

AI-Driven Salient Soccer Events Recognition Framework for Next Generation IoT-Enabled Environments

Khan Muhammad, *Member, IEEE*, Hayat Ullah, Mohammad S. Obaidat, *Life Fellow, IEEE*, Amin Ullah, *Student Member, IEEE*, Arslan Munir, *Senior Member, IEEE*, Muhammad Sajjad, Victor Hugo C. de Albuquerque, *Senior Member, IEEE*

Abstract— The salient events recognition of soccer matches in next-generation Internet of things (Nx-IoT) environment aims to analyze the performance of players/teams by the sports analytics and managerial staff. The embedded Nx-IoT devices carried by the soccer players during the match capture and transmit data to an Artificial Intelligence (AI)-assisted computing platform. The interconnectivity of data acquisition devices with an AI-assisted computing platform in the Nx-IoT environment will not only allow the spectators to track the formation of their favorite players during a soccer match but will also enable the managerial staff to evaluate the players' performance in the soccer match as well as in practice sessions. This Nx-IoT-enabled salient event detection feature can be provided to spectators and sports' managerial staff as a financial technology (FinTech) service. In this paper, we propose an efficient deep learning-based framework for multi-person salient soccer events recognition in IoT-enabled FinTech. The proposed framework performs event recognition in three steps: Firstly, image frames are extracted from video streams and resized in the preprocessing step to match the input of the deep network. Secondly, frame-level discriminative features are extracted using a pre-trained convolutional neural network (CNN) architecture. Thirdly, we employ a multi-layer long short-term memory (MLSTM) network to recognize high-level events in soccer videos by exploiting the sequential relation between adjacent frames. Moreover, we introduce a new soccer video events (SVE) dataset containing videos of six salient events of soccer game. To provide a strong baseline, we evaluate our newly created SVE dataset using different traditional machine learning and deep learning algorithms. We also perform event recognition on untrimmed soccer videos using our proposed framework and compare the results with state-of-the-art methods. The obtained results validate the suitability of our proposed framework for salient events recognition in Nx-IoT environments.

Index Terms— Event recognition, CNN, MLSTM, edge computing, next-generation Internet of Things (Nx-IoT), computer vision.

I. INTRODUCTION

The noticeable advancement in networking technologies and efficient deep learning algorithms over edge devices has allured the sports industry to adopt next generation Internet of things (Nx-IoT) enabled financial technology (FinTech) services for a wide range of applications. The IoT is a network where multiple edge devices/sensors are interconnected via the Internet. The Nx-IoT devices in sports communicate and transmit data with other IoT/edge devices for edge-centric distribution and processing of sports data. Mostly sports organizations, especially soccer officials, provide edge based IoT environments, which can significantly improve the quality of sport analytics systems, and enable the spectators to have a more enjoyable interactive experience. The video data captured by the vision sensors can be instantly processed over edge devices, and then transmitted to AI-assisted edge computing platforms for a variety of applications such as events detection/recognition, player identification, and players

Manuscript received February 25, 2021; Revised June 16, 2021, and July 18, 2021; Accepted August 19, 2021; Published: XXXX. This research is supported by: (1) the tenure of an ERCIM 'Alain Benoussan' Fellowship Programme under the Contract 2019-40 and (2) Color and Visual Computing Lab, Department of Computer Science, NTNU, Gjøvik. The work of Mohammad S. Obaidat was supported in part by PR of China Ministry of Education Distinguished Possessor Grant given to him under number: MS2017BJKJ003. (*Corresponding authors: Khan Muhammad and Muhammad Sajjad*)

Khan Muhammad and Hayat Ullah (co-first authors) are with the Department of Software, Sejong University, Seoul 143-747, Republic of Korea. (Email: khan.muhammad@ieee.org, khanh9474@gmail.com)

Mohammad S. Obaidat, Fellow of IEEE and Fellow of SCS, Founding Dean and Professor, is with the College of Computing and Informatics, University of Sharjah, Sharjah 27272, UAE, with King Abdullah II School of Information Technology, University of Jordan, Amman 11942, Jordan, and with School of Computer and Communication Engineering, University of Science and Technology Beijing 100083, China. (Email: m.s.obaidat@ieee.org)

Amin Ullah is with the CORIS Institute, Oregon State University, Corvallis, OR 97331, USA (Email: ullaham@oregonstate.edu)

Arslan Munir is with the Department of Computer Science, Kansas State University, Manhattan, Kansas, 66506, USA. (Email: amunir@ksu.edu)

Muhammad Sajjad is with the Digital Image Processing Laboratory, Islamia College Peshawar, Peshawar 25000, Pakistan, and also with the Norwegian Colour and Visual Computing Laboratory, Department of Computer Science, Norwegian University of Science and Technology (NTNU), Gjøvik, 2815, Norway. (Email: muhammad.sajjad@icp.edu.pk)

Victor Hugo C. de Albuquerque is with the Department of Teleinformatics Engineering, Federal University of Ceara, Fortaleza, Fortaleza/CE, Brazil. (Email: victor.albuquerque@ieee.org)

formation tracking in the field of soccer stadium equipped with Nx-IoT. Due to the exponential growth of fans following, soccer has become the world's most watchable sport with more than 4.0 billion audience worldwide [1]. According to a recent report (Google: Watch time for YouTube sports highlights jumps 80%), 90% of online viewers search for soccer videos highlights or prefers to access salient sports events (such as goal, penalties, fouls, and corner-shots, etc.) rather than watching full matches [2]. Furthermore, the live spectators inside the stadium are very excited to support their favorite team/players and cheer for their best performance. The Nx-IoT-enabled edge-based events recognition service in soccer stadiums can improve the experiences of live spectators by providing an interactive entertainment environment.

Considering the complex game rules and players with different field formations, soccer is the most difficult game to analyze. Researchers around the globe are contributing to different aspects of soccer events detection and recognition. Event recognition is an essential component for high-level sports video analytics tasks, such as event-aware highlights generation [3], sports videos retrieval [4], and indexing of sports videos [5]. However, soccer events are different from other sports events, where a video clip contains fascinating contents for a random time interval with semantically starting and ending boundaries of events rather than a fixed time interval. For instance, in counterattack, the brilliant assist before the goal and the celebrations after the goal, are the

complete soccer events. All these events are high-level semantics, which can be recognized with multi-scale deep features (i.e., extract CNN features at different layer with varying spatial dimensionality). For instance, in soccer match, a goal is an event which involves different movements of human body, such as running, jumping, passing ball, and shooting ball towards goalmouth.

The earlier research studies [6, 7] for soccer events detection and recognition were based on low-level or handcrafted visual features and traditional machine learning algorithms. These low-level feature-based methods rely on global features including texture, edge, color, shape, and motion information. Although these methods made some achievements in the last decade, they could only detect few soccer events and usually failed for complex type of events with clutter background. To interrelate the semantic gap between low-level semantics and high-level semantics, several traditional soccer events detection methods have been proposed. These methods utilized mid-level features to obtain the intermediate representation of soccer events, including field view classification, player tracking, scoreboard detection, and play-break segmentation. For instance, Zhao et al. [8] used mid-level features for video segmentation into play-break segments. They segmented the video based on color, contour, and histogram. The works presented in [3, 9, 10] first extracted the excitement clip from lengthy videos, and then detected the salient events using histogram and color computations. In addition, these earlier methods often required additional information such as text from score boards and audio commentary related to the game play. Despite, acquiring additional information about game, the results achieved by these methods still suffered from misclassification for complex events.

Recently, deep learning has gained tremendous attention in computer vision, and has significantly improved the performance of event detection and action recognition systems. Various CNN-assisted approaches have been proposed for soccer events detection and annotation [11, 12], which extract both the spatial and temporal features from the video frames and analyze the event's type and boundary for specific time interval. While, some researchers have adopted 3D-CNN based soccer events detection approaches that extract spatial and temporal features from video frames [13, 14], other studies [15, 16] have presented the combined CNN and recurrent neural network (RNN) frameworks for soccer event detection and achieved state-of-the-art results. However, most of the contemporary deep learning methods are limited to certain types of events and single-person events detection and cannot be deployed in IoT-enabled environments. Considering the availability of embedded edge devices and efficient deep learning architectures, there is a need to develop an efficient edge computing-based approach for salient soccer events recognition in Nx-IoT-enabled environments. Besides the existence of smart embedded devices and energy-friendly deep learning architectures, it is very crucial to have sufficient amount of data for the problem/task under consideration. The availability of problem-related datasets greatly helps the researchers to train and evaluate their proposed systems without devoting considerable efforts on generating new datasets. However, to the best of our knowledge, there are very few

soccer videos datasets [17-19], and further the existing datasets are very specific and do not cover salient events of soccer.

