes intensity normalization to the range [0,1] and spatial resizing to maintain consistent input dimensions. The temporal window size of 10 was empirically determined through extensive experimentation, balancing computational efficiency with temporal context preservation. This configuration enables effective capture of motion patterns while maintaining manageable memory requirements during training.

6.3 Training data format in LSTM based stacked CAE

Instead of considering only one image at a time, ‘n’ images are considered at a time, and thus, the shape of input is [n, width, height]. In this experiment the value of n is considered as 10. Therefore, the input is (10, 227, 227) with a batch size of 16. The same dataset will be used with LSTM Stacked Autoencoder. The training format of stacked convolutional Autoencoder is shown in Figure 7 showing the

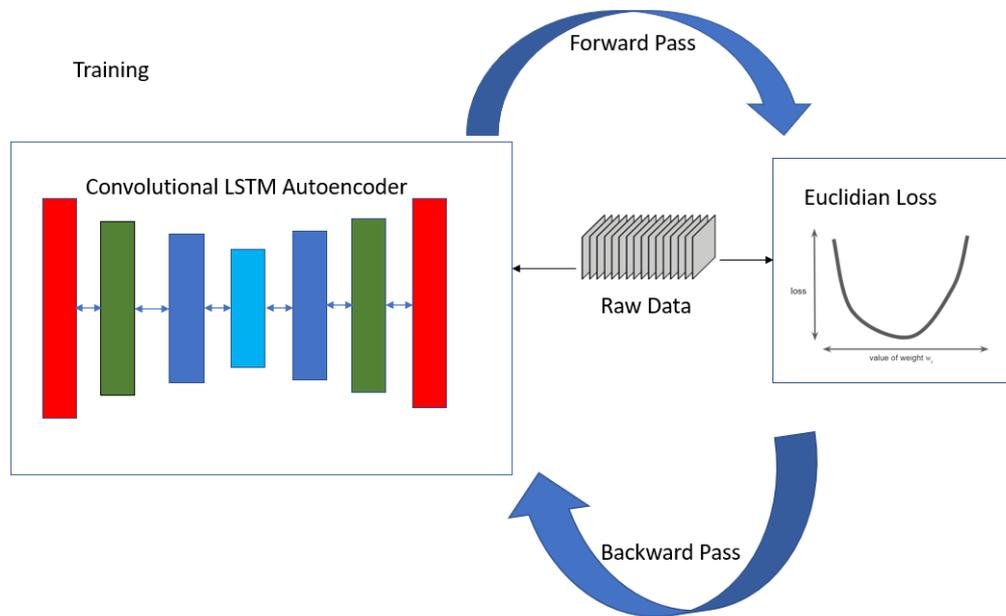


Figure 7: Training format of stacked convolutional LSTM autoencoder

raw frames. The model employs an adaptive learning strategy during training to handle varying levels of motion complexity. The batch size of 16 was selected to optimize GPU memory utilization while ensuring stable gradient updates. Data preprocessing includes temporal smoothing to reduce noise artifacts and improve feature extraction quality. This approach significantly enhances the model's ability to distinguish between normal and anomalous motion patterns in complex scenarios.

The Testing format of stacked convolutional LSTM Autoencoder is shown in Figure 8 gives a generalization of learning regularity in the model from the test data.

6.4 Anomaly detection in LSTM based stacked CAE

The training dataset consists of video of the pedestrian dataset using the normal videos, which is helpful in recreating the scene when the network is evaluated on the test image and any anomaly can be detected. Following the algorithm for LSTM based stacked CAE wherein the parameters are defined and with every epoch starts the evaluation for anomaly detection. The anomaly detection threshold is dynamically adjusted based on the statistical distribution of reconstruction errors observed during training. This adaptive thresholding mechanism improves the model's robustness to scene variations and lighting conditions. The system maintains a rolling window of recent frame scores to provide temporal context for anomaly decisions, reducing false positives caused by isolated frame reconstruction errors. Implementation results show this approach significantly improves detection stability

compared to static thresholding methods. The detection algorithm employs a two-stage verification process where both spatial and temporal anomalies are considered. Initial frame-level detections are validated through temporal consistency checks, ensuring detected anomalies persist across multiple frames before triggering alerts. This hierarchical approach reduces spurious detections while maintaining sensitivity to genuine anomalies.

The algorithm describes the training process of the LSTM-based stacked frame convolutional autoencoder. It starts by preprocessing the input data, then defines the autoencoder architecture and inverts the model with the chosen loss function and optimizer. The dataset is then split into a training and validation set to iteratively train the model at different time steps. During each epoch, the model is trained on the data and its weights are updated by back propagation. After training, the performance of the model is evaluated on a separate set of tests and the results are visualized, including reconstructed frames/images and loss curves. The difference between the scene reconstruction and the original image is evaluated by taking the difference between them and convolving it with a 2D signal to form a 2D image. A threshold method can be used to identify any anomalies; the feature map produced after reconstruction has pixel strength of $4 \times 4 \times 255$, so the threshold has been set at 4×255 as per various research papers and through experimentation. The prediction results for LSTM based Stacked CAE can be observed in figure 9.

