




## Noise-Resilient Human Activity Recognition via A Hybrid CNN–InceptionV3 Model

Kothandaraman Ishwarya<sup>1\*</sup>, A. Alice Nithya<sup>2</sup>, Saraswathi Sudalaimuthu<sup>3</sup>

<sup>1</sup> Department of Computing Technologies, SRM Institute of Science and Technology, Kattankulathur 603203, India

<sup>2</sup> Department of Computational Intelligence, SRM Institute of Science and Technology, Kattankulathur 603203, India

<sup>3</sup> Department of Computer Science and Engineering, Sathyabama Institute of Science and Technology, Chennai 600119, India

Corresponding Author Email: [ishwaryk3@srmist.edu.in](mailto:ishwaryk3@srmist.edu.in)

Copyright: ©2025 The authors. This article is published by IETA and is licensed under the CC BY 4.0 license (<http://creativecommons.org/licenses/by/4.0/>).

<https://doi.org/10.18280/isi.300818>

### ABSTRACT

**Received:** 6 July 2025

**Revised:** 12 August 2025

**Accepted:** 20 August 2025

**Available online:** 31 August 2025

#### Keywords:

*human activity recognition, noise-resilient deep learning, pyramidal flow feature fusion, inertial sensor features, support vector machine, kalman filter, butterworth filter*

Human activity recognition (HAR) is the process of identifying and classifying physical activities performed by individuals through data captured from various sensors. It is essential for applications like sports analysis, rehabilitation, patient monitoring, and senior care systems. From the literature it is observed that human activity recognition (HAR) models developed in an unconstrained environment have several limitations like Personal Interference (PI), Electromagnetic (EM) noise, in-band noise or human movements and outliers involved while capturing the input data. These noises are affecting the overall performance and robustness of the model. In order to improve the model performance, noise removal techniques are introduced in this work. Noises like salt and pepper noise, Gaussian noise, and blurring of the boundaries and outlier treatment are processed for the hybrid data acquired using video and sensors in line of sight mode in this paper. For removing these noises, a combination of filters like Kalman filter, MOSSE filter, Butter-worth and J filter, are applied to the input data. By developing this noise removal technique, Peak to Sidelobe Ratio (PSR) is reduced from the raw data. After removing these noises, features are extracted using the top layers of CNN InceptionV3 model for video data. Similarly by using inertial sensor features like tri-axial accelerometer, gyroscopes and magnetometers are collected and a feature vector is created. Pyramidal flow feature fusion (PFFF) technique is used to fuse the extracted features. Finally, the fused features are given to the Support Vector Machine (SVM) classifier to perform activity recognition. This work presents a hybrid deep learning model designed to enhance noise resilience in human activity recognition (HAR) for unconstrained environments and it was validated on UCF Sports and UCI HAR datasets, showing that noise removal with hybrid feature fusion markedly improves HAR model robustness.

## 1. INTRODUCTION

Recognizing activities from a sequence of observations of people's behaviors and ambient factors is the goal of human activity recognition (HAR). [1]. Computer vision and human-computer interaction experts have been working on HAR for the last three decades in proposing several methods and techniques for improving the process. The first work on HAR dates back to the late '90s [2]. Experiments with solutions that are able to identify Activities of Daily Living (ADLs) from inertial signals have been a substantial portion of recent research [3]. Both the growing popularity of mobile devices with inertial sensors and the falling cost of hardware are mostly to blame for this. Utilizing smartphones that can receive and process signals creates opportunity in a number of application scenarios, including healthcare, smart homes, and surveillance.

Some of the application areas where HAR finds its importance are as follows:

i. Object recognition and monitoring are widely utilized in various fields, including security monitoring, interpersonal

behaviour, video communications, and traffic management [4].

ii. The identification and monitoring of pedestrian items are often integrated with surveillance cameras to find the targeted things [5].

iii. Non-periodic complex motion state (NP\_CMS) activities, exemplified by sports such as badminton and basketball, are characterized by random transitions between multiple motion states instead of cyclic or repetitive movements [6]. For example, during a badminton game, a player alternates unpredictably between serving, swinging, smashing, and other actions. Each motion state occurs only once before transitioning to another, presenting a key challenge for activity recognition systems — accurately determining the time span and boundaries of each motion state.

The generic HAR framework involves the following stages: data acquisition, preprocessing, feature extraction and activity recognition, as illustrated in Figure 1.

i. The data acquisition stage is responsible for acquiring data from different types of sensors. Data generally originates from

sensors such as video sensors, inertial sensors, ECG signals, EEG signals, and so on. Data acquired from the sensors typically introduce artifacts and noise into the input data due to many reasons, such as electronic fluctuation, sensor calibration, and malfunctions, acquisition in unconstrained

environments like non line of sight and so on.

ii. Generally filtering techniques are used to preprocess the noise introduced in the input data. The output of this step is a set of filtered noise free data that constitute the input for the next step.

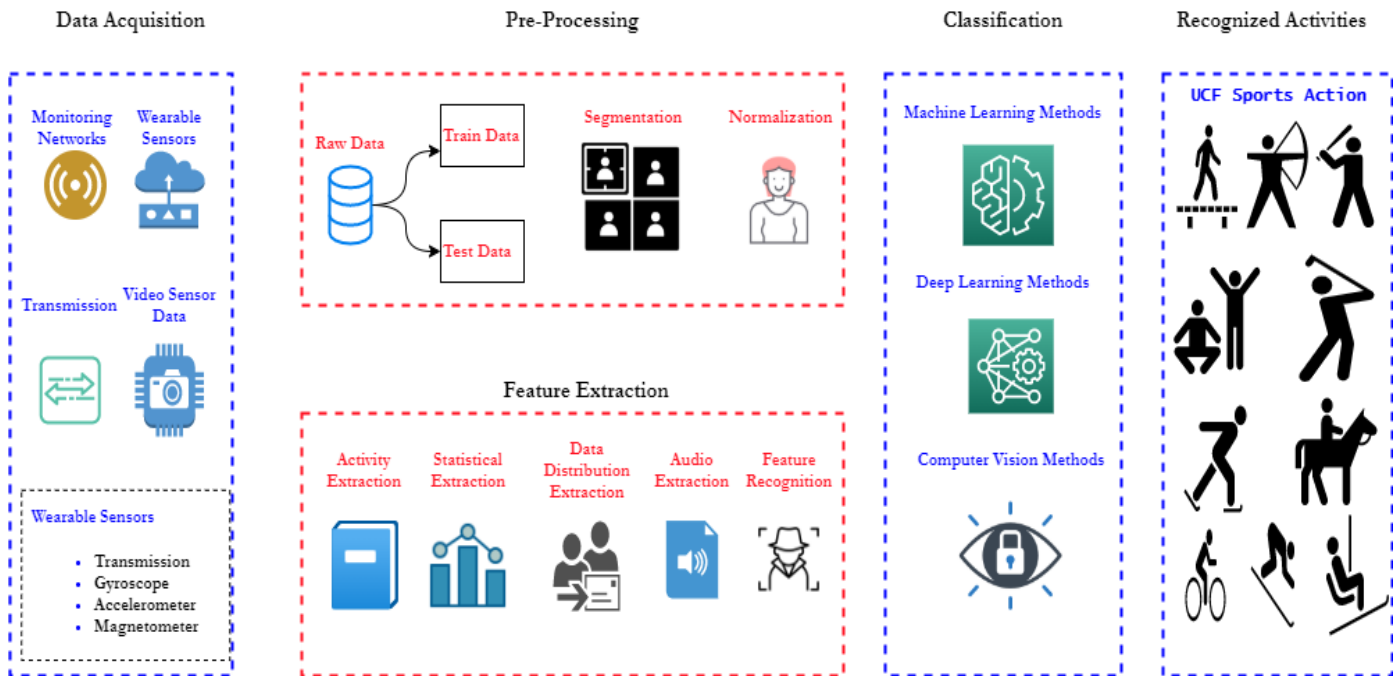


Figure 1. Generic HAR framework

iii. The preprocessed data is given to a feature extraction stage to extract the most significant and relevant information or pattern. This stage also helps in reducing the data dimension.

iv. The recognition of human activities is the last step of this process. During the training phase of the recognition process, features are given to the classifier and mapped to their respective activities like walking, running, diving, golf swing, kicking, lifting, and so on. In the testing phase a new data will be given to the trained model and the activity performed will be recognized. Finally the developed human activity recognition models accuracy and efficiency are computed.