Therefore, in this paper we present an efficient deep learning-based framework for salient soccer events recognition over edge-centric FinTech computing platform and our newly created soccer videos events (SVE) dataset. The proposed framework performs event recognition process in three steps: (i) preprocessing, (ii) features extraction, and (iii) sequence learning. In the preprocessing step, image frames are extracted from video streams and resized to match the input of the deep network. For feature extraction, our framework uses a pretrained CNN architecture which extracts deep discriminative features from the video frames. While for sequence learning, a multi-layer LSTM is used to analyze the video stream by capturing the temporal relation between adjacent frames. Our newly created SVE dataset contains short duration clips of six different soccer events. To better understand the problem, we first evaluate the soccer event recognition with handcrafted features and a well-known machine learning classifier (HOG+SVM). Later, we investigate the soccer event recognition problem using MLSTM along with different CNN on our SVE dataset. The key contributions of our scheme can be summarized as follows:

- 1) To recognize the salient events in soccer matches, we investigated traditional machine learning (HOG+SVM) and deep learning-based approaches (CNN+MLSTM) for FinTech-enabled soccer events recognition service and propose an energy-efficient CNN+LSTM framework. Our proposed framework strikes a tradeoff between computational complexity and model accuracy and is a suitable solution for edge-centric FinTech computing platforms and similar domains associated with Nx-IoT environments, showing its flexibility and scalability.
- 2) The literature contains very few datasets for soccer events detection/recognition. However, there is no benchmark dataset of key events, which defines the interest of live/offline spectators. We have created our own SVE dataset, which contains salient events of soccer matches captured from multiple views. The SVE dataset will be publicly available for further research to mature the event detection/recognition systems for soccer videos.
- 3) We have conducted comprehensive experiments on our newly created SVE dataset to evaluate the performance of our framework. Further, we have tested the proposed framework on relevant events from other datasets and have conducted a comparative study. The obtained results reveal that the proposed framework generalizes well and performs better than existing methods.

The rest of the paper is organized as follows: Section II presents the overview of the related works. The proposed framework is presented in Section III followed by experimental evaluation of the proposed framework in Section IV. Finally, Section V concludes the paper with possible future directions.

II. RELATED WORK

In this section, we briefly describe the event recognition literature and critically discuss the soccer event recognition

approaches that are reported in the recent works along with their strengths and limitations. Generally, the soccer event recognition methods can be categorized into two parts low-level features based and deep learning-based methods.

2.1 Low-level features-based event recognition approaches

Event recognition has played a very important role in different domains of sports video contents analysis including highlight generation, event-based sports video retrieval, and statistical summary generation of sports videos (e.g., soccer, hockey, etc.). A variety of methods have been proposed to automate the event recognition process in sports videos. Most of the early methods [6, 7, 20] were based on low-level features. These methods usually used hand-crafted descriptors and machine learning classifiers for feature extraction and classification, respectively. For instance, Tavassolipour et al. [21] proposed automatic event detection in soccer videos for highlights generation. They used hidden Markov model for the segmentation of video into meaningful segments named as play-back patterns followed by mid-level features extraction from each segment. Finally, they extracted discriminative features using a Bayesian network. Kolekar and Sengupta [22] proposed an automatic highlight generation system that could generate highlight from sports TV broadcasts. First, they detected the exciting clips using audio features and then segmented the clips into different scenes. Next, they assigned concept-score to each scene within a clip using a probabilistic Bayesian belief network (PBBN) and selected the scene with higher concept-score. Wang et al. [23] proposed a soccer video annotation framework based on coarse-grained time information. They annotated the soccer events by synchronizing the video clips and external text information (match reports) with coarse time constraints. Fakhar et al. [24] presented a learning-based soccer event detection approach based on two main concepts. First, they analyzed the frame and estimated the saliency of each frame regarding soccer events. Second, the event-oriented and discriminative dictionary was learned using their proposed K-SVD algorithm. Furthermore, Babbitt et al. [25] proposed a technique for video-based talent identification of the youth for soccer using smart devices. Hosseini and Moghadam [26] presented a fuzzy rule-based system for soccer event detection and annotation. They used statistical information quantized from audiovisual features and rule-based reasoning classifier, which constructed the semantic perception for the occurred events. However, these hand-crafted or low-features based methods are less effective and time consuming for detecting high-level soccer events. These limitations can create issues when processing lengthy videos or dealing with sports TV broadcasts. Besides these traditional hand-crafted-based methods, numerous learning-based event detection/recognition methods have been proposed which significantly improve the event detection and recognition task and overcome the limitations of traditional methods.

2.2 Deep learning-based event recognition approaches

Recently, CNN-oriented methods have achieved greater success and have improved the performance of various computer vision task including image classification [27, 28], object detection [29, 30], image enhancement [31, 32], speech

recognition [33, 34], and activity recognition [35, 36]. For instance, Jiang et al. [15] proposed a deep learning-based approach for soccer video event detection. Their method utilized the combination of CNN and RNN, where they segmented the soccer video in play-break segments by determining the event boundary, and then extracted CNN features of key frames from play-break segments. Finally, RNN was deployed to recognize the salient soccer events, including goal, goal attempt, card, and corner. Tsunoda et al. [37] presented a hierarchical RNN for analyzing the understanding between players of team sports activity. They integrated multiple person-centered features with LSTM cell output over temporal sequences. Further, Fani et al. [38] introduced a parallel feature fusion (PFF) network for automatic event detection and classification in soccer broadcast videos. The PFF combined local as well as full scene features for zoom in and zoom out scene classification. Hidden observable Markov model was deployed to determine play/break status of the scenes in soccer videos. Giancola et al. [39] introduced a benchmark SoccerNet dataset for action spotting in soccer videos. The duration of the dataset was 764 hours and consisted of goal, yellow/red card and substitution. To prevent violence incidents in football stadiums, Samuel et al. [40] proposed a real-time violence detection system for recognizing violence using human intelligence simulation. Their proposed method processed enormous amounts of real-time video streams from different sources, where Histogram of Oriented Gradients (HOG) was used as a feature descriptor followed by Bidirectional Long Short-Term Memory (BDLSTM). Liu et al. [41] proposed a soccer event detection method based on temporal action localization and play-break segmentation. First, they localized the action in soccer videos using 3D CNN and then employed play-break rules for organizing actions into corresponding events. These deep learning-based approaches have shown remarkable performance and have effectively overcome the limitations of low-level features-based methods. On the other hand, these deep learning-based approaches require high computation power for training purposes. Different from the existing methods, our proposed framework efficiently reduces the computational complexity by adopting transfer learning and frame skip strategies.

III PROPOSED FRAMEWORK

Human action-oriented events typically involve sequences of specific human postures evolving in video frames, which demonstrate variations in both spatial and temporal domains. For example, the goal event in soccer consists of more than one action where each action is the combination of different postures of human body. While analyzing the soccer video, these high-level actions can be visualized as hidden sequential patterns, which can be detected and recognized with strong representation (high-level features). In this paper, we propose a deep learning-based salient event recognition framework, which analyzes the input soccer video using CNN with MLSTM. The proposed framework comprises of three steps: **Step 1 – Preprocessing:** The preprocessing step extracts the image frames from the video stream and resize the frames to match the input of the deep neural network.