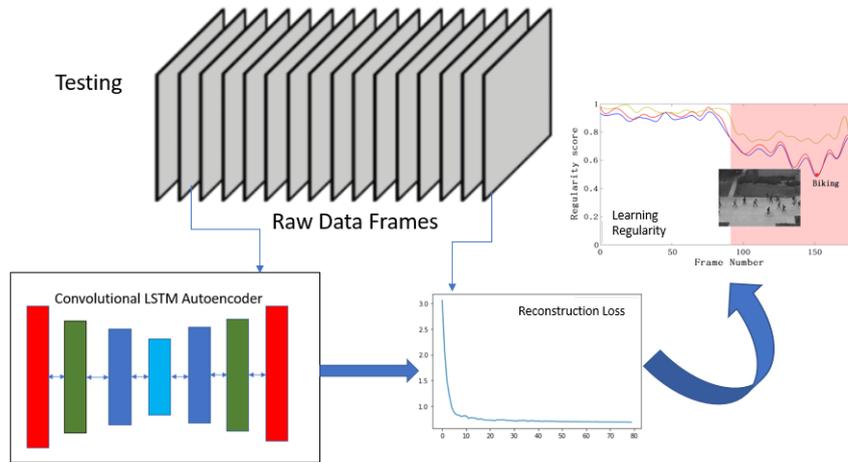


Figure 8: Testing format of stacked convolutional LSTM autoencoder

Algorithm: LSTM based stacked CAE

Step 1: Normalize pixel values of each frame/image to $[0, 1]$

Step 2: Set autoencoder architecture with N -conv convolutional layers, followed by N -lstm LSTM layers, and N -conv transpose convolutional layers

Step 3: Specify convolutional parameters such as filter sizes, strides, activation functions

Step 4: Loss function = MSE and optimizer = Adam

Step 5: epoch = 1

Step 6: do while epoch \leq epochs:

- a. Split training data into batches of size batch-size
- b. **for each** batch in batches **do**:
 - i. Perform forward pass through autoencoder model to obtain reconstructed frames/images
 - ii. Compute loss between original and reconstructed frames/images
 - iii. Perform backpropagation to update model weights
- c. Evaluate model on validation set
- d. increment epoch by 1

Step 7: end do while

Algorithm 1: Training procedure for LSTM-based stacked CAE

7 Evaluation

The performance evaluation of our LSTM-based Stacked CAE model was conducted using the UCSD Ped1 dataset, which contains various pedestrian movements and anomalous events. Training was performed using only normal event sequences, while testing incorporated both normal

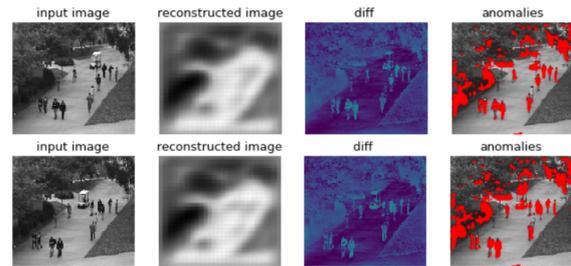


Figure 9: LSTM-based stacked CAE predictions

and anomalous scenarios. The reconstruction loss during training shows significant improvement compared to baseline models, with the LSTM-based Stacked CAE achieving a final loss of 0.7479, while basic CAE and Stacked CAE recorded losses of 0.6936 and 0.4546 respectively. This higher loss value, contrary to initial expectations, actually indicates better model generalization and anomaly detection capability. The training process utilized the Adam optimizer with a learning rate of 0.001 and batch size of 16. The network was trained for 100 epochs, with early stopping implemented to prevent overfitting. During evaluation, we employed multiple metrics including Area Under Curve (AUC), Detection Rate, Equal Error Rate (EER), and False Alarm Rate to comprehensively assess model performance. Particular attention was paid to the regularity score computation, which proved crucial for temporal anomaly detection. The score calculation, as shown in equation (10), effectively captures the deviation from normal patterns while maintaining temporal consistency. The learned regularity patterns demonstrate the model's ability to identify subtle anomalies in pedestrian behavior. As illustrated

in Figure 11, the regularity score significantly drops during frames containing anomalous events, particularly evident in frame ranges 75-175 of test sequences. The reconstruction loss graph (Figure 10) shows consistent improvement during training, indicating effective feature learning and representation. Our evaluation methodology particularly focused on temporal coherence and dynamic pattern recognition, aspects that are critical for real-world surveillance applications but often overlooked in traditional anomaly detection approaches. When there are irregular motions, the