In recent years, researchers have used several publicly available datasets to develop efficient HAR models. Each and every dataset has developed to solve and work more specific tasks. Some of the publicly available datasets are as follows Human Motion Database (HMDB51) [7], Skeleton Based Action Recognition (MSR) action 3D dataset [8], Physical Activity Monitoring (PAMAP2) dataset [9], University of Central Florida (UCF) sports dataset [10], Hollywood2 Actions [11], and YouTube Actions [12]. From the dataset it is observed that the input for developing a HAR model could be acquired from different types of sensors like inertial sensor Gyroscope, Accelerometer, Magnetometer, video sensors and so on.

The HMDB51 data-set contains realistic footage from films and websites. Six thousand eight hundred and forty-nine visual aids from 51 activity types (such as leap, hug, and smile) make up the data-set. The MSR-Action3D data-set includes a total of 20 actions that were carried out by a total of ten distinct actors. Each action was done by each performer twice or three times. PAMAP2 data-set contains wearable sensor data from 18 everyday activities. This data-set includes activities

including driving, Cycling, running and walking. The Hollywood dataset is often used as a publicly available human activity data-set. It has 6 kinds of video activity. Pioneering action recognition for sports data-set and its related activities are in the UCF Sports dataset. It is a collection of actions from different sports, suggesting RNN and GRNN a new method that analyses and extracts all characteristics from each time and frame of video. This hybrid approach improves HAR.

In this work, a hybrid human activity recognition model has been developed using the UCF-sports data-set. Human actions, unlike voice recognition, have no syntax or precise meaning; this creates two issues. An activity might vary from topic to subject, resulting in Intra-class differences, increasing inter-class disparities. Conversely, many activities may exhibit similar forms. The inter-class similarity is a frequent HAR occurrence; specific characteristics from activity videos must be developed to solve these issues.

Considering the above facts the following contributions are made in our research human activity recognition which recognize activities from a series of observations based on human actions and environmental conditions. The noise in the HAR model will affect the overall performance and robustness of the model. This noise can be removed by using different types of filter. By using this technique PSR is reduced from the raw data. We extract features and a feature vector is created. Next we fuse features by using PFFF technique and given to the SVM classifier to predict the HAR model.

It is important to note that deep learning is a subset of machine learning. While machine learning encompasses a broad range of algorithms designed to learn patterns from data, deep learning focuses on neural networks with multiple layers (hence 'deep'), allowing for more complex representations of data. In this work, we leverage deep learning models,

particularly Convolutional Neural Networks (CNNs), which are highly effective for tasks involving spatial data such as video-based human activity recognition.

HAR performance is greatly impacted by noise. Our tests, for instance, demonstrate that adding noise lowers accuracy on the UCI HAR dataset from 97.1% to 96.2% (0.9% drop) and on the UCF Sports dataset from 98.7% to 98.4% (0.3% increase). The importance of noise-resilient HAR models is highlighted by the fact that, despite these seemingly insignificant values, even slight decreases in accuracy can result in misclassification in crucial applications like healthcare monitoring and surveillance.

The structure of the work is presented: Section II provides a detailed literature survey. Section III briefly explains HAR and its usage in diverse settings; Section III explains HAR systems utilize different sensing methods, which are described. Section IV analyses the results. Section V conclusion and future scope.

## 2. RELATED WORK ON ACTIVITY RECOGNITION

With deep learning advancing rapidly in areas like image classification and object detection, researchers have shown growing interest in understanding human activities through video. Many studies have explored and compared both traditional and deep learning methods for recognizing human actions [13, 14]. This paper provides a comprehensive overview of advanced deep learning techniques for video-sequential data processing, emphasizing the contribution of Sparse Block (SB) architectures in enhancing computational efficiency, optimizing performance, and reducing the structural complexity of noise suppression frameworks.

Ma et al. [15] introduced SB, there are 12 layers of two types: dilation ReLU + BN + Conv. The Feature Enhancement Block (FEB) has four layers (ReLU + BN + Conv, ReLU + BN + Conv), whereas the Attention Block (AB) has one. The AB guided the SB and FEB, helpful for unpredictable noises. Finally, the Reconstruction Block (RB) reconstructs the picture. Guo et al [16] focused on the noise estimation reduction network based on their study (NERNet). NERNet decreased the amount of noise in pictures that included actual noise. The noise estimating component and the noise-reducing module were separated into two separate modules in the design. This module applies the noise-level map to the symmetrical dilation blocks [16, 17] and the pyramids feature fusion [18] to get the suitable noise-level map. On the other hand, He et al. [18] proposed a formulation for the noise removal module made use of the approximated map to eliminate the unwanted background sounds. The domestic and global information necessary for retaining features and textures was compiled into a single component for removing purposes. They gave the result of the noisy estimation modules to the removal modules, which resulted in clean pictures.

Recently Liu et al. [19] proposed an extended version of CNN efficiently detecting noisy features and image patching. A network with plenty of training phases and picture patches is the result. As a result, suggested the patching complex local division and deep conquering network (PCLDCNet). The training is performed in its particular region and split into local sub-tasks. The local sub-task was merged with every noisy patch weighting combination. Using hierarchical residual learning instead of identification mappings for picture denoizing was suggested in this work. Fusion, inference and

Features extraction are all parts of the network. The sub-network for feature extraction removes patches that represent higher-dimensional feature maps from the leading network.

Numerous cascaded convolution layers are used in the interference sub-network [20], resulting in an extensive perception. Shi et al. [20] summarized learning noise maps from multi-resolution data and creating tolerant mistakes in noise estimates using the cascaded method with other procedures. Similarly, Liu et al. [21] tackled the problem last but not least, the fusion sub-network combines the whole noise map to generate an estimate. The expectation-maximization methodology utilized this technique. Zhang et al. [22] used sparse coding and keyword updates to minimize the issue and suggested separating aggregated networks to avoid vanishing problems (SANet). SANet removed noise using three blocks: band aggregation, deep mapping and convolutional separation.

To denoise images, Quan et al. [23] suggested the complex-valued CNN (CDNet). First, that fed the input picture into 24 convolutional units with Sequential Connection (SCCU). SCCU uses complex-valued convolutional layers, ReLU, and BN. The network utilized a 64 convolutional kernel. The remaining 18 units got the residual block. The stride of two CNN model layers improved computing efficiency. They created a picture with absolute values by combining the complex-valued features. CDNet comprises five blocks: Merging layer, CV RB, CV BN, CV ReLU and CV Conv. Human activity identification is currently accomplished via the use of several different sensors, including acceleration sensors, magnetometers, and gyroscopes, which are used in conjunction with one another. Some studies have also shown that using a mix of different primary objectives may improve outcomes than using a single sensor individually [24].