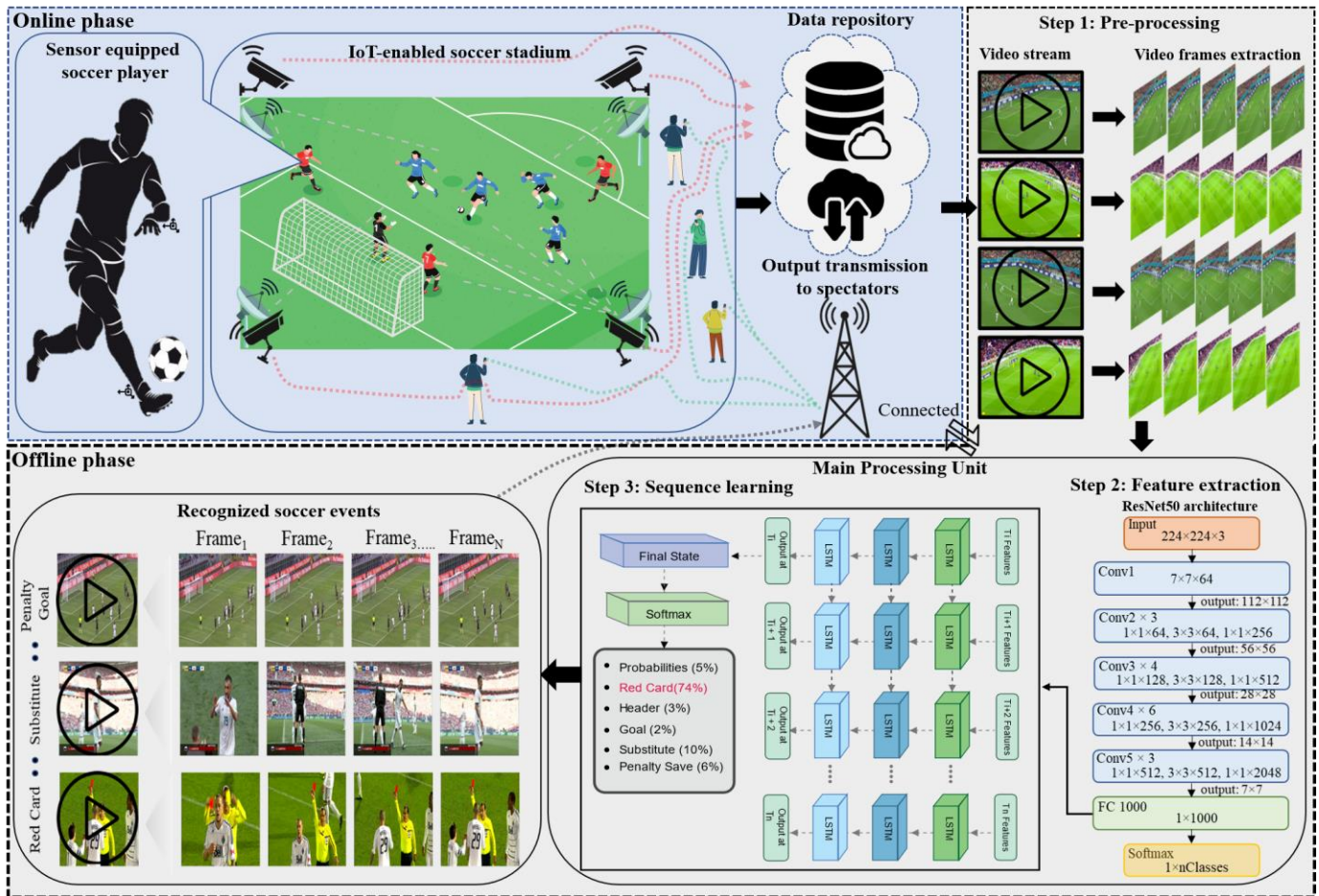


Fig. 1. The detailed overview of our proposed framework for salient soccer events recognition over edge devices in Nx-IoT-enabled environment. **Step 1:** This step involves frames extraction and frame resizing. **Step 2:** This step involves features extraction, where we employ a pretrained CNN (ResNet50) architecture to extract discriminative features from video frames. **Step 3:** This step receives CNN features and trains multi-layer LSTM, which outputs a confidence score for an event detected in sequences of video frames.

Step 2 – Feature Extraction: In the feature extraction step, our framework extracts deep features using a pretrained CNN network from the sequence of frames.

Step 3 – Sequence Learning: In the third step, the extracted features are fed into an MLSTM to retrieve high-level abstraction and temporal information for the event recognition task. The workflow of the proposed framework and its main components are illustrated in Fig. 1. Algorithm 1 presents the stepwise implementation of soccer event recognition process. The input and output parameters used in our proposed framework are presented in Table 1.

3.1 Preprocessing

In the preprocessing step, frames are extracted from the video with three different frame skip strategies and then resized each frame to match the input of the deep neural network. In our proposed framework, the input is resized to $224 \times 224 \times 3$ to match the input layer dimensionality of ResNet50 architecture.

3.2 Features Extraction ResNet50 Architecture

Soccer event is the combination of multiple actions (running, jumping, and passing, etc.), where each action is itself the integration of different poses. The event recognition workflow starts from low-level semantics extraction (actions) to high-level semantics (events). To represent these sequences of

TABLE I
DESCRIPTION OF PARAMETERS USED IN OUR PROPOSED FRAMEWORK FOR INPUT AND OUTPUT OPERATION

Symbols	Description
fc-1000	Fully connected layer of ResNet50 architecture.
f_s	Number of frame skip during feature extraction.
x_t	Input to LSTM at time t .
f_t	Output of forget gate.
i_t	Output of input gate.
o_t	Output of output gate.
c_t	Output of current state of LSTM cell.
c_{t-1}	Previous state of LSTM cell.
w_f	Weights of forget gate of LSTM cell.
w_i	Weights of input gate pf LSTM cell.
w_o	Weights of output gate of LSTM cell.
b_f	Biases of forget gate.
b_i	Biases of input gate.
b_o	Biases of output gate.
h_t	Final output of LSTM cell.

actions, CNN features from each frame are expressed as an individual CNN feature vector. Similarly, the complete video can be represented as a set of feature vectors having sequential relation

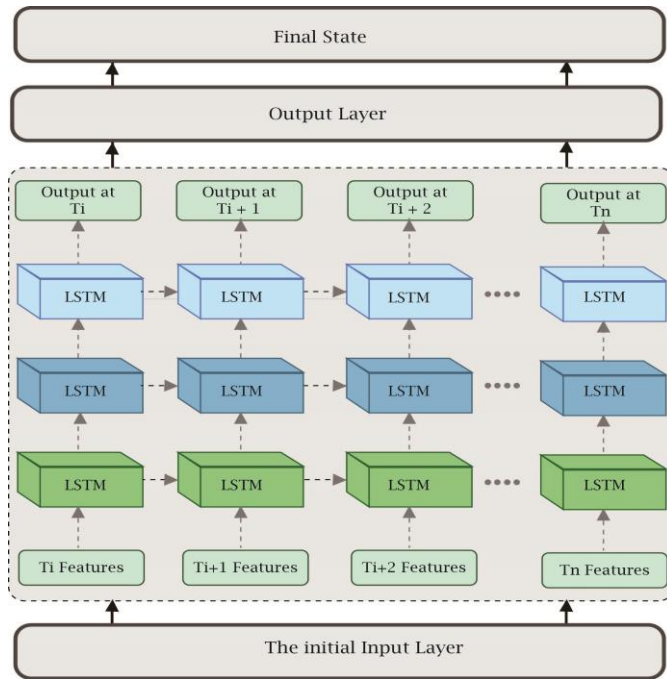


Fig. 2. The external structure of a multi-layer LSTM.

between the adjacent feature vectors. On the other hand, traditional machine learning algorithms use hand-crafted features for event detection and recognition tasks, which require a lot of efforts for feature engineering and scaling. Despite the extra-ordinary efforts for feature engineering, traditional soccer events recognition approaches are still unable to detect complex and long duration events.

CNN is originally introduced for image classification task [42] and has achieved state-of-the-art results. It has the ability to extract features of different scales and is equipped with a classifier at the end of architecture. CNNs are widely used for a variety of high-level computer vision tasks. The main reason behind the success and achievements of CNNs is the hierarchical nature of architecture that contains series of layers, including convolution, pooling, and fully connected layers. The convolutional layer generates different representation of the same image by convolving different kernels with different sizes. The pooling layer sub-samples the input feature maps by selecting the high activations values, while fully connected layer learns high-level representations and reshapes the input feature maps to a one-dimensional feature vector. Training a new CNN architecture from the scratch requires a huge amount of image data along with powerful machines for execution such as GPUs (Graphics Processing Units) or TPUs (Tensor Processing Units). This problem can be addressed using transfer learning strategy where a pretrained model is utilized for another computer vision problem. To this end, our proposed framework uses a pre-trained ResNet50 architecture [43], which is trained on the ImageNet dataset containing more than 20,000 categories (door, chair, and car, etc.). The first layer of the ResNet50 architecture is input layer with dimensionality of $224 \times 224 \times 3$, the second layer is convolution layer with kernel size 7×7 . Rest of the architecture has four residual blocks, fully connected layers, and a softmax layer. Each residual block contains three convolutional layers with kernel sizes of 1×1 and 3×3 followed by ReLU activation and Batch Normalization.

Algorithm 1 Event Recognition in Soccer Videos

Input: Soccer video V_{soccer}

Preparation:

1. Load pretrained ResNet50 CNN network M_f
2. Load trained multi-layer LSTM network M_c

Steps:

while (V_{soccer})

1. Read frames $\leftarrow (f_i, V_{\text{soccer}})$
2. Pass frame f_i to ResNet50 CNN
3. Extract feature $fv_i \leftarrow M_f(f_i)$
4. Forward feature vector fv_i to trained LSTM M_c , and Predict event class $\leftarrow M_c(fv_i)$
5. Display predicted event with confidence score

end while

Output: Display event with predicted label and confidence score

layers. We have used a fully connected layer (fc-1000) as a generic feature descriptor. Each feature vector represented a single frame of video, these features are then fed into MLSTM in the form of features block for a fixed time interval. MLSTM processes these features and learns the hidden sequential patterns from the input feature data. The detailed explanation of RNN and its variants are presented in the next section.

3.3 Event-Specific Sequence Learning using Multi-layer LSTM

Despite the powerful characteristics and flexibility, CNNs can only be used for tasks where inputs and outputs have fixed dimensionality and mostly fail while dealing with the data having different input and output dimensionality. Along with this limitation, CNNs are restricted to static data and cannot be used for problems dealing with time series and sequential data. Most of the problems such as speech recognition, machine translation, and activity recognition in videos are efficiently expressed with sequences having variable lengths. To solve sequential pattern learning problems or predicting time-series data, the need of such a method becomes crucial that can precisely map sequences and learn its hidden patterns from input time-series data. To meet the needs of such kind of systems, a special kind of neural network named RNN has been introduced which has the ability to learn from temporal features and map the temporal relation of a given time-distributed data. RNNs are specially designed for classification of time-series and sequential data. RNNs analyze the hidden sequential patterns in both spatial and temporal dimension by connecting the previous information with the current information and predict the future output. The suitability and efficiency of RNN for temporal data analytics has attracted the research community to investigate it for various time-series prediction and sequence classification problems and achieved incredible results. Although, RNNs can decipher the hidden sequential patterns in time-series data (i.e., video, audio, or numerical data), RNNs fail to remember earlier information while interpreting long term sequences. Such type of problem is known as a vanishing gradient or gradient exploding, which can be resolve by using a special variant of RNN known as LSTM, which has the ability to remember the earlier input information for a long-time interval.