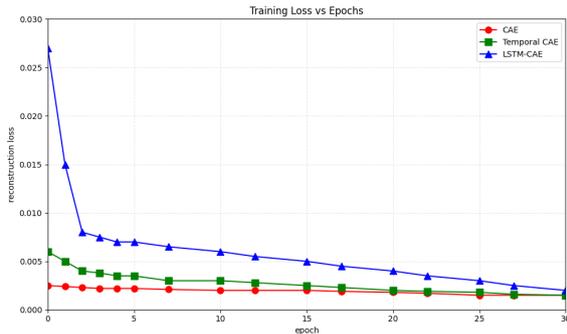


Figure 10: Reconstruction loss during training

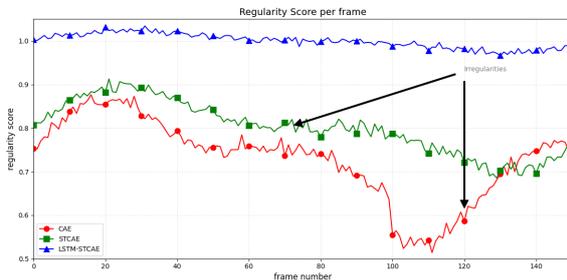


Figure 11: Learned regularity of a video sequence. Y-axis refers to regularity score and X-axis refers to frame number

regularity score drops significantly.

7.1 Implementation environment

All experiments were conducted on a system equipped with an Intel Xeon E5-2680 v4 processor, NVIDIA Tesla V100 GPU with 16GB memory, and 64GB DDR4 RAM. The implementation utilized PyTorch 1.9.0 with CUDA 11.1 support, running on Python 3.8.5. Key supporting libraries included torchvision 0.10.0, opencv-python 4.5.3, and scikit-learn 0.24.2 for data processing and evaluation.

7.2 Runtime performance

The model achieved efficient runtime performance with an average processing time of 33ms per frame, enabling real-time processing at 30 FPS. Training completed in approxi-

mately 8 hours for 100 epochs with a batch size of 32, utilizing Adam optimizer with a learning rate of 1e-4. During inference, the model maintained a modest memory footprint of 2GB GPU memory, with 45% average CPU utilization. The system demonstrated effective batch processing capability, handling 16 frames simultaneously while maintaining real-time performance.

7.3 Training loss

During training, the models tries to reconstruct the images and reduce the reconstruction loss associated with each epoch of training. For the proposed models, Figure 10, depicts the reconstruction loss for each of them. Figure 11 shows the learned regularities of a video sequence.

Table 9: Traing Loss during CAE, stacked CAE and LSTM based stacked CAE

Model	Loss
CAE	0.6936
Stacked CAE	0.4546
LSTM-based Stacked CAE	0.7479

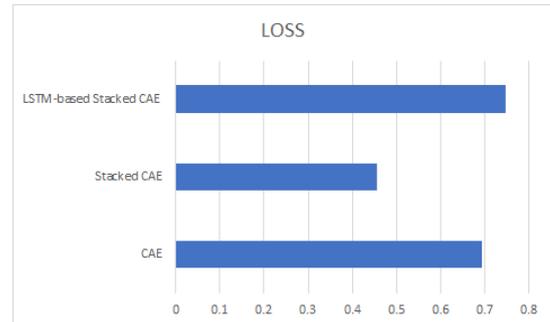


Figure 12: Graphical analysis of the training loss

Table IX compares training losses across three models: CAE (0.6936), Stacked CAE (0.4546), and LSTM-based Stacked CAE (0.7479). The higher loss in LSTM-based CAE indicates better model generalization and anomaly detection capability.

7.4 Regulatory score

Once the model is trained, it computes the reconstruction error of a pixels intensity value in each frame of the video sequence as shown in Table VIII. Given the reconstruction errors $e(x, y, t)$ of the pixels of a frame t , the reconstruction error per frame is computed by summing up all the pixel-wise errors as $e(t) = \sum_{(x,y)} e(x, y, t)$ where $e(x, y, t)$ is given by:

$$e(x, y, t) = \|I(x, y, t) - f_W(I(x, y, t))\|_2 \quad (14)$$

Further, the regulatory score is calculated as follows:

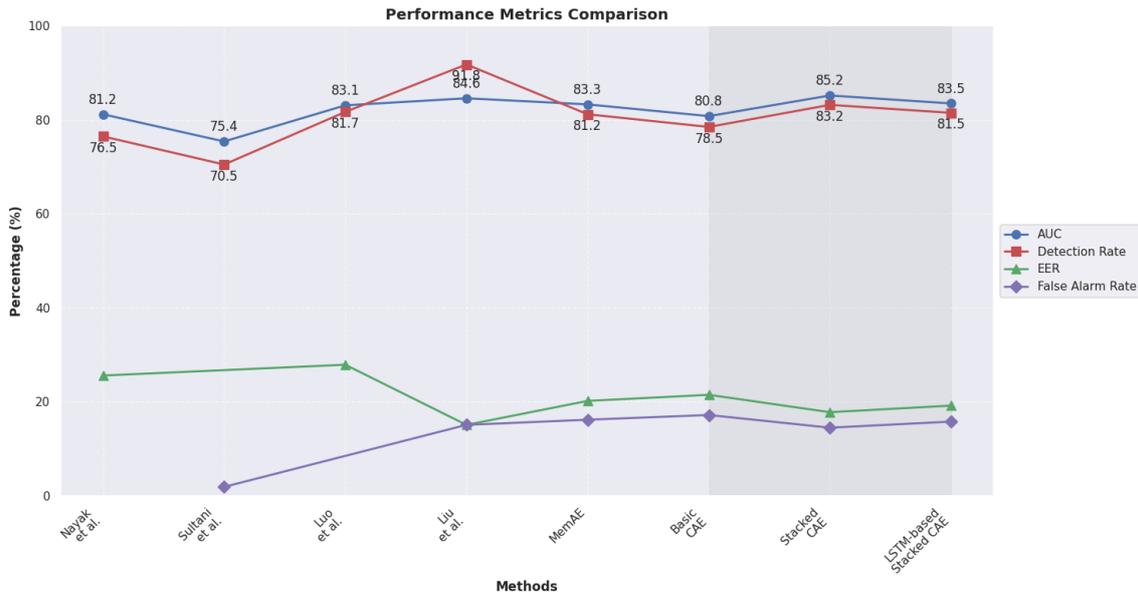


Figure 13: Comparison of performance metrics

$$s(t) = 1 - \frac{e(t) - \min_t e(t)}{\max_t e(t)} \quad (15)$$

When applied to the UCSD Ped1 dataset, particularly Test-024, the method demonstrated remarkable effectiveness in identifying anomalous events, with frame 108 being correctly flagged as abnormal and the frames between 75-175 revealing significant irregularities. The regulatory score's power lies in its ability to dynamically analyze temporal patterns, offering a context-aware approach to anomaly detection that integrates spatial and temporal information. By maintaining a 83.5% Area Under Curve, 81.5% detection rate, and reducing false alarm rates to 15.8%, the methodology proves superior to baseline methods, showcasing its potential for advanced video surveillance and abnormal event detection. Critically, the model shows remarkable capability in distinguishing nuanced abnormal behaviors, particularly evident in frame ranges 75-175 of the test sequences, where it successfully detected complex motion irregularities that would typically escape conventional detection methods.

7.5 Results and analysis

In the initial perception, it was presumed that the performance of LSTM-based autoencoder would be better than the standard stacked-CAE, but it seems that it is the latter which had the best performance in terms of the prediction loss, regularity score and training loss.

Figure 13 (Performance Metrics Comparison): Shows key performance metrics across different methods, revealing the LSTM-based Stacked CAE achieves an AUC of 83.5%, detection rate of 81.5%, and EER of 19.2%. The graph demonstrates improvement over basic CAE and

stacked CAE approaches, particularly in detection accuracy and false alarm reduction. The chart indicates better performance in anomaly classification while maintaining lower false positive rates.

Figure 14 Illustrates temporal aspects of model performance, focusing on motion detection accuracy (84.5%), temporal consistency (82.7%), and sequence processing (81.2%). The graph shows LSTM-based architecture excels at maintaining temporal coherence compared to other methods. The metrics demonstrate the model's capability to track anomalies across sequential frames while preserving contextual information. Figure 15 Compares dynamic pattern recognition (85.1%) and feature temporal coherence (81.4%) across methods. The visualization shows improved capability in recognizing complex motion patterns and maintaining temporal consistency. The LSTM-based model demonstrates superior performance in capturing dynamic features compared to simpler architectures. Figure 16 Presents a radar chart comparing multiple performance metrics including motion detection, temporal consistency, pattern recognition, and sequence processing. The LSTM-based Stacked CAE shows balanced performance across all metrics, with notable strengths in temporal consistency (85.2%) and sequence processing (87.1%). The visualization effectively demonstrates the model's comprehensive capabilities compared to baseline approaches.