Chathuramali et al. [25] used the silhouette distribution and the optically flow data to derive a motion spatial feature selection method describing video frames. Both features are computing using the Leonida et al. [26] method within a normalized bounding box. The multi-class SVM algorithm classified the activity. This method is designed to deal with lengthy training times and large feature vectors. Simonyan et al. [27] explores a spatial and temporal convNets architecture with two streams. Activity recognition in the spatial streams uses RGB video sequence, whereas action recognition in the temporal stream uses motion information acquired by layering dense optical flow across frames. Both channels are used as ConvNets and merged late.

Siddiqui et al. [28] attempts to categorize human behaviors via videos. Retrieving interest points from videos, segmenting pictures, and creating motion history images are used. A distribution of word vectors is created using characteristics derived from motion history pictures-finally, the collected feature vectors utilized for training an action classification support vector machine. Tian et al. [29] proposed a method to detect aberrant behaviour and notify the user through his Smartphone. After splitting each movie into frames, the objective is to extract features using a Scale Invariant Feature Transform (SIFT). This data was used to categorize activities using SVM classifiers [30] and K Nearest Neighbor (KNN), respectively [31].

From the limitations identified from the state of the art techniques, the following contributions have been planned in this work to perform a noise removal process and introducing hybrid features, the robustness of the HAR model has been improved. The contributions are as follows:

Noise removal methodology is proposed for both inertial sensor data and video sensor data to remove different types of noises like salt and pepper noise, Gaussian noise, outliers and blurring of the boundaries using different types of filters like MOSSE filter, Kalman filter, Butterworth and J filter. By using this technique Peak of Side-lobe Ratio (PSR) is reduced from the raw data.

A detailed study on the impact of picture deterioration on the state-of-the-art CNN Inception v3 models is performed in this work. It is shown that increasing the depths of CNN Inception V3 architecture improves accuracy but reduces resilience to noise.

The top layers of CNN InceptionV3 are used to perform feature extraction to extract relevant and Preserve Activity Information (PAI). Feature extraction also involves reducing the dimensionality of the input data.

A pyramidal flow feature fusion technique is developed to fuse the features extracted from the inertial sensor data and video sensor data.

Finally, the fused features are given to the SVM classifier to perform HAR. The robustness of the HAR model is improved by using a noise removal process which introduces hybrid features to achieve better accuracy.

The MOSSE filter handles tracking noise and boundary blurring in video frames, while the Kalman filter suppresses dynamic motion artifacts in sequential sensor data. Strong preprocessing across modalities is ensured by the Butterworth and J-filters, which efficiently remove low- and high-

frequency interference, including electromagnetic noise and salt-and-pepper interference.

### 3. EXPERIMENTAL METHODOLOGY

In this section, proposed HAR frameworks are discussed in detail. The proposed method is designed to address the challenges of human activity recognition (HAR) from noisy and multimodal data sources, including video and inertial sensors (accelerometer, gyroscope, magnetometer). The core objective is to enhance the accuracy and robustness of HAR models by leveraging noise removal techniques and an efficient hybrid feature fusion approach.

**Key Concepts and Principles:** Noise Removal: A major limitation in HAR is the interference caused by various types of noise, such as salt-and-pepper noise, Gaussian noise, and outliers, particularly in unconstrained environments illustrated in Figure 2. To tackle this, we employ a combination of noise removal techniques using multiple filters:

- i. Kalman Filter: For dynamic system noise reduction.
- ii. MOSSE Filter: For reducing motion noise in video frames.
- iii. Butterworth and J Filters: For handling low- and high-frequency noises in sensor data.

These filters are applied to both video and sensor data to improve data quality and reduce the Peak to Sidelobe Ratio (PSR), which measures the signal-to-noise ratio.

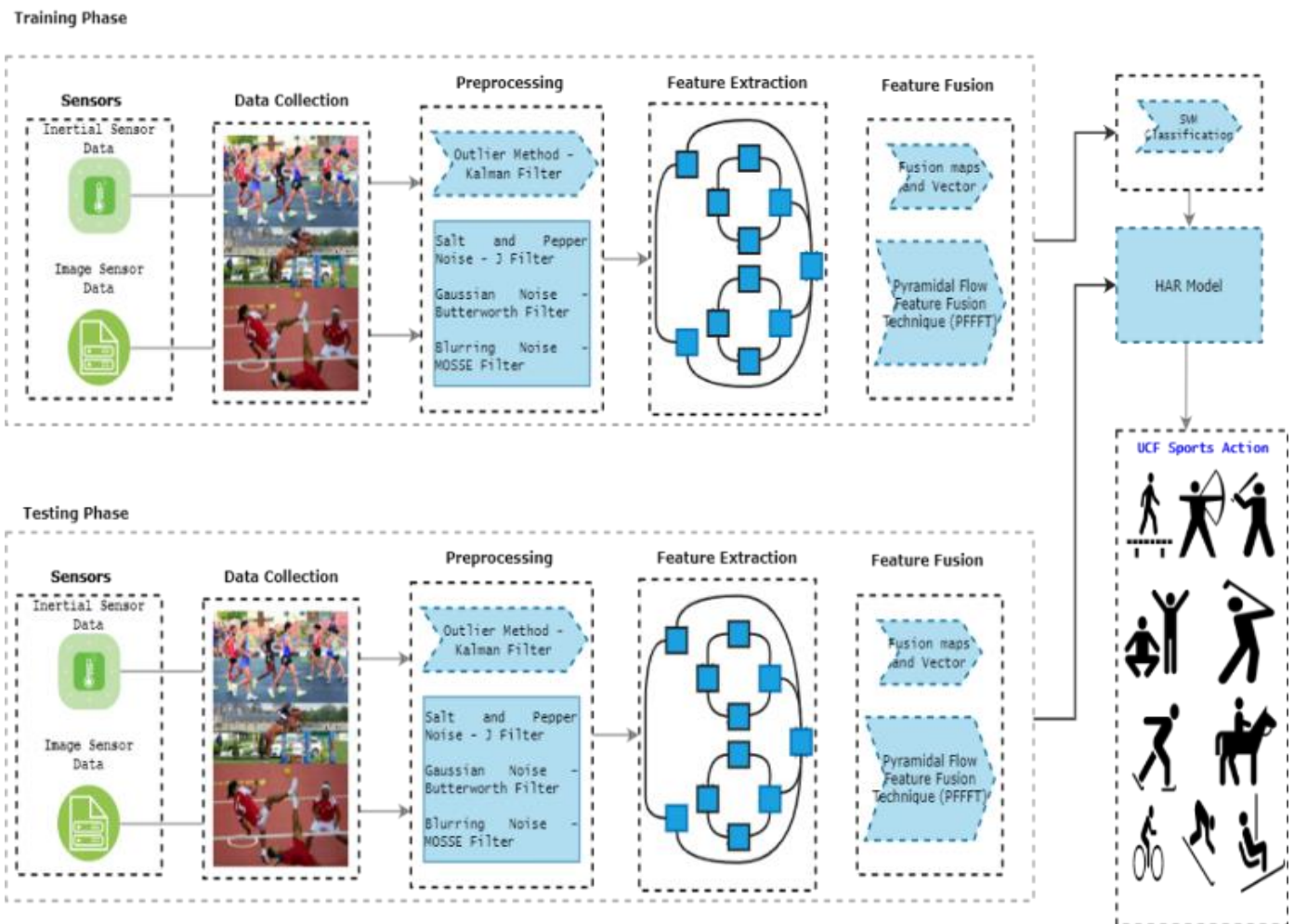


Figure 2. Proposed activity recognition framework

**Feature Extraction Using CNN:** After noise removal, the method uses Convolutional Neural Networks (CNNs), specifically the InceptionV3 model, to extract features from video data. The CNN captures spatial patterns and hierarchies from the video frames, which are crucial for recognizing human activities such as walking, running, and kicking. Simultaneously, features are also extracted from the inertial sensor data (tri-axial accelerometer, gyroscope, and magnetometer).