3.3.1 Multi-layer LSTM Network

The LSTM [44] is an extension of the RNN architecture, which is specially designed for interpreting long-term sequences, thereby resolving the problem of vanishing gradient and gradient exploding faced by RNNs. The internal structure

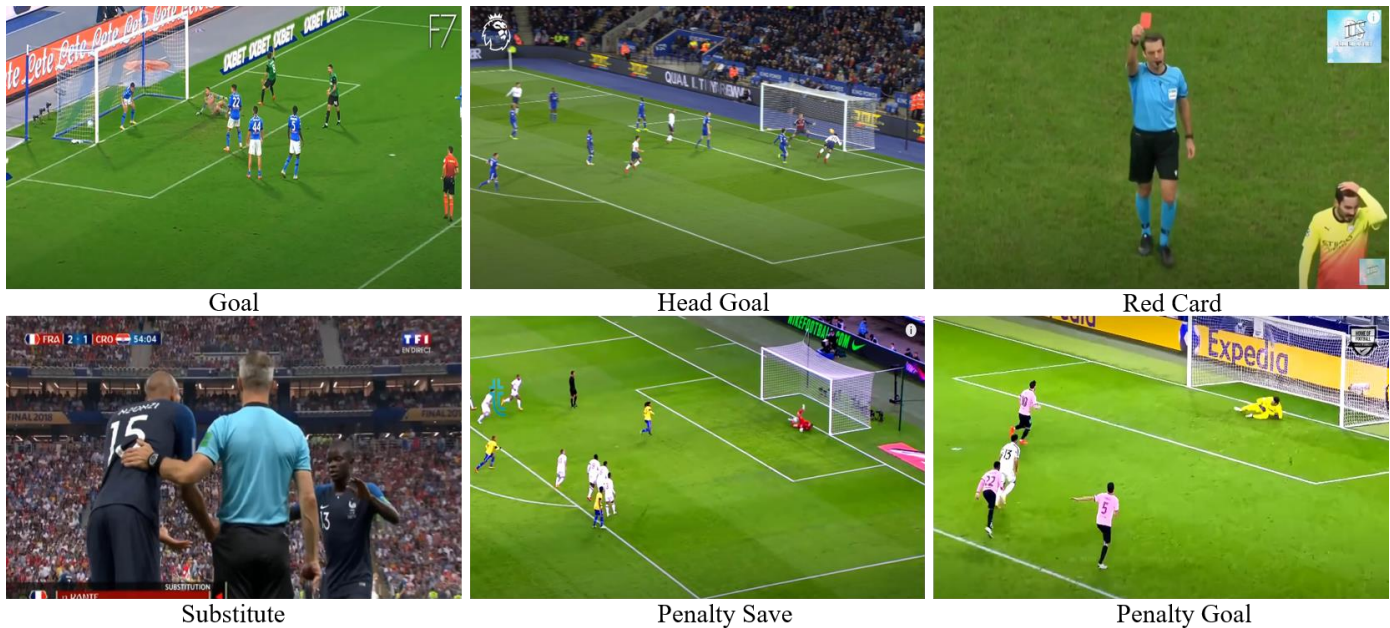


Fig. 3. Sample event classes from our newly created soccer videos dataset.

of LSTM consists of several cell units, where each cell unit contains special gates (input, output, and forget gate) that switch the flow of information and control the sequential pattern recognition process. These gates are configured in such a way that each gate receives input from the previous stage and forwards the computed output to the next gate. All these gates are controlled by a sigmoid function. For instance, the input gate i_t decides that what portion of information should be updated, whereas the output gate o_t stores the information of the coming sequence. The forget gate f_t processes the information from the input gate and the previous cell state and removes the previous information from the memory when needed. The recurrent unit g computes the previous cell state c_{t-1} and the current input x_t using \tanh activation function, whereas h_t can be compute by multiplying the value of output gate with the \tanh of current cell state c_t . The final output can be obtained by passing the h_t to softmax classifier. The mathematical equations of the operations perform by these gates are given in Eq. (1)-(7).

$$i_t = \sigma(w_i * [h_{t-1}, x_t] + b_i) \quad (1)$$

$$o_t = \sigma(w_o * [h_{t-1}, x_t] + b_o) \quad (2)$$

$$f_t = \sigma(w_f * [h_{t-1}, x_t] + b_f) \quad (3)$$

$$g = \tanh(w_g * (x_t + c_{t-1}) + b_g) \quad (4)$$

$$c_t = ((c_{t-1} * f_t) + (g * i_t)) \quad (5)$$

$$h_t = (\tanh(c_t) * o_t) \quad (6)$$

$$output = softmax(h_t) \quad (7)$$

In general, the performance of any deep neural network can be improved by stacking more and more layers; similarly, the hidden sequential pattern leaning capability of an LSTM can be enhanced by increasing the number of layers in the network. Therefore, we add three layers to our LSTM network thereby increasing the ability to analyze the given input data at different

time scales and produce good results as compared to a standard LSTM. Unlike the standard LSTM, when data is fed to the MLSTM, the input data is processed in several layers in a hierarchical fashion where each layer in the network receives the hidden state of the previous layer as an input and forwards the output to the next layer. The computation process of memory cell of the MLSTM is same as the standard LSTM as explained by Eq. (1) to Eq. (7). Fig. 2 depicts the building block of MLSTM, where the first hidden layer receives data from the input layer and the input of the second hidden layer is the output of the first hidden layer. Similarly, the input of the third hidden layer is the output of the second hidden layer. The final output is obtained by computing the output of the final last hidden layer using softmax.

IV EXPERIMENTAL RESULTS AND DISCUSSION

This section presents comprehensive experimental evaluation of our proposed framework and details about the SVE dataset used in experiments. We have first performed the salient events recognition using SVM with HOG descriptor, and then assessed the performance of different state-of-the-art architectures with MLSTM for salient events recognition. Further, the proposed framework is implemented using Matlab 2018 with Matconvnet on PC equipped with 3.60 GHz Intel Core i7 processor and NVIDIA GTX 1080 with 4GB GPU. We have initialized the training with random weight initializer for 60 epochs with batch size of 32. For weights optimization, we have used Adagrad optimizer with learning rate 0.0001. For performance evaluation, we have used five evaluation metrics including Precision, Recall, True Positive Rate (TPR), False Positive Rate (FPR), and F1-score.

4.1 Details of the Dataset

For any advanced computer vision problem, the data acquisition phase is very crucial because without appropriate and sufficient amount of data, one cannot achieve desirable results. Further, the collected data must be labeled properly according to the type of data and nature of the problem. Since

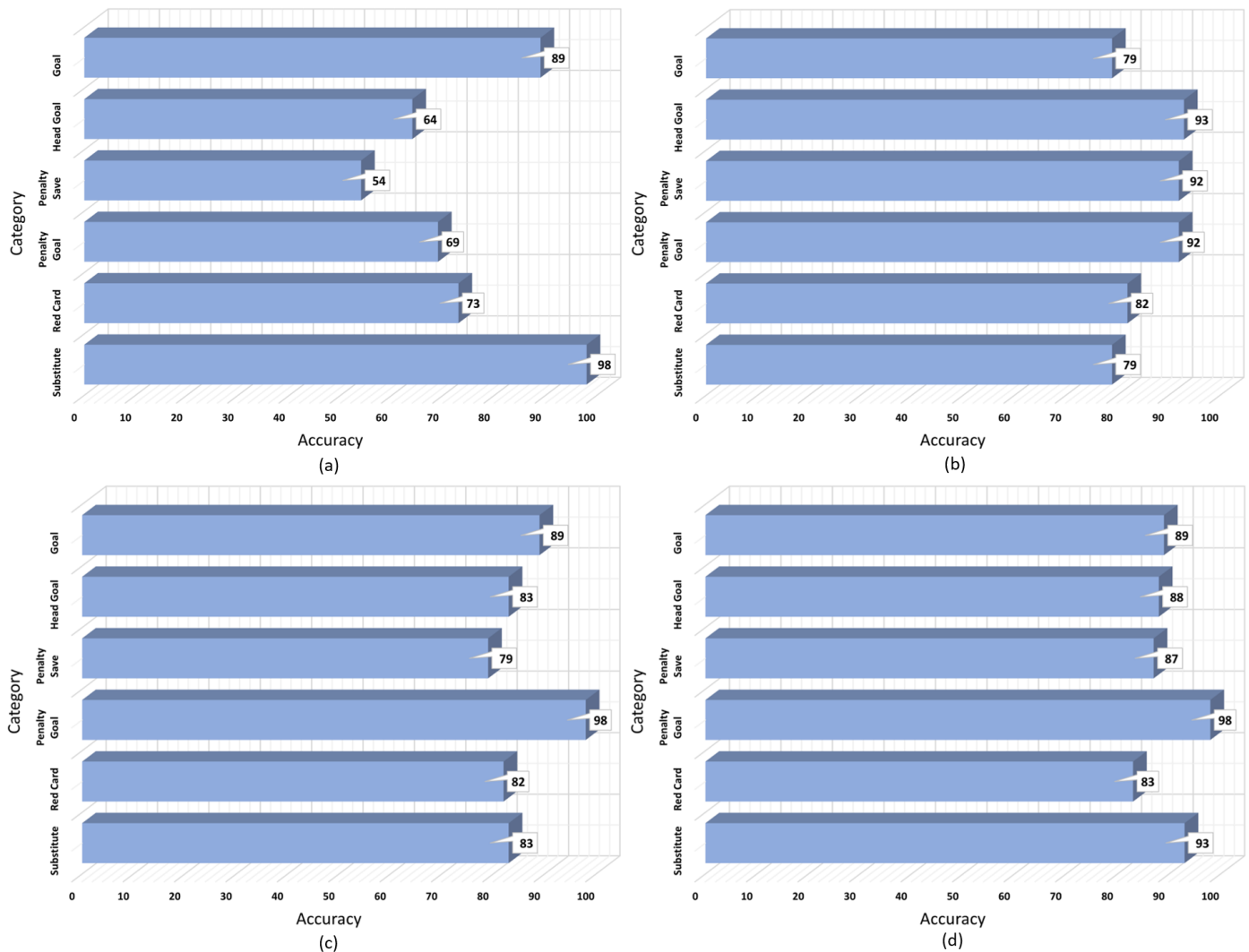


Fig. 4. Category-wise accuracy of four different experimental evaluations of our dataset. (a) HOG+SVM, (b) AlexNet+MLSTM, (c) GoogleNet+MLSTM and (d) ResNet50+MLSTM.