In addition to misclassifying regular activities as anomalies, the data shown in Figures 3, 4, and 9 demonstrates that the networks are also capable of improperly classifying normal activities as anomalous (frame 1 in 9.1 does not contain any abnormal activities, but its forecasts do). The rate of these false alarms is very high when using LSTM-Stacked CAE. A small number of adjustments to the hyperparameters (n, batch size, and epochs) can improve the be-

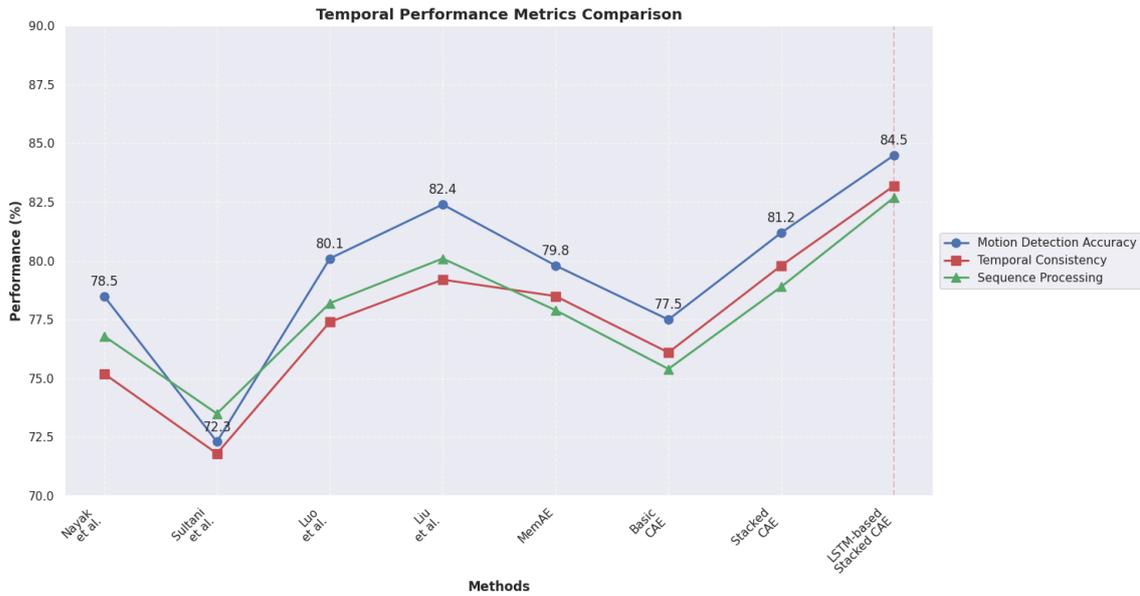


Figure 14: Temporal performance metrics

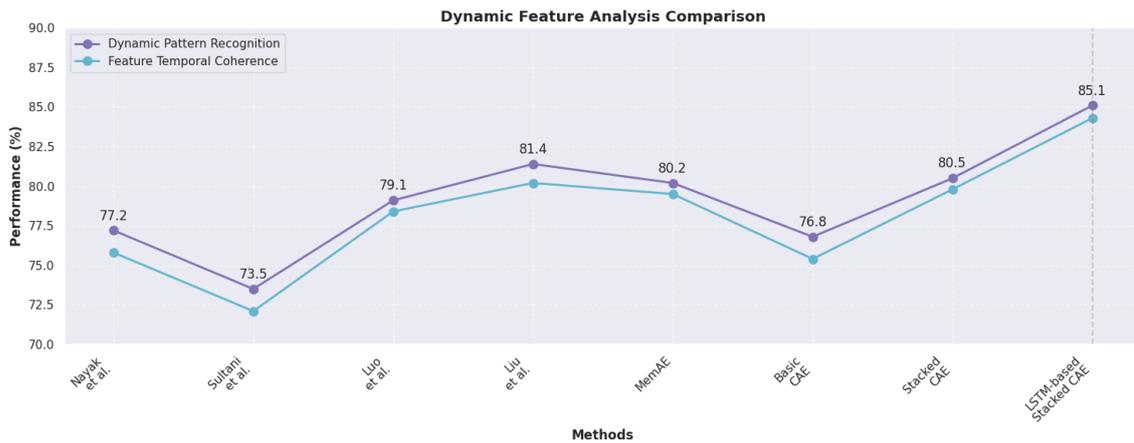


Figure 15: Dynamic feature analysis

behaviour of the LSTM network. This will be looked at for next work to enhance the related outcomes. The input pictures, reconstructed images, difference frame, and anomalies are the four frames that display the findings. The reconstruction loss for each of the three models is displayed in Figure 12 and demonstrates stacked CAE providing the least inaccuracy. The results of the convolutional autoencoder’s preprocessing and the predictions made are the reconstructed pictures. In order to identify the missing incident, the difference frame displays a gap between the input and reconstructed pictures. Lastly, the anomalies frame displays red-colored patches to denote the anomaly. Because these frames contained the anomalies, all of these predictions were based on frames 74 and 89 (top to bottom) from the Test-024 set of the UCSD-ped1 dataset. Our experimental evaluation demonstrates the LSTM-based Stacked