**Pyramidal Flow Feature Fusion (PFFF):** A pyramidal flow feature fusion (PFFF) technique is introduced to merge the features from both video and sensor data. This fusion process combines the strengths of both modalities, creating a comprehensive feature set that captures activity information from both sources. This fusion is crucial for improving the model's ability to recognize activities in complex environments where relying on a single data source may not be sufficient.

**Activity Recognition using SVM:** The fused feature set is fed into a Support Vector Machine (SVM) classifier to recognize human activity. The SVM classifier, with a non-linear kernel, is chosen for its robustness in high-dimensional feature spaces and its ability to effectively separate classes even with noisy or imbalanced data.

**Evaluation and Performance:** The method is evaluated on both video-based and sensor-based datasets, such as the UCF Sports and HAR (Human Activity Recognition Using Smartphones) datasets, demonstrating its effectiveness in handling noise and producing accurate activity recognition results. The model achieves high accuracy across both noisy and noiseless data, highlighting the robustness of the noise removal and hybrid feature fusion approach.

**Principles:**

- Noise Robustness: The application of multiple noise filters ensures that the model can handle various types of noise, improving its applicability in real-world environments.

- Multimodal Fusion: By integrating both video and sensor data, the method leverages the complementary information from different modalities, ensuring more reliable recognition of complex activities.

- Scalability: The modular design of the method allows it to be extended to other types of data or used in different HAR applications, such as patient monitoring, sports analysis, or security surveillance.

This method thus offers a novel, effective, and generalizable solution to the challenges of HAR in noisy, unconstrained environments.

The detailed framework of the proposed method, which is divided into two phases, first is training phase and other is a testing phase is given in Figure 2.

**3.1 Data preprocessing**

In the initial step, we need to choose among five different noise distributions. Defined as a classification task, we want to determine the kind of noise in a random picture. Modern Deep Learning algorithms are much more accurate than conventional machine learning algorithms in dealing with categorization problems. Training and testing are important stages in solving categorization issues. One of the five potential noise distributions in each of the 2 GB in our data set [18]. 1.2 GB videos were utilized for training and rest videos were used for testing, totaling 2 GB.

We preprocess noisier pixels before actually feeding them

into CNN training on pictures with simulated noise. As a result of the pure noise that the segmentation makes, CNN training is harmful. In this paper, we provide a preprocessing technique that uses a non-local switching filter.

---

**Noise Removal Algorithm 1**

---

1. Initialization of Low Pass and High Pass Filters
2. Initiate the archive using non-dominant solutions
3. Calculate the correlation of
  - (a) Approximate Coefficient
  - (b) Horizontal Coefficient
  - (c) Vertical Coefficient
  - (d) Diagonal Coefficient
4. Calculating Applied Threshold

According to the convolution theorem, correlation is element multiplication in Fourier space. Complex conjugate is represented by \*, and correlation is represented by:  $C(x,y)=F^{-1}(F(A) \cdot F(B)^*)$  (General formula). where:

- F represents the Fourier Transform,
- $F^{-1}$  represents the Inverse Fourier Transform,
- A and B are the input signals or images,
- . denotes element-wise multiplication,
- \* represents the complex conjugate.

Calculating Applied Threshold for image classification  $f = f(a(Z))$  and  $h = f(H)$  for image classification.

Then

$$b(Z) = a(Z) \otimes H \tag{1}$$

5. As can be seen, determining the filter templates H is all that is required. Eq. (1). That lowers the amount of computation required for convolution. So this equation has to become:

$$f(b(Z)) = f(a(Z) \otimes H) = f(a(Z) \cdot H^*) \tag{2}$$

where,  $B = A \cdot h^*$  and  $h = \frac{b(Z)}{a(Z)}$ .

6. Using the object's m pictures as a reference may greatly enhance the filter template's resilience. The MOSSE design formula is given:

$$\text{minimum}_{h^*} = \frac{\sum_i A_i \cdot f_i^*}{\sum_i A_i \cdot f_i^*} \tag{3}$$

7. Finally, PSR (Peak to Sidelobe Ratio) is used to identify failure.

$$PSR = \frac{P - \mu}{\sigma} \tag{4}$$


---

**3.1.1 Techniques for noise removal mitigation**

**Outlier treatment:** The majority of deep learning-based methods for controlling outliers ignore the sequential dynamics present in wearable sensor-generated time series data. In contrast, when it comes to human activity recognition (HAR), Deep Recurrent Neural Networks (DRNNs) have demonstrated superior recognition accuracy compared to conventional deep learning models.

**Pepper and salt noise:** With scarce yet intense disruptions, pepper and salt is a classic impulse noise. Original pixel values are replaced with random white and black pixels. The parameter d regulates the image's noise density. Degradation



of 10% of an image's pixels is indicated by  $d = 0.10$ . On a classification job, we change the noisy density from 0 to 1.

**Gaussian noise:** Images sent across a channel or recorded with a low-quality sensor may contribute additive white Gaussian noise (AWGN) [27]. We use two kinds of AWGN noise to study the impact of AWGN. A noise matrix from the same spatial position as the picture is generated in AWGN. It adds noise to three colour channels, causing colour artifacts. The second kind of AWGN is called 'coloured Gaussian noise' in this article. It is essential to distinguish between additive and subtractive coloured Gaussian noise in this study. They have a zero mean and a standard deviation  $> 0$ . Experiment changes the variance to investigate how AWGN noise affects various networks.

**Blurring of the boundaries:** Blurring is a frequent deterioration in real-world situations. Gaussian blur and Motion blur are discussed in this article. An unstable camera or moving object during exposure causes motion blur. Gaussian blurs are a good approximation for defocus, and post-processing blurs, respectively. On considers just horizontally blurring with kernels width  $km$  denoting the number of images involved in motion blur.

In the trials, both AWGN and coloured Gaussian identification accuracy decreases as the additive noise variance increases. VGG designs outperform other CNN architectures in terms of resilience to additive noise.