this paper is focusing on recognition of salient events in soccer

TABLE II
STATISTICS OF TRAINING, VALIDATION AND TESTING SETS OF SVD DATASET

	Number of videos clips	Video clip type	Duration in seconds	Frame rate (fps)
Training data	360	MP4	3s – 6s	29
Validation data	120	MP4	3s – 6s	29
Testing data	120	MP4	3s – 6s	29

videos, this research problem requires sufficient amount of labeled soccer video data. To the best of our knowledge, very few soccer videos datasets are presented so far for specific type of tasks including ball tracking, player position, and movement tracking. However, these datasets do not consist of generic type of events, such as, Goal, Substitute, and Red Card, etc. Therefore, in this paper, we present newly created balance SVE dataset of soccer videos, which comprises of short video clips of six different events including goal, penalty save, penalty goal, card, head goal, and substitute. Also, our newly created SVE dataset contains event videos captured from different views with both far and close field of views that offer great variety in the data. The SVE dataset is created in three distinct phases: (i) we collect soccer videos from multiple sources such

as (UEFA Champions League, English Premier League, FIFA World Cup 2018, Bundesliga and Primera Division); (ii) extract event-specific short video clips from the downloaded soccer videos; (iii) annotate the event-specific video clips with start and end boundary of event. The SVE dataset contains a total of 600 short video clips, which are divided into three subsets including train, validation, and test set with split ratio of 60%, 20%, and 20%, respectively. The detailed information of dataset is presented in Table II, where the representative images for each event of our SVD dataset are depicted in Fig. 3.

4.2 Experimental Analysis

1) Experiment 1: SVM Classifier with Histogram of Oriented Gradients Features

We have evaluated the SVE dataset with conventional machine learning technique where we have used a Histogram of Oriented Gradient (HOG) as a feature descriptor and SVM as a classifier to detect the salient events in soccer videos. The HOG descriptor represents gradient orientation and magnitude of objects in a particular region of an image and captures shape-relevant information of detected objects in a video frame. After feature extraction process, we have trained the SVM classifier on extracted features and have evaluated the trained classifier





Event specific video frames	Ground Truth	Predictions	Confidence Score
	Goal	Goal	0.78
	Head Goal	Head Goal	0.63
	Penalty Save	Penalty Goal	0.29
	Penalty Goal	Penalty Goal	0.59
	Red Card	Red Card	0.81
	Substitute	Substitute	0.87

Fig. 5. The predictions of our proposed framework for event recognition in soccer videos, where classifications with high confidence values are indicated in green color, classifications with moderate confidence values are in blue color, and misclassified predictions are in red color.

TABLE III
CONFUSION MATRIX OF SVM(HOG)

Actual Class	Predicted Class					
	Goal	Head Goal	Penalty Save	Penalty Goal	Red Card	Substitute
Goal	18	7	2	2	0	0
Head Goal	2	13	2	1	0	0
Penalty Save	0	0	11	3	0	0
Penalty Goal	0	0	4	14	3	2
Red Card	0	0	0	0	15	1
Substitute	0	0	1	0	5	20
Recall (%)	90.0	65.0	55.0	70.0	75.0	100.0
Precision (%)	62.0	72.0	78.0	77.0	100.0	76.0

on test dataset. The results obtained from the test dataset are shown in Table III, where the diagonal values represent the true positive produce by the SVM classifier. The precision and recall scores from Table III reveal that there is still a considerable room for the improvement of event recognition rate, especially for Head Goal, Penalty Goal, and Penalty save. These scores can be significantly improved using deep learning techniques such as CNN and RNN.

2) Experiment2: Integration of MLSTM with AlexNet Architecture

We have analyzed the soccer events recognition using MLSTM with AlexNet architecture. First, we have extracted discriminative CNN features using a pretrained AlexNet CNN architecture, and then classified the event types by inputting the extracted features to an MLSTM. For feature extraction, we have used the fully connected layer fc-7 of the pretrained AlexNet model as a generic feature descriptor, which converts

TABLE IV
CONFUSION MATRIX OF OUR SOCCER DATASET FOR MLSTM+ALEXNET

Actual Class	Predicted Class					
	Goal	Head Goal	Penalty Save	Penalty Goal	Red Card	Substitute
Goal	16	1	0	0	0	0
Head Goal	4	19	1	1	0	0
Penalty Save	0	0	19	0	0	0
Penalty Goal	0	0	0	19	3	2
Red Card	0	0	0	0	17	2
Substitute	0	0	0	0	0	16
Recall (%)	80.0	95.0	95.0	95.0	85.0	80.0
Precision (%)	94.1	76.0	100.0	79.2	89.5	100.0

a video frame into a 1×4096 feature vector. After features extraction process, MLSTM is trained on extracted features. Finally, we have evaluated our trained model on the test dataset, where 20 video clips per class are given to the trained model for the event recognition task. The obtained results using this approach are presented in Table IV. From Table IV, we can observe that the recognition rate for head goal, penalty goal, penalty save, and red card is improved as compared to the results obtained by SVM(HOG) in experiment 1.

3) Experiment 3: Integration of MLSTM with GoogleNet Architecture

In this set of experiments, we have replaced the AlexNet with a deeper CNN architecture named GoogleNet. It is a deeper architecture with 22 convolution layers and extracts more useful features, which significantly improves the performance of MLSTM for event recognition task. To extract features, we