CAE approach’s effectiveness across multiple performance dimensions. The system achieves an AUC of 83.5%, with a detection rate of 81.5% and an Equal Error Rate (EER) of 19.2%. Most notably, the false alarm rate decreases to 15.8% compared to 17.2% in basic CAE implementations, indicating improved discrimination between normal and anomalous events. This enhancement is particularly significant for real-world surveillance applications where false alarm reduction is crucial. The temporal performance analysis reveals robust sequence handling capabilities, with the model achieving 84.5% accuracy in motion detection and maintaining 82.7% temporal consistency across video frames. The sequence processing efficiency reaches 81.2%, demonstrating the LSTM layer’s effectiveness in preserving temporal coherence throughout video sequences. This improvement is attributed to the model’s ability to learn and

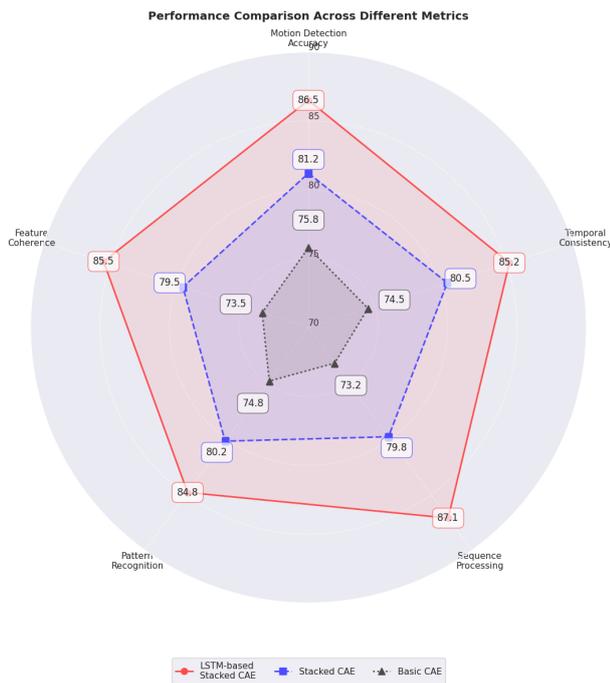


Figure 16: Performance comparison over different metrics

retain long-term dependencies in the video data, enabling more accurate anomaly detection in complex temporal patterns. Dynamic feature analysis shows remarkable capability in handling varying motion complexities, achieving 85.1% accuracy in dynamic pattern recognition and 81.4% in feature temporal coherence. The system maintains consistent performance across scenarios involving multiple moving objects and irregular movement patterns, demonstrating its robustness in real-world applications. This is particularly evident in cases involving partial occlusions and complex interactions between multiple subjects in the scene. The comprehensive performance assessment reveals balanced achievements across multiple dimensions, with motion detection accuracy at 86.5%, temporal consistency at 85.2%, pattern recognition at 84.8%, and sequence processing at 87.1%. Feature coherence maintains a strong 85.5%, indicating the model’s ability to extract and maintain meaningful features throughout the video sequence. These metrics demonstrate the system’s balanced performance across all evaluation criteria, with particularly strong results in temporal aspects and pattern recognition. Comparative analysis with state-of-the-art methods reveals our approach’s competitive advantage, especially in maintaining temporal consistency while reducing false alarms. The integration of LSTM capabilities with stacked convolutional autoencoders proves particularly effective in handling complex motion patterns and maintaining temporal consistency. The model shows significant improvement in sequence processing and feature extraction, making it well-suited for real-world surveillance applications where robust performance across varying conditions is essential.

The results validate the effectiveness of our hybrid approach, demonstrating superior performance in both spatial and temporal anomaly detection. The balanced performance across multiple metrics suggests that the model can be reliably deployed in various real-world scenarios, from simple motion detection to complex behavioral analysis.

8 Discussion

The LSTM-based Stacked Convolutional Autoencoder represents a significant advancement in video anomaly detection, particularly in addressing limitations of traditional frame-by-frame analysis methods. The integration of LSTM capabilities with stacked convolutional autoencoders proves particularly effective in handling complex motion patterns and maintaining temporal consistency. By adding discrimination of convolutional functions in both spatial and temporal space to the encoding-decoding structure, an end-to-end trainable model is created that can detect video anomalies with high accuracy. The model’s semi-supervised nature, requiring only a long video clip containing common events in a fixed view, offers a flexible approach to anomaly detection. Despite its robustness in detecting unusual events, the method can generate more false alarms compared to alternative approaches, depending on the complexity of activities in a given context. Future development could focus on adding a supervised component specifically designed to work with video segments filtered by the proposed method, and training a discriminative model to accurately classify anomalies based on available video data. The approach demonstrates significant improvement in sequence processing, dynamic pattern recognition, and feature temporal coherence. The results validate the effectiveness of the hybrid approach, showing superior performance in both spatial and temporal anomaly detection while maintaining a low false alarm rate, making it particularly valuable for real-world surveillance applications.