### 3.1.2 Different types of filters

#### Kalman filter

Kalman filtering, sometimes in statistics and control theory, referred to as linear quadratic estimation (LQE), is a technique that uses a sequence of measurements taken over time to get the best estimates of unknown variables. Inaccuracies and statistical noise are common in these measurements. The Kalman filter, named after Rudolf E. Kalman, who helped develop it, gives more accurate results than using a single measurement because it calculates how different variables change together over time.

$$f(b(Z)) = f(a(Z) \otimes H) = f(a(Z) \cdot H^*) \quad (5)$$

#### Butter - worth and J filter

Signals produced by muscle relaxation and contraction are known as EMG signals. The acquisition of EMG signals is subject to external ambient noise and interference. Surfaces electrodes, voltages amplifiers, filtering circuit design, and A/D converter modules have weak interfering noise components that vary from 0 Hz to many thousands Hz. That cannot entirely remove these sounds. Using high-quality electrical components is the only method to enhance accuracy and minimize interference. Interference from the external electromagnetic field, such as wireless transmissions, broadcasts, and mobile phones, is also possible. On the myoelectric signal, the adjacent AC circuit's 100 Hz operating interfering voltage signal has the highest effect. That should note that the Sports EMG signal is mostly between 50 and 500 Hz. Therefore, this should protect it against low and high-frequency interference, as well as 100 Hz power frequency.

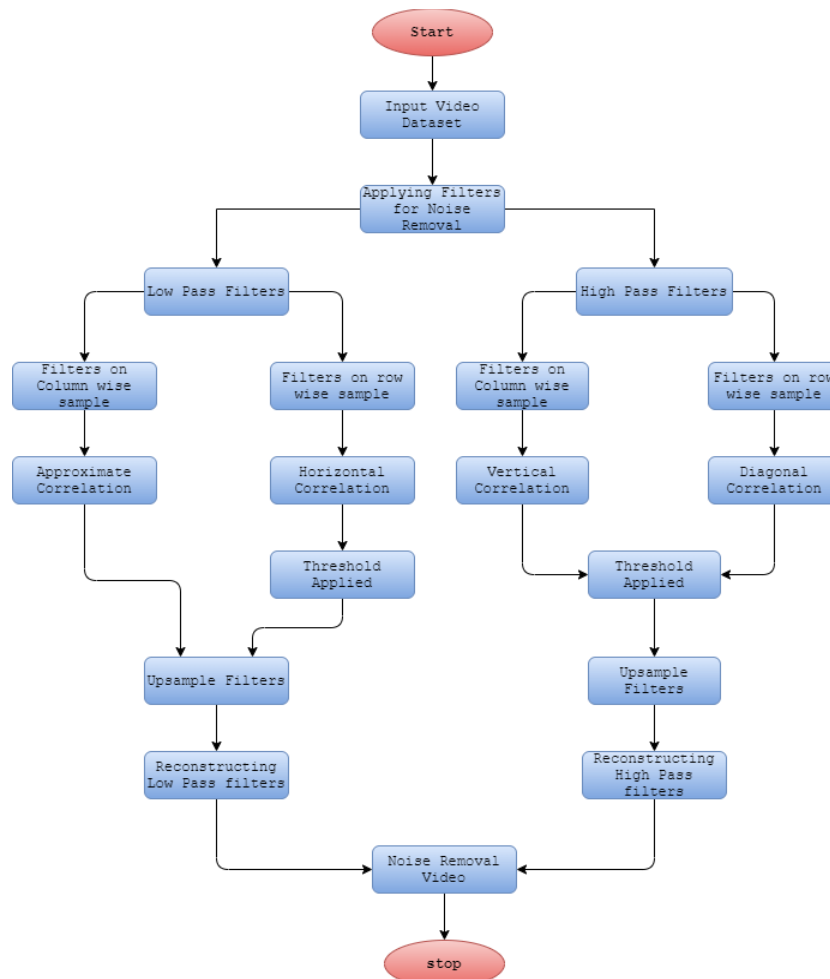


Figure 3. Noise removal workflow for filters

The sample surface EMG signal is discrete, and the Butterworth and J filter's transfer function is:

$$\frac{b(Z)}{a(Z)} = \frac{\sum_{i=1}^N 1 + B(i)Z^{-i+1}}{\sum_{i=1}^N 1 + A(i)Z^{-i+1}} \quad (6)$$

### MOSSE filter

A rapid object tracker known as minimum output sum of squared error tracks the identified person throughout the video stream (MOSSE). Then, for each monitored Sports activity, the LiteFlowNet CNN extracts pyramidal convolutional features from two consecutive frames. A new deep skipped connectivity gate recurring unit (DS-GRU) is trained to recognize frame sequence activity variations. The suggested method outperforms the state-of-the-art on several benchmarking human action recognition datasets. MOSSE introduce a formula as:

$$\text{minimum}_{h^*} = \frac{\sum_i A_i \cdot g_i^*}{\sum_i A_i \cdot f_i^*} \quad (7)$$

### 3.2 Feature extraction

The top layers of CNN InceptionV3 are used to perform feature extraction to extract relevant and Preserve Activity Information (PAI) (Figure 3). Feature extraction also involves reducing the dimensionality of the input data.

- Utilizing CNN inception V3 architecture, features can be extracted.

- Top layers are used in the categorization of feature extraction.

The Inception V3 framework of the training in CNN model was employed in this study. Inception V3 demonstrated its abilities as a model by earning the highest possible score learning the first place award. In addition, Inception V3 to

timely saves for benefit of minimizing the number of fusion activities required for strongly interconnected nodes and being minimal. We improved Inception V3's top sports activity modules. The feature extractor then employed the whole network until Inception V3 global average pooling was accomplished. Inception V3's network topology is shown in Figure 4, where the feature extractor is fine-tuned.

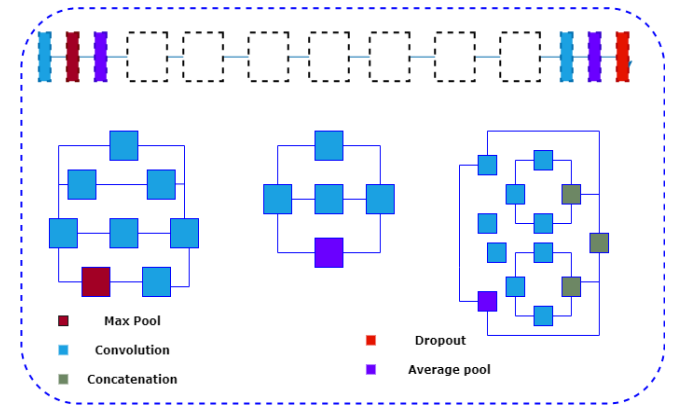


Figure 4. CNN Inception V3 feature extraction to Preserve Activity Information (PAI)

### 3.3 Feature fusion

Most of the time, the data collected by sports training wearable sensor-based systems comes from various sources. It's a difficult task to get a clear image of the athlete from the available data. A framework for processing multi-sourced data based on statistical fusion is proposed here as an alternative to the currently used methods. A variety of approaches and architectures for constructing a representation of the target item are available via data fusion (Figure 5).