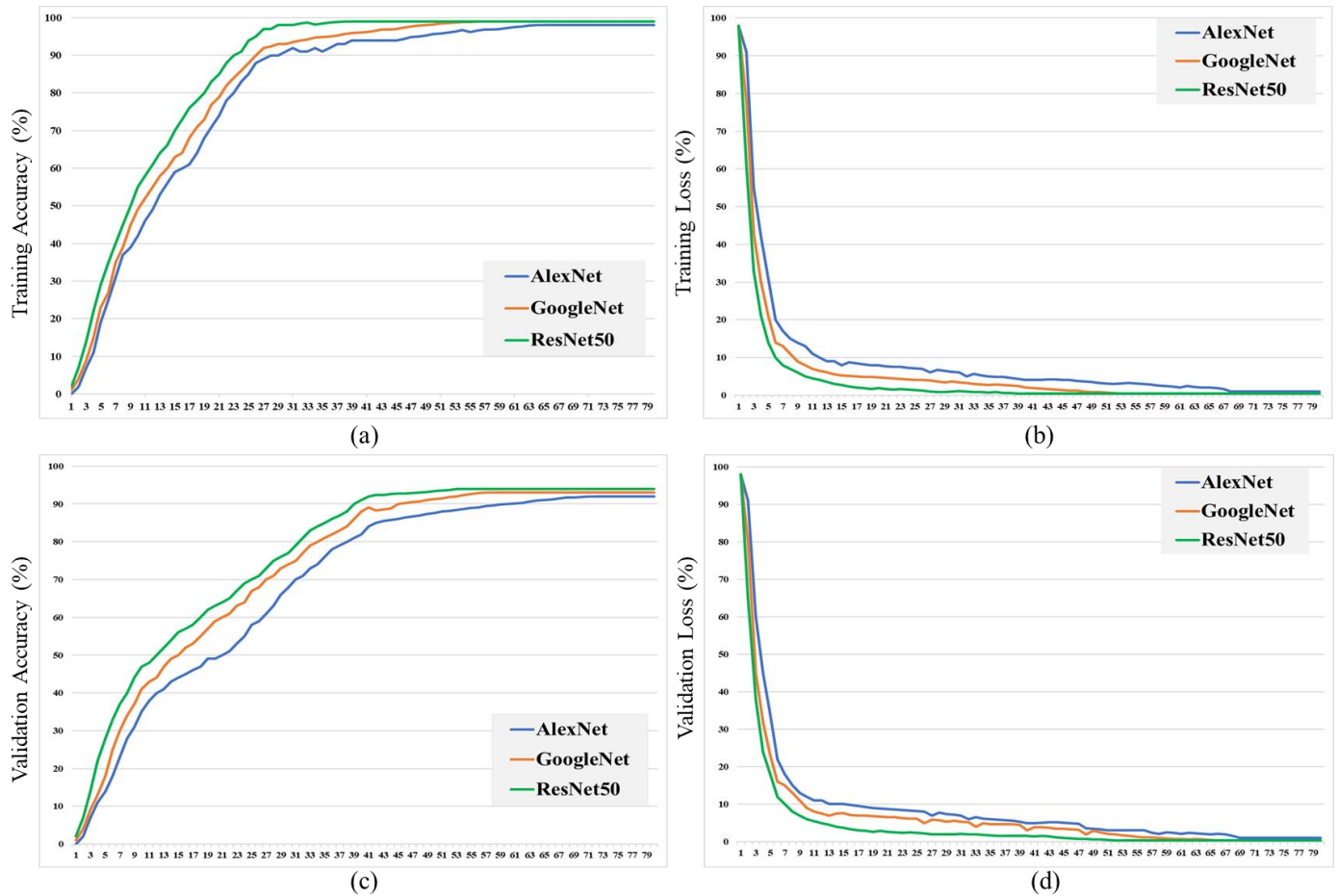


Fig. 6. The training and validation evaluation of our proposed framework along with different investigated techniques for recognition of salient events in soccer videos. (a) Training Accuracy, (b) Training Loss, (c) Validation Accuracy, and (d) Validation Loss.

TABLE V
CONFUSION MATRIX OF OUR SOCCER DATASET FOR
MLSTM+GOOGLENET

Actual Class	Predicted Class					
	Goal	Head Goal	Penalty Save	Penalty Goal	Red Card	Substitute
Goal	18	3	2	0	0	0
Head Goal	2	17	1	0	0	0
Penalty Save	0	0	16	0	0	0
Penalty Goal	0	0	1	20	3	2
Red Card	0	0	0	0	17	1
Substitute	0	0	0	0	0	17
Recall (%)	90.0	85.0	80.0	100.0	85.0	85.0
Precision (%)	78.3	85.0	100.0	76.9	91.8	100.0

have used loss3-classifier as a feature descriptor and have obtained feature vector having a length of 1×1000 . Further, MLSTM is trained on extracted features of length 1×1000 , and then the performance of trained model is evaluated on the test data. We test 20 video clips per event using trained model. The obtained results are presented in Table V. It can be inferred from Table V that the GoogleNet with MLSTM achieves more or less similar results in terms of precision and recall as obtained in our second set of experiments (i.e., Experiment 2) but improved the event recognition rate with minor increment of 0.24%.

4) Experiment 4: Integration of MLSTM with ResNet50 Architecture

Finally, we have evaluated the performance of our ultimate framework which combined ResNet50+MLSTM for recognition of salient events in soccer videos. Our proposed event recognition framework first performs feature extraction process where we extract CNN features from video frames using fully connected layer fc-1000 of ResNet50. After feature extraction, we have trained MLSTM on extracted features. Further, we have performed our model testing where 20 videos per event have been tested on the trained model. The obtained results are presented in Table VI. It can be observed from Table VI that MLSTM with ResNet50 architecture not only achieves the best results in terms of precision and recall, but also improves the overall accuracy on our SVE dataset instances, which are correctly classified. This is indicated by diagonal values, whereas the precision and recall scores are listed at the bottom of Table VI.

4.3 Overall Performance Comparison

In this section, we compared the performance of investigated approaches for salient events recognition in experiments 1 to 4 on test data. Fig. 4 depicts the category-wise accuracy of our

TABLE VI
CONFUSION MATRIX OF OUR SOCCER DATASET FOR
MLSTM+RESNET50

Actual Class	Predicted Class					
	Goal	Head Goal	Penalty Save	Penalty Goal	Red Card	Substitute
Goal	18	1	1	0	0	0
Head Goal	2	18	1	0	0	0
Penalty Save	0	1	18	0	0	0
Penalty Goal	0	0	1	20	2	0
Red Card	0	0	0	0	17	1
Substitute	0	0	0	0	1	19
Recall (%)	90.0	90.0	90.0	100.0	85.0	95.0
Precision (%)	90.0	85.0	94.7	90.9	94.4	95.0

TABLE VII
THE OVERALL COMPARISON OF THE INVESTIGATED METHODS

Method	FPR (%)	TPR (%)	Precision (%)	Recall (%)	F1-Score (%)
SVM (HOG)	0.13	0.81	75.83	77.5	76.65
MLSTM+ AlexNet	0.08	0.92	88.33	89.80	88.05
MLSTM+ GoogleNet	0.07	0.93	87.50	89.11	88.29
MLSTM + ResNet50	0.05	0.96	91.66	91.78	91.74

TABLE VIII
ACCURACY OF THE PROPOSED FRAMEWORK ON DIFFERENT
FRAME SKIP STRATEGIES

Experiments	Frame Skip	Average Time Complexity (Sec)	Accuracy (%)
8 frame skip strategy	8	0.97	89.31
6 frame skip strategy	6	1.19	90.06
Proposed (4 frame skip strategy)	4	1.43	91.74

four different experiments. Results reveal that ResNet50+MLSTM, as adopted in our framework, outperforms other approaches for salient soccer events recognition in terms of accuracy. We further compare the presented approaches in terms of True Positive Rate (TPR), False Positive Rate (FPR), Precision, Recall and F1-score. The results obtained by our proposed framework are presented in Table VII. The evaluation metrics TPR and FPR represent the predicted true positive rate and false positive rate of each investigated method in our soccer dataset. Further we calculated the precision and recall measure of each method. Finally, the F1-Score is computed using precision and recall. The accuracy measures presented in Table VII validate that our proposed solution (LSTM+ResNet50) dominates all pervious investigated methods in terms of FPR, TPR, Precision, Recall, and F1-score. Further, the training and validation performance of our proposed framework and other investigated techniques are given in Fig. 6. Moreover, we have investigated the performance of the proposed framework using three different fame skip schemes for events recognition in soccer videos. Table VIII presents the statistics of the experiments conducted based on different frame skip strategies. It can be noticed in Table VIII that our proposed four frame-skip strategy shows overwhelming performance improvement over other frame skip strategies (i.e., 8 frame skip strategy and 6 frame skip strategy). Therefore, we adopt four frame skip strategy in our approach which enables us to achieve reasonable accuracy with acceptable time complexity.

4.4 Visual Results

We have further evaluated our proposed framework on random videos with predefined events. During evaluation, the proposed framework makes prediction for each video that could be correct or incorrect. While testing input video, our method extracted CNN features with four frame skip strategy. The extracted features are then fed to Multi-layer LSTM for analyzing the video sequences and predict the type of event present in the video. In Fig. 5, each row represents specific events, row 3 is misclassified, where ‘‘Penalty Save’’ is classified as ‘‘Penalty Goal’’. This misclassification is due to the visual similarity of contents, motion of player (i.e., running), and similar background.

4.5 Comparison with Existing Soccer Events Recognition Methods

This section presents the comparative study of our proposed framework with existing soccer events recognition approaches [15, 21, 45]. The results of our proposed framework are evaluated on test videos of our SVE dataset. To validate the effectiveness of our proposed framework, we have compared the proposed framework with five existing soccer events recognition methods. The obtained results are shown in Table IX. For comparison with state-of-the-art methods, we have used accuracy as the evaluation metric. From Table IX, it can be observed that the performance of each method varies from one event to another. For instance, method [21] has the best accuracy for ‘‘Red card’’ event. Whereas our proposed framework dominates the existing soccer events recognition methods, in particular, for detecting ‘‘Goal’’ and ‘‘Penalty Save or Penalty attempt’’ events. Our proposed framework increases the recognition accuracy for ‘‘Goal’’ and ‘‘Penalty save or penalty attempt’’ events by 1.13% and 3.57%, respectively, on average as compared to the existing methods.

V CONCLUSION AND FUTURE WORK

The advent of smart cameras, Nx-IoT, and efficient learning algorithms for sports video analytics will enhance the performance of players as well as facilitate the live spectators inside the stadium. The smart cameras in Nx-IoT soccer environment are interconnected through wireless networks, that capture and transmit the data to an AI-assisted computing platform. Majority of the spectator are very enthusiastic to watch and celebrate the better performance of their favorite teams. The IoT-enabled soccer environment will provide the spectators with live information (visual, and textual) related to the important events of the match and will allow them to share and discuss the match situation in real-time, which can be provided to the spectators as a FinTech service. Therefore, in this article, we have proposed an efficient deep learning-based framework for multi-person salient soccer events recognition in Nx-IoT-enabled environments. The proposed framework recognizes the salient events in soccer video, including goal, substitute, penalty save, penalty goal, red card, and head goal. Our proposed framework examines different CNN architectures with multi-layer LSTM and proposes an efficient CNN+LSTM approach for soccer events learning and recognition in Nx-IoT-enabled environments. Further, we have developed a new soccer dataset SVE, containing six most salient soccer events (i.e., Goal, Red card, Penalty save, Penalty goal, Substitute,

TABLE IX
ACCURACY COMPARISON WITH STATE-OF-THE-ART SOCCER EVENTS RECOGNITION METHODS

Method	Goal (%)	Head goal (%)	Penalty goal (%)	Penalty save/attempt (%)	Red card (%)	Substitute (%)
[45]	88.03	-	-	-	90.35	-
[21]	90.29	-	-	86.48	96.42	-
[15]	89.68	-	-	91.66	86.66	-
Our Method	90.81	87.42	92.29	95.23	89.45	95.0

Head goal). The results obtained from the experimental evaluation validate the suitability and accuracy of our proposed framework for soccer events recognition in FinTech-enabled Nx-IoT environments.