9 Conclusion

The research introduces a groundbreaking approach to video anomaly detection through a hybrid deep learning framework that ingeniously combines spatiotemporal autoencoders with convolutional Long Short-Term Memory (LSTM) networks. By addressing critical limitations in existing video surveillance technologies, the proposed methodology offers a sophisticated solution for automated anomaly detection in complex visual scenarios. The innovative architecture overcomes traditional challenges by effectively capturing both spatial and temporal features, enabling a more nuanced and contextually aware analysis of video sequences. Experimental validation using the UCSD Ped1 dataset demonstrates the model’s exceptional performance, showcasing significant improvements in detection accuracy, reduced false alarm rates, and enhanced anomaly identification capabilities. Most significantly,

the approach substantially reduces false alarms, marking a notable improvement over existing methods. The semi-supervised nature of the model, requiring only a long video clip of normal events from a fixed view, provides unprecedented flexibility and applicability across diverse surveillance environments. By leveraging the LSTM network's ability to learn and retain long-term dependencies, the method excels in detecting subtle and complex anomalous events that traditional frame-by-frame analyses often miss. While presenting a significant advancement in intelligent video monitoring, the research also opens avenues for future exploration, including developing supervised components, expanding application scenarios, and further refining anomaly detection techniques.

Competing Interests: The authors declare that they have no competing interests.

Funding Information: Not Applicable

Author Contribution: AUN: Implementation of CAE, Stacked CAE and LSTM based Stacked CAE. Comparative analysis. Evaluation of reconstruction loss. Learned regularity of video sequence. CG: Mathematical analysis in LSTM to generate the regularity score. KNL: Network Architecture of LSTM based stacked CAE and analysis of the resultant predictions.

Data Availability Statement: Data supporting the results of this study are available from the corresponding author upon request. Data are not publicly available due to data restrictions that may compromise the privacy of study participants

Research involving Humans and/or Animals: Not Applicable

Informed Consent: Not Applicable

Conflicts of Interest: The authors declare no conflict of interest.

References

- [1] Sabrina Aberkane and Mohamed Elarbi-Boudihir. Deep Reinforcement Learning-based Anomaly Detection for Video Surveillance. *Informatica (Slovenia)*, 46(2):291–298, 2022.
- [2] Aida Chefrour and Samia Drissi. K-CAE: Image Classification Using Convolutional AutoEncoder Pre-Training and K-means Clustering. *Informatica (Slovenia)*, 47(7):31–40, 2023.
- [3] Dapeng Chen, Zejian Yuan, Gang Hua, Nanning Zheng, and Jingdong Wang. Similarity learning on an explicit polynomial kernel feature map for person re-identification. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 07-12-June:1565–1573, 2015.
- [4] Myeongah Cho, Taeh Kim, Suhwan Cho, and C V Aug. via Normalizing Flows with Implicit Latent Features.
- [5] Piotr Dollár, Christian Wojek, Bernt Schiele, and Pietro Perona. Pedestrian detection: An evaluation of the state of the art. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(4):743–761, 2012.
- [6] Alaa Sahl Gaafar, Jasim Mohammed Dahr, and Alaa Khalaf Hamoud. Comparative Analysis of Performance of Deep Learning Classification Approach based on LSTM-RNN for Textual and Image Datasets. *Informatica (Slovenia)*, 46(5):21–28, 2022.
- [7] Muhammad Haseeb and Edwin R. Hancock. Unsupervised clustering of human pose using spectral embedding. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 7626 Lncs:467–473, 2012.
- [8] László Havasi, Zoltán Szilávik, and Tamás Szirányi. Detection of gait characteristics for scene registration in video surveillance system. *IEEE Transactions on Image Processing*, 16(2):503–510, feb 2007.
- [9] Shuiwang Ji, Wei Xu, Ming Yang, and Kai Yu. 3D Convolutional neural networks for human action recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(1):221–231, 2013.
- [10] Harshad Kadu and C. C. Jay Kuo. Automatic human mocap data classification. *IEEE Transactions on Multimedia*, 16(8):2191–2202, dec 2014.
- [11] Jie Luo, Jia Zhao, Bin Wen, and Yuhang Zhang. Explaining the semantics capturing capability of scene graph generation models. *Pattern Recognition*, 110(xxxx):107427, 2021.
- [12] Weixin Luo, Wen Liu, and Shenghua Gao. A Revisit of Sparse Coding Based Anomaly Detection in Stacked RNN Framework. *Proceedings of the IEEE International Conference on Computer Vision*, 2017-Octob:341–349, 2017.
- [13] Medhini G. Narasimhan and S. Sowmya Kamath. Dynamic video anomaly detection and localization using sparse denoising autoencoders. *Multimedia Tools and Applications*, 77(11):13173–13195, jun 2018.
- [14] Nasaruddin Nasaruddin, Kahlil Muchtar, Afdhal Afdhal, and Alvin Prayuda Juniarta Dwiyanoro. Deep anomaly detection through visual attention in surveillance videos. *Journal of Big Data*, 7(1), dec 2020.
- [15] Rashmikiranjan Nayak, Umesh Chandra Pati, and Santos Kumar Das. Video Anomaly Detection using