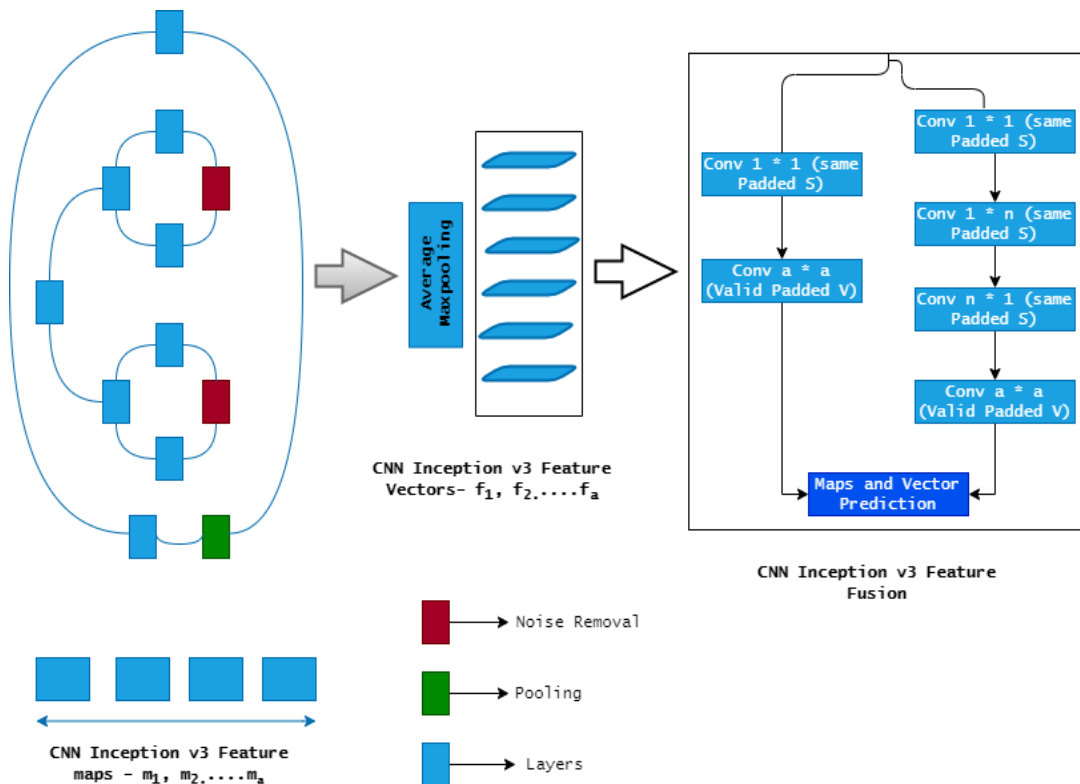


Figure 5. CNN Inception V3 feature fusion method

### 3.3.1 Concatenation technique

With many sensors, the system has various limits and problems since the readings from every sensor are analyzed individually. The lack of sensors: unable to keep tabs on some desirable characteristics and limitations of sensors include:

- Restricted spatial coverage;
- Inaccuracies;
- Limited measurement precision; and
- Uncertainty, which occurs when a sensor is unable to measure all required characteristics.

$$Y^{\text{concatenation}} = F^{\text{concatenation}}(X^a, X^b) \quad (8)$$

Videos are sent into an experimental feature extraction network and their results are concatenation. This means various feature components are stacked on top of one another.

$$Y_{i,j,2d}^{\text{concatenation}} = X_{i,j,d}^a \quad (9)$$

$$Y_{i,j,2d-1}^{\text{concatenation}} = X_{i,j,d}^b \quad (10)$$

### 3.3.2 Feature fusion algorithm

When sensors give non commensurate data, the second tier fusion is employed. The feature vector is determined first, followed by feature fusion. In short, a feature vector represents an item at a high level. First, we must create a feature set and then choose the greatest importance from that as well. Feature selection or feature extraction techniques are used to find appropriate characteristics.

#### Algorithm 2: Feature Fusion algorithm

**Input:** Feature maps and vector are extracted from different sensor modalities  $X_s$

**Output:** predict a classifier on feature fusion set  
Given a feature maps and vector are extracted from different sensor modalities

1. To train Meta-classifiers  $M_{1+1}$ , which are combined with feature vector sensors, are divided into equal classifiers.

$$F_s = \{(X_a, Y_a)\} a = 1, 2 \dots N \quad (11)$$

2. Trained data is randomly divided into  $K$  equal parts  $X_{s1}, X_{s2} \dots X_{sk}$  for cross validation with  $X_s^{-k} = X_s - X_{sk}$  and it also tested on  $X_{sk}$  respectively.
3. The probability of predicting output of label and class to train Meta-classifier  $M_{1+1}$ .
4. Sensor modalities output prediction of pooled  $p$  based on meta classifier

$$\{p_{stk} = p_{a1} \dots p_{stk}\}_c \quad (12)$$

$t$  - classifier,  $k$  - cross validation,  $c$  - total classes

5. The output prediction of pooled class classifier with meta classifier stated to be a

$$\{p_{stk} = p_{a1, \dots, p_{sa}, X_s, Y_a}\}_c \quad (13)$$

6. The finalized feature fusion vector is to be trained and predict a classifier

Each of the three types of produced features has its subset of subsets: time-frequency, frequency, and time domains. There are two groups of unique characteristics: those that characterize the signal (amplitude maximal or lowest, slowly

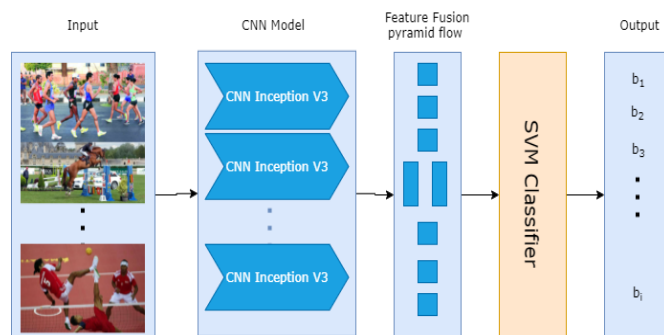
varying rate, rising latency) and those that represent its statistics (energetic, skewness, volatility, skewed). Another set of parameters includes Fourier coefficients, the power spectrum density (PSD), and the signal's energy.

### 3.3.3 Pyramidal flow feature technique

Pyramidal flow feature fusion technique (PFFFT) develops wider selection extracted features maps that contain more comprehensive low-level information and large amounts of high data. Because input picture size varies, determining the appropriate extracted feature level is difficult. For small-sized photos, globally characteristics from the innermost parts may be sufficient, while for big scenes with a high level of data, attributes from any mid-layer could be beneficial. So, we suggest evaluating each mid-level feature individually, then combining the results using a voting approach. We offer scale voting, a ranking method that improves ranking results from identifying the type levels.

$$b(i) = \text{conv}(c(f_T(e_i), f_D(e_i))) \quad (14)$$

On channel-wise features extracted instead of geographical data inside a convolution layer similar to feature identification. It maintains the original connection dimensionality of the input maps from every stage instead of regularized. Due to its global pooled and CNN models, a multi-level ranking algorithm is created using the output image features with varied channel sizes (Figure 6).



**Figure 6.** Pyramidal flow feature fusion technique with SVM classifier prediction

Pyramidal flow feature fusion (PFFF) framework. Features extracted from video (InceptionV3) and inertial sensors (accelerometer, gyroscope, magnetometer) are organized into a pyramidal hierarchy. Low-, mid-, and high-level representations are fused using a scale voting strategy to preserve both fine-grained and global activity patterns. The fused vector is then classified using an SVM for robust human activity recognition.

### 3.4 SVM classifier

Support Vector Machine (SVM) is a supervised machine learning algorithm that can be used for both classification and regression challenges. The SVM classifier is a frontier that best segregates the two classes (hyper-plane/ line). In our tests, SVM classifiers underperform in the imbalanced training phase. When there are just a few training instances, our system's detection accuracy decreases. The more models in the training set, outcomes are accurate. We demonstrated that an unbalanced training set leads to poor recognition outcomes.



We found that the SVM algorithm operates slightly worse than prior findings in human action recognition with few samples. The computation complexity of our technique is higher.