This paper mainly focuses on the recognition of salient soccer events in IoT environment using combined CNN+MLSTM deep learning framework. Although, the current approach uses an efficient CNN architecture in terms of feature enrichment, the series of residual blocks employed in this approach increase the overall computation complexity. Also, the proposed system has no suitable mechanism for ball tracking and player position tracking in the ground field. Moreover, the current system sometime misclassifies Penalty Save event as a Penalty Goal. Considering these limitations of our current method, in future we are aiming to use a light-weight CNN architecture having lower computation complexity. Further, we have intentions to extend our proposed framework by introducing more robust and discriminative features for efficient event recognition task, such as optical flow, motion saliency, and C3D features (C3D features are utilized by 3D ConvNets) along with more robust sequence learning algorithm such as Gated Recurrent Unit (GRU). Furthermore, we plan to integrate other modules, such as player position tracking, player identification, and soccer ball tracking in our current framework for efficient soccer events recognition and streaming in Nx-IoT-enabled environments.

REFERENCES

[1] Worldatlas. "The most popular sports in the world." Accessed: Feb. 13, 2020. [Online]. Available: <https://www.worldatlas.com/articles/what-are-the-most-popular-sports-in-the-world.html>.

[2] M. Dive. "Google: Watch time for YouTube sports highlights jumps 80%." Accessed: Feb. 14, 2020. [Online]. Available: <https://www.marketingdive.com/news/google-watch-time-for-youtube-sports-highlights-jumps-80/516281/>.

[3] M. H. Kolekar and S. Sengupta, "Bayesian network-based customized highlight generation for broadcast soccer videos," *IEEE Transactions on Broadcasting*, vol. 61, no. 2, pp. 195-209, 2015.

[4] M. G. I. Rathod and M. D. A. Nikam, "Review on event retrieval in soccer video," *Int. J. Comput. Sci. Inf. Technol.*, vol. 5, pp. 5601-5605, 2014.

[5] A. Ghosh and C. Jawahar, "SmartTennisTV: Automatic indexing of tennis videos," in *National Conference on Computer Vision, Pattern Recognition, Image Processing, and Graphics*, 2017: Springer, pp. 24-33.

[6] M. Xu, N. C. Maddage, C. Xu, M. Kankanhalli, and Q. Tian, "Creating audio keywords for event detection in soccer video," in *2003 International Conference on Multimedia and Expo. ICME'03. Proceedings (Cat. No. 03TH8698)*, 2003, vol. 2: IEEE, pp. II-281.

[7] Q. Ye, Q. Huang, W. Gao, and S. Jiang, "Exciting event detection in broadcast soccer video with mid-level description and incremental learning," in *Proceedings of the 13th annual ACM international conference on Multimedia*, 2005, pp. 455-458.

[8] W. Zhao, Y. Lu, H. Jiang, and W. Huang, "Event detection in soccer videos using shot focus identification," in *2015 3rd IAPR Asian Conference on Pattern Recognition (ACPR)*, 2015: IEEE, pp. 341-345.

[9] A. Raventos, R. Quijada, L. Torres, and F. Tarrés, "Automatic summarization of soccer highlights using audio-visual descriptors," *SpringerPlus*, vol. 4, no. 1, pp. 1-19, 2015.

[10] M.-H. Sigari, H. Soltanian-Zadeh, and H.-R. Pourreza, "A framework for dynamic restructuring of semantic video analysis systems based on learning attention control," *Image and Vision Computing*, vol. 53, pp. 20-34, 2016.

[11] B. Fakhar, H. R. Kanan, and A. Behrad, "Event detection in soccer videos using unsupervised learning of Spatio-temporal features based on pooled spatial pyramid model," *Multimedia Tools and Applications*, vol. 78, no. 12, pp. 16995-17025, 2019.

[12] Z. Wang, J. Yu, and Y. He, "Soccer video event annotation by synchronization of attack-defense clips and match reports with coarse-grained time information," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 27, no. 5, pp. 1104-1117, 2016.

[13] W. Huang and Z. Wang, "Soccer Video Event Detection Using 3D Convolutional Networks and Shot Boundary Detection via Deep Feature Distance," in *Neural Information Processing: 24th International Conference, ICONIP 2017, Guangzhou, China, November 14-18, 2017, Proceedings*, 2017, vol. 10635: Springer, p. 440.

[14] M. Z. Khan, S. Saleem, M. A. Hassan, and M. U. G. Khan, "Learning deep C3D features for soccer video event detection," in *2018 14th International Conference on Emerging Technologies (ICET)*, 2018: IEEE, pp. 1-6.

[15] H. Jiang, Y. Lu, and J. Xue, "Automatic soccer video event detection based on a deep neural network combined cnn and rnn," in *2016 IEEE 28th International Conference on Tools with Artificial Intelligence (ICTAI)*, 2016: IEEE, pp. 490-494.

[16] G. Yaparla, S. Allaparthi, S. K. Munnangi, G. Ramamurthy, and G. Canaria, "A Novel framework for Fine Grained Action Recognition in Soccer," in *proceedings of International Work-Conference on Artificial Neural Networks 2019*, Springer, pp. 137-150.

[17] S. A. Pettersen *et al.*, "Soccer video and player position dataset," in *Proceedings of the 5th ACM Multimedia Systems Conference*, 2014, pp. 18-23.

[18] J. Yu, A. Lei, Z. Song, T. Wang, H. Cai, and N. Feng, "Comprehensive dataset of broadcast soccer videos," in *2018 IEEE Conference on Multimedia Information Processing and Retrieval (MIPR)*, 2018: IEEE, pp. 418-423.

[19] N. Feng *et al.*, "SSET: a dataset for shot segmentation, event detection, player tracking in soccer videos," *Multimedia Tools and Applications*, vol. 79, no. 39, pp. 28971-28992, 2020.

[20] H. Ullah and M. Sajjad, "Salient Event Detection in Soccer Videos using Histogram of Oriented Gradient" in *Proceedings of the 4th International Conference on Next Generation Computing (ICNGC)*, 2018.

[21] M. Tavassolipour, M. Karimian, S. J. I. T. o. c. Kasaei, and s. f. v. technology, "Event detection and summarization in soccer videos using bayesian network and copula", *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 24, no. 2, pp. 291-304, 2013.

[22] M. H. Kolekar and S. J. I. T. o. B. Sengupta, "Bayesian network-based customized highlight generation for broadcast soccer videos," *IEEE Transactions on Broadcasting*, vol. 61, no. 2, pp. 195-209, 2015.

[23] Z. Wang, J. Yu, Y. J. I. T. o. C. He, and S. f. V. Technology, "Soccer video event annotation by synchronization of attack-defense clips and match reports with coarse-grained time information", *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 27, no. 5, pp. 1104-1117, 2016.

[24] B. Fakhar, H. R. Kanan, A. J. J. o. V. C. Behrad, and I. Representation, "Learning an event-oriented and discriminative dictionary based on an adaptive label-consistent K-SVD method for event detection in soccer videos", *Journal of Visual Communication and Image Representation*, vol. 55, pp. 489-503, 2018.

[25] K. J. Bennett, A. R. Novak, M. A. Pluss, A. J. Coutts, J. J. J. o. s. Fransen, and m. i. sport, "Assessing the validity of a video-based

decision-making assessment for talent identification in youth soccer", *Journal of Science and Medicine in sport*, vol. 22, no. 6, pp. 729-734, 2019.

[26] M.-S. Hosseini and A.-M. Eftekhari-Moghadam, "Fuzzy rule-based reasoning approach for event detection and annotation of broadcast soccer video," *Applied Soft Computing*, vol. 13, no. 2, pp. 846-866, 2013.

[27] S. Khan, K. Muhammad, S. Mumtaz, S. W. Baik, and V. H. C. de Albuquerque, "Energy-efficient deep CNN for smoke detection in foggy IoT environment," *IEEE Internet of Things Journal*, vol. 6, no. 6, pp. 9237-9245, 2019.

[28] K. Muhammad, J. Ahmad, Z. Lv, P. Bellavista, P. Yang, and S. W. Baik, "Efficient deep CNN-based fire detection and localization in video surveillance applications," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 49, no. 7, pp. 1419-1434, 2018.

[29] J. Jia *et al.*, "EMBDN: An Efficient Multiclass Barcode Detection Network for Complicated Environments," *IEEE Internet of Things Journal*, vol. 6, no. 6, pp. 9919-9933, pp. 1483-1498, 2019.

[30] Z. Cai and N. Vasconcelos, "Cascade r-cnn: High quality object detection and instance segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no.5, pp. 1483-1498, 2019.