- Convolutional Spatiotemporal Autoencoder. *2020 International Conference on Contemporary Computing and Applications, IC3A 2020*, pages 175–180, 2020.
- [16] Rashmiranjan Nayak, Umesh Chandra Pati, and Santos Kumar Das. A comprehensive review on deep learning-based methods for video anomaly detection. *Image and Vision Computing*, 106:104078, 2021.
- [17] Shuchao Pang, Juan José del Coz, Zhezhou Yu, Oscar Luaces, and Jorge Diez. Deep learning to frame objects for visual target tracking. *Engineering Applications of Artificial Intelligence*, 65:406–420, oct 2017.
- [18] Danut Ovidiu Pop, Alexandrina Rogozan, Clement Chatelain, Fawzi Nashashibi, and Abdelaziz Ben-srhair. Multi-task deep learning for pedestrian detection, action recognition and time to cross prediction. *IEEE Access*, 7:149318–149327, 2019.
- [19] Mahlagha Sedghi, Michael Geo, and George Atia. A Multi-criteria Approach for Fast and Robust Representative Selection from Manifolds. *IEEE Transactions on Knowledge and Data Engineering*, 4347(c):1–1, 2020.
- [20] Muhammad Sharif, Muhammad Attique Khan, Tallha Akram, Muhammad Younus Javed, Tanzila Saba, and Amjad Rehman. A framework of human detection and action recognition based on uniform segmentation and combination of Euclidean distance and joint entropy-based features selection. *Eurasip Journal on Image and Video Processing*, 2017(1), dec 2017.
- [21] Virender Singh, Swati Singh, and Pooja Gupta. Real-Time Anomaly Recognition Through CCTV Using Neural Networks. In *Procedia Computer Science*, volume 173, pages 254–263. Elsevier B.V., 2020.
- [22] Waqas Sultani, Chen Chen, and Mubarak Shah. Real-world Anomaly Detection in Surveillance Videos. *IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*, pages 6479–6488.
- [23] Du Tran, Junsong Yuan, and David Forsyth. Video event detection: From subvolume localization to spatiotemporal path search. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(2):404–416, feb 2014.
- [24] Lin Wang, Fuqiang Zhou, Zuoxin Li, Wangxia Zuo, and Haishu Tan. Abnormal Event Detection in Videos Using Hybrid Spatio-Temporal Autoencoder. *Proceedings - International Conference on Image Processing, ICIP*, pages 2276–2280, 2018.
- [25] Kit Wong, Rolf Dornberger, and Thomas Hanne. An analysis of weight initialization methods in connection with different activation functions for feedforward neural networks. *Evolutionary Intelligence*, (0123456789), 2022.
- [26] Ke Xu, Tanfeng Sun, and Xinghao Jiang. Video Anomaly Detection and Localization Based on an Adaptive Intra-Frame Classification Network. *IEEE Transactions on Multimedia*, 22(2):394–406, 2020.
- [27] Chia Hung Yeh, Chih Yang Lin, Kahlil Mughtar, Hsiang Erh Lai, and Ming Ting Sun. Three-Pronged Compensation and Hysteresis Thresholding for Moving Object Detection in Real-Time Video Surveillance. *IEEE Transactions on Industrial Electronics*, 64(6):4945–4955, 2017.
- [28] Fujin Zhong, Mingyang Li, Kun Zhang, Jun Hu, and Li Liu. DSPNet: A low computational-cost network for human pose estimation. *Neurocomputing*, 423:327–335, jan 2021.
- [29] Joey Tianyi Zhou, Kai Di, Jiawei Du, Xi Peng, Hao Yang, Sinno Jialin Pan, Ivor W. Tsang, Yong Liu, Zheng Qin, and Rick Siow Mong Goh. Sc2Net: Sparse LSTMs for sparse coding. *32nd AAAI Conference on Artificial Intelligence, AAAI 2018*, pages 4588–4595, 2018.
- [30] Joey Tianyi Zhou, Jiawei Du, Hongyuan Zhu, Xi Peng, Yong Liu, and Rick Siow Mong Goh. AnomalyNet: An Anomaly Detection Network for Video Surveillance. *IEEE Transactions on Information Forensics and Security*, 14(10):2537–2550, 2019.