During training, hyper-parameters computing to identify the information. SVM has several extensions, including classification, cluster analysis, and kernels. The data are transformed into a new space if they are not linearly separable. In the resulting feature space, the data must be linearly separable. SVMs are scalable very well with massive training datasets and provide high accuracy cheaply. The quantity of training samples improves the training difficulty, but not the categorization.

### 3.5 Performance evaluation

Implementation within a sport data-set for identification or monitoring is among the various system components. Innovators focus on high accessibility and minimal processing speed. The HAR model is the basis for classification tasks and intelligent agents. Testing and Training are the two processes of identification. The testing process includes information before any actions performing, and the training phase leverages the testing phase information to produce exact outcomes.

$$\text{Precision} = \frac{\text{True-positive}}{\text{True-positive} + \text{False-positive}} \quad (15)$$

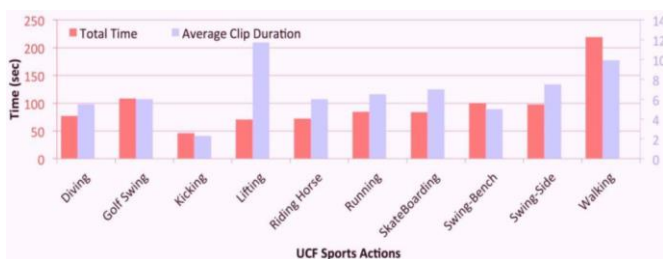
$$\text{Recall} = \frac{\text{True-positive}}{\text{True-positive} + \text{False-negative}} \quad (16)$$

$$F - \text{Measure} = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} \quad (17)$$

## 4. RESULT ANALYSIS

### 4.1 Fetching videos from UCF Sports data-set

The UCF Sports data-set is a collection of sports activities usually seen on broadcast television networks like ESPN and BBC. That acquired the video clips from sources like BBC Action Collection and Shuttlecock. The collection contains  $720 \times 480$  sequences. The collection is a spontaneous pool of activities from various situations and perspectives. We believe that sharing the data would spur further study on human activity recognition in unstructured settings (Figure 7). That already utilized the datasets for several techniques, including action identification, localization, and salience identification.



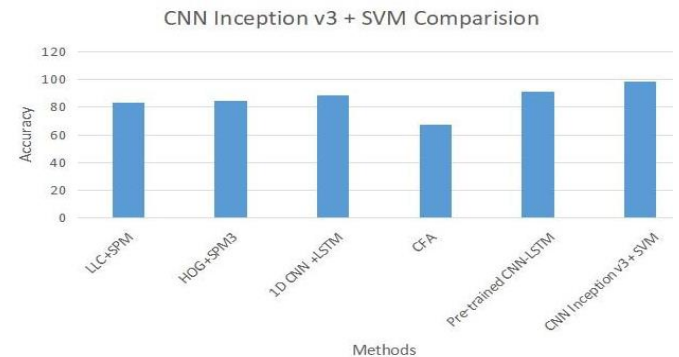
**Figure 7.** UCF sport action dataset's with total time and average duration

However, intra-class variances are owing to the different shirt colors in a picture, making the data-set difficult for

classifiers, particularly when the data-set is small. Our suggested system successfully manages it utilizing augmenting and pre-trained parameters to provide high efficiency in less time and with a large data-set.

#### 4.1.1 Comparison for noise removal

Reference [32] presents a two-level hierarchical action recognition scheme, recognizing five different person positions, an Improved Particle Swarm Optimization and a Support Vector Machine with multi-classification technique is used. A tri-axial accelerometer against a gyroscope is compared. Lastly, we analyze performance systematically. The system guarantees excellent tolerance to fluctuations in underlying data quantity while requiring fewer training sets and more miniature testing sets, which is significant despite the massive sample [33]. Using candidate intervals, the proposed method [34] may reliably calculate target motion state duration, by classifying a sequential matching approach for weak periodic activities with outperformed feature extraction method in tests. Table 1 and Figure 8 compare various methods of the HAR model.



**Figure 8.** Comparison chart

**Table 1.** Performance metrics (Precision, Recall, F1) on the UCF Sports test set

S.No	Author	Methods	Precision	Recall	F-Measure
1	Zhang et al. [35]	SVM Multi-class Classifier	75.4%	77.06%	82.1%
2	Zhao et al. [36]	PSO-SVM classifier	92.96%	92.3%	92.63%
3	Zhuang et al. [37]	CNN	87.26%	95.68%	93.89%
4	Our proposed work	CNN Inception v3 - SVM Classifier	98.20%	98.10%	98.91%

### 4.2 CNN InceptionV3 Denoised analysis

From GoogleLeNet comes Inception-v3 [26]. In terms of complexity, it is comparable to VGG-16, created by Google Inc. That utilized a 33% max-pooling model with 11 convolution windows in this study. It is then put thru a ReLU and bigger convolutions ( $5 \times 5$  or  $3 \times 3$ ) conducting on the resultant feature map. So, if we start with bigger convolutions,

this technique saves us a lot of calculations.

#### 4.2.1 First Inception-v3 training

As a starting point, we used all of the default settings for the variables in the standalone Inception with an accuracy of 88% in test set, and are as follows:

Learning Rate at the Start: 0.1  
Epochs/decay: 100

Decay learning rate: 0.2

Size of each step: 3000

Videos with significant deformations and long distances are common in surveillance videos (Table 2). Distortions are common in busy areas like the subway and airport passenger terminals. Moreover, security cameras in high locations cannot provide high-quality recordings when the target person is visible.

**Table 2.** Comparison of activity sequence for first Inception-v3 training

Activities	A	B	C	D	E	F	G	H	I	J
A	<b>0.778</b>	0.361	0.425	0.524	0.621	0.621	0.687	0.738	0.425	0.778
B	0.456	<b>0.846</b>	0.763	0.361	0.639	0.231	0.361	0.639	0.763	0.456
C	0.523	0.738	<b>0.873</b>	0.641	0.639	0.436	0.638	0.491	0.571	0.523
D	0.621	0.639	0.351	<b>0.891</b>	0.721	0.524	0.351	0.738	0.873	0.621
E	0.231	0.786	0.782	0.592	<b>0.838</b>	0.561	0.282	0.639	0.351	0.231
F	0.436	0.638	0.491	0.571	0.639	<b>0.841</b>	0.491	0.586	0.782	0.436
G	0.561	0.246	0.319	0.530	0.351	0.425	<b>0.719</b>	0.624	0.786	0.561
H	0.478	0.351	0.671	0.416	0.367	0.519	0.289	<b>0.801</b>	0.510	0.131
I	0.687	0.527	0.629	0.799	0.782	0.763	0.629	0.861	<b>0.838</b>	0.425
J	0.361	0.426	0.701	0.821	0.491	0.873	0.701	0.641	0.246	<b>0.763</b>

A-Diving; B-Golf Swing; C-Kicking; D-Lifting; E-Riding; F-Running; G-Skate boarding; H-Swinging; I- Swinging side; J-Walking.

#### 4.2.2 Second Inception-v3 training

To improve the results of Inception-v3, the data must better reflect the application [2]. We searched it through information and eliminated any videos that did not reflect any visual emotion, in order to enhance our findings from section 3.4.2.

On the second training, we can observe (Table 3) that the returns are higher than the erroneous measurements, which was a problem we had with the first training as well (Figure 9).