[31] T. Liu, H. Liu, Y.-F. Li, Z. Chen, Z. Zhang, and S. Liu, "Flexible FTIR spectral imaging enhancement for industrial robot infrared vision sensing," *IEEE Transactions on Industrial Informatics*, vol. 16, no. 1, pp. 544-554, 2019.

[32] W. Ren *et al.*, "Low-light image enhancement via a deep hybrid network," *IEEE Transactions on Image Processing*, vol. 28, no. 9, pp. 4364-4375, 2019.

[33] L. Liu, G. Feng, D. Beutemps, and X.-P. Zhang, "Re-synchronization using the Hand Preceding Model for Multi-modal Fusion in Automatic Continuous Cued Speech Recognition," *IEEE Transactions on Multimedia*, vol.23, pp. 292-305, 2020.

[34] A. Chowdhury and A. Ross, "Fusing MFCC and LPC Features using 1D Triplet CNN for Speaker Recognition in Severely Degraded Audio Signals," *IEEE Transactions on Information Forensics and Security*, vol. 15, pp. 1616-1629, 2019.

[35] A. Ullah, K. Muhammad, J. Del Ser, S. W. Baik, and V. H. C. de Albuquerque, "Activity recognition using temporal optical flow convolutional features and multilayer LSTM," *IEEE Transactions on Industrial Electronics*, vol. 66, no. 12, pp. 9692-9702, 2018.

[36] A. Ullah, K. Muhammad, I. U. Haq, and S. W. Baik, "Action recognition using optimized deep autoencoder and CNN for surveillance data streams of non-stationary environments," *Future Generation Computer Systems*, vol. 96, pp. 386-397, 2019.

[37] T. Tsunoda, Y. Komori, M. Matsugu, and T. Harada, "Football action recognition using hierarchical LSTM," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2017, pp. 99-107.

[38] M. Fani, M. Yazdi, D. A. Clausi, and A. J. I. A. Wong, "Soccer video structure analysis by parallel feature fusion network and hidden-to-observable transferring markov model", *IEEE Access*, vol. 5, pp. 27322-27336, 2017.

[39] S. Giancola, M. Amine, T. Dghaily, and B. Ghanem, "Soccernet: A scalable dataset for action spotting in soccer videos," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018, pp. 1711-1721.

[40] E. Fenil, G. Manogaran, G. Vivekananda, T. Thanjaivadeivel, S. Jeeva, and A. J. C. N. Ahilan, "Real time violence detection framework for football stadium comprising of big data analysis and deep learning through bidirectional LSTM", *Computer Networks*, vol. 151, pp. 191-200, 2019.

[41] T. Liu *et al.*, "Soccer video event detection using 3D convolutional networks and shot boundary detection via deep feature distance," in *International Conference on Neural Information Processing*, 2017: Springer, pp. 440-449.

[42] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278-2324, 1998.

[43] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770-778.

[44] S. Hochreiter and J. J. N. c. Schmidhuber, "Long short-term memory", *Neural Computation*, vol. 9, no. 8, pp. 1735-1780, 1997.

[45] C.-L. Huang, H.-C. Shih, and C.-Y. J. I. T. o. M. Chao, "Semantic analysis of soccer video using dynamic Bayesian network", *IEEE Transactions on Multimedia*, vol. 8, no. 4, pp. 749-760, 2006.



KHAN MUHAMMAD (Member, IEEE) received the Ph.D. degree in digital contents from Sejong University, Republic of Korea, in February 2019. He has been working as an Assistant Professor with the Department of Software since March 2019. He is currently the Director of the Visual Analytics for Knowledge Laboratory (VIS2KNOW Lab), Sejong University, Seoul, South Korea. He has registered eight patents in South Korea (seven)/Australia (one) and has contributed more than 150 articles in peer-reviewed journals and conference proceedings in his areas of research. He was recently selected among top 100 000 scientists around the globe by Stanford Researchers List. His research interests include intelligent video surveillance, medical image analysis, information security, video summarization, multimedia data analysis, computer vision, the IoT/IoMT, and smart cities.



Hayat Ullah received the B.S. and MS degree in Computer Science from Islamia College University, Peshawar, Pakistan and Sejong University, Seoul, Republic of Korea, respectively. He secured a scholarship-based PhD Fellowship at Intelligent Systems, Computer Architecture, Analytics, and Security (ISCAAS) Laboratory at the Department of Computer Science, Kansas State University, Manhattan, USA and will officially start his PhD soon. His research interests include image processing, sports video analytics, deep learning, computer vision, image enhancement, and image/video quality assessment.



Mohammad S. Obaidat (Fellow of IEEE 2005, and SCS Fellow 2000) is an internationally known academic/researcher/scientist/scholar. He received his Ph.D. degree in Computer Engineering with a minor in Computer Science from the Ohio State University, Columbus, USA. He has received extensive research funding and published to date (2019) about One Thousand and Two Hundred (1,200) refereed technical articles. About half of them are journal articles, over 95 books, and about 70 Book Chapters. He is now the Founding Dean and Professor, College of Computing and Informatics at the University of Sharjah, UAE. He has chaired numerous (Over 175) international conferences and has given numerous (Over 175) keynote speeches worldwide. He received many best paper awards for his papers. He also received Best Paper awards from IEEE Systems Journal in 2018 and in 2019 (2 Best Paper Awards). In 2020, he received 4 best paper awards from IEEE Systems Journal. In 2021, he also received the IEEE Systems best paper award. In 2021, he was ranked by Guide2Research as Number 1 Computer Scientist in UE in terms of Number of Publications.



Amin Ullah received Ph.D. degree in digital contents from Sejong University, Seoul, South Korea. He is currently working as a Postdoc Researcher at the CoRIS Institute, Oregon State University, Corvallis, Oregon, USA. His major research focus is on human action and activity recognition, sequence learning, image and video analytics, content-

based indexing and retrieval, 3D point clouds, IoT and smart cities, and deep learning for multimedia understanding. He has published several papers in reputed peer reviewed international journals and conferences including IEEE Transactions on Industrial Electronics, IEEE Transactions on Industrial Informatics, IEEE Transactions on Intelligent Transportation Systems, IEEE Internet of Things Journal, IEEE Access, Elsevier Future Generation Computer Systems, Elsevier Applied Soft Computing, International Journal of Intelligent Systems, Springer Multimedia Tools and Applications, Springer Mobile Networks and Applications, and IEEE Joint Conference on Neural Networks.



Arslan Munir (M'09, SM'17) is currently an Associate Professor in the Department of Computer Science at Kansas State University. He was a postdoctoral research associate in the Electrical and Computer Engineering department at Rice University, Houston, Texas, USA from May 2012 to June 2014. He received his M.A.Sc. in

Electrical and Computer Engineering from the University of British Columbia, Vancouver, Canada, in 2007 and his Ph.D. in Electrical and Computer Engineering from the University of Florida, Gainesville, Florida, USA, in 2012. From 2007 to 2008, he worked as a software development engineer at Mentor Graphics Corporation in the Embedded Systems Division.

Munir's current research interests include embedded and cyber-physical systems, secure and trustworthy systems, parallel computing, artificial intelligence, and computer vision. Munir received many academic awards including the doctoral fellowship from Natural Sciences and Engineering Research Council (NSERC) of Canada. He earned gold medals for best performance in electrical engineering, gold medals and academic roll of honor for securing rank one in pre-engineering provincial examinations (out of approximately 300,000 candidates). He is a Senior Member of IEEE.



Muhammad Sajjad received the master's degree from the Department of Computer Science, College of Signals, National University of Sciences and Technology, Rawalpindi, Pakistan in 2012, and the Ph.D. degree in digital contents from Sejong University, Seoul, South Korea in 2015. He is currently working as an ERCIM Research

Fellow at NTNU, Norway. He is an Associate Professor with the Department of Computer Science, Islamia College University Peshawar, Pakistan. He is also the Head of the Digital Image Processing Laboratory with Islamia College University Peshawar, where many students are involved in different research projects under his supervision, such as Big

data analytics, medical image analysis, multi-modal data mining and summarization, image/video prioritization and ranking, fog computing, the Internet of Things, autonomous navigation, and video analytics. His primary research interests include computer vision, image understanding, pattern recognition, robotic vision, and multimedia applications, with current emphasis on economical hardware and deep learning, video scene understanding, activity analysis, fog computing, the Internet of Things, and real-time tracking. He has published more than 65 papers in peer-reviewed international journals and conferences. He is serving as a professional reviewer for various well-reputed journals and conferences. Currently, he is the associate editor at IEEE Access and acting as a guest editor at IEEE Transactions on Intelligent Transportation Systems.



Victor Hugo C. de Albuquerque [M'17, SM'19] is Professor and senior researcher at the Department of Teleinformatics Engineering/Graduate Program on Teleinformatics Engineering at the Federal University of Ceará, Brazil. He has a Ph.D in Mechanical Engineering from the Federal University of Paraíba (UFPB, 2010), an MSc

in Teleinformatics Engineering from the Federal University of Ceará (UFC, 2007), and he graduated in Mechatronics Engineering at the Federal Center of Technological Education of Ceará (CEFETCE, 2006). He is a specialist, mainly, in Image Data Science, IoT, Machine/Deep Learning, Pattern Recognition, Automation and Control, and Robotics.