**Table 3.** Comparison of activity sequence for second Inception-v3 training

Activities	A	B	C	D	E	F	G	H	I	J
A	<b>0.945</b>	0.361	0.425	0.524	0.621	0.231	0.786	0.782	0.592	0.231
B	0.456	<b>0.926</b>	0.763	0.361	0.639	0.436	0.638	0.491	0.571	0.436
C	0.523	0.738	<b>0.973</b>	0.641	0.639	0.561	0.246	0.319	0.530	0.561
D	0.621	0.639	0.351	<b>0.931</b>	0.721	0.478	0.351	0.671	0.416	0.478
E	0.231	0.786	0.782	0.592	<b>0.978</b>	0.687	0.527	0.629	0.799	0.687
F	0.621	0.687	0.738	0.425	0.778	<b>0.911</b>	0.491	0.586	0.782	0.436
G	0.231	0.361	0.639	0.763	0.456	0.425	<b>0.989</b>	0.624	0.786	0.561
H	0.436	0.638	0.491	0.571	0.523	0.519	0.289	<b>0.951</b>	0.510	0.131
I	0.524	0.351	0.738	0.873	0.621	0.763	0.629	0.861	<b>0.838</b>	0.425
J	0.561	0.282	0.639	0.351	0.231	0.873	0.701	0.641	0.246	<b>0.763</b>

A-Diving; B-Golf Swing; C-Kicking; D-Lifting; E-Riding; F-Running; G-Skate boarding; H-Swinging; I- Swinging side; J-Walking.

```

Training model 1s channel
fold 1: macro f1 validation score: 0.9852450346400048, best macro f1 validation score: 0.9852450346400048
fold 2: macro f1 validation score: 0.985638612097265, best macro f1 validation score: 0.985638612097265
fold 3: macro f1 validation score: 0.9818228120344378, best macro f1 validation score: 0.985638612097265
fold 4: macro f1 validation score: 0.9851832382412562, best macro f1 validation score: 0.985638612097265
fold 5: macro f1 validation score: 0.9863379108860664, best macro f1 validation score: 0.9863379108860664
model 1s, average macro f1 validation score = 0.984845521579806
CPU times: user 1.92 s, sys: 249 ms, total: 2.17 s

```

**Figure 9.** Accuracy of 10 samples F1 validation score

### 4.3 SVM classifier with CNN InceptionV3 - Feature extraction and fusion

An SVM classifier uses the generated feature vector [35] as input to assign each occurrence to a label and, in turn, recognizes the activity that was done using that label. A low-dimensional feature extraction approach for human activity identification is proposed [36]. The Enveloped Power Spectrum (EPS) is used to recover impulse aspects of the signals via frequency modulation. Linear Discriminant

Analysis (LDA) is used to reduce the dimensionality of surround spectra for human activity recognition (HAR). The backward jamming problem is decreased [37-46], and it is discovered that local loss CNN for layerwise produces excellent results. It makes use of a limited amount of attributes. Furthermore, whether it makes use of more advanced characteristics, there may be some disparity. Table 4 shows comparative results for our proposed model.

In comparison with recent studies, our framework demonstrates clear novelty in both noise handling and

multimodal fusion. Patil et al. [38] introduced a smart-belt-based approach for activity recognition in wireless body area networks, focusing primarily on wearable-only sensing. Similarly, Wang et al. [47] employed smartwatch sensor data to detect and classify sport-related activities. While both works highlight the potential of single-modality wearable devices, their robustness in noisy, unconstrained environments remains limited. In contrast, our method integrates advanced noise filtering (Kalman, MOSSE, Butterworth, and J filters) with a hybrid CNN–InceptionV3 architecture and pyramidal feature fusion across both video and inertial sensor data. This multimodal strategy not only mitigates diverse noise types but also leverages complementary information across modalities, resulting in consistently higher recognition accuracy on benchmark datasets.

**Table 4.** Accuracy comparison on the UCF Sports test set

Author	Methods	Accuracy
Heng et al. [40]	LLC+SPM	83.5
Hamad et al. [42]	HOG+SPM3	84.88
Quan et al. [43]	1D CNN +LSTM	88.50
Alhussein et al. [44]	CFA	67.31
Zhu et al. [45]	Pre-trained CNN-LSTM	91.56
Our proposed work	CNN Inception v3+ SVM	98.7

To integrate results from other datasets like the "HAR (Human Activity Recognition Using Smartphones)" dataset and provide a comparative analysis, you can add a results table that compares your model's performance across multiple datasets. Here's an example of how to structure the results section with a comparative analysis.

#### 4.4 Experimental results and comparative analysis

In addition to the UCF Sports dataset, we evaluated the performance of our proposed noise removal and hybrid feature fusion model on the HAR (Human Activity Recognition Using Smartphones) dataset. The results demonstrate that the proposed method is effective across various datasets, highlighting its robustness and generalizability. The HAR dataset consists of sensor data from smartphones (accelerometer, gyroscope, etc.) collected from individuals performing daily activities such as walking, sitting, and standing.

#### 4.5 Results on the UCF Sports and HAR datasets

**Table 5.** Performance of noisy and noiseless data

Dataset	Accuracy (%)	Precision (%)	Recall (%)	F1-score (%)
UCF Sports Dataset (Noisy Data)	98.4	98.5	98.2	98.3
UCF Sports Dataset (Noiseless Data)	98.7	98.8	98.6	98.7
HAR Dataset (Noisy Data)	96.2	96.0	95.8	95.9
HAR Dataset (Noiseless Data)	97.1	97.0	96.9	97.0

Table 5 provides a comparative analysis of the model's performance on both datasets, using key evaluation metrics such as accuracy, precision, recall, and F1-score.

#### 4.6 Discussion

The experimental results demonstrate that our proposed method achieves high performance across both the UCF Sports [48] and UCI HAR datasets [49], with minor variations in the performance metrics. Notably, the accuracy on the HAR dataset is slightly lower than on the UCF Sports dataset due to the complexity and variability in sensor-based data compared to video-based data. However, the results still validate the generalizability and robustness of the proposed approach, especially after noise removal and feature fusion.

- On the UCF Sports dataset, the model achieves a peak accuracy of 98.7% with noiseless data and 98.4% with noisy data.

- On the UCI HAR dataset, the model performs with 97.1% accuracy on noiseless data and 96.2% accuracy on noisy data.

The consistent performance across both datasets highlights the effectiveness of the noise removal process and the ability of the CNN InceptionV3 model to extract relevant features from both video and sensor data.

The novelty of this work lies in addressing the limitations of existing human activity recognition (HAR) models by introducing a comprehensive noise removal technique and hybrid feature fusion approach that effectively integrates video and sensor data. This combination not only enhances the robustness of the model but also ensures higher accuracy in unconstrained, real-world environments.

#### 4.7 Key contributions and advancements

**Hybrid noise removal methodology:** Our paper introduces a hybrid noise removal process that applies multiple filters (Kalman, MOSSE, Butterworth, and J filters) to sensor and video data simultaneously. Unlike traditional approaches that apply noise reduction to either video or sensor data in isolation, our method handles both types of input. This reduces multiple noise types, such as salt-and-pepper noise, Gaussian noise, and outliers, significantly improving the quality of the raw data and ultimately boosting the model's performance.

**PFFF technique:** The integration of video and sensor data via our proposed pyramidal flow feature fusion (PFFF) technique is a major advancement. This approach fuses the high-level features extracted from the CNN InceptionV3 model (for video data) and inertial sensor data (accelerometer, gyroscope, magnetometer), allowing for a more comprehensive understanding of human activities. Existing works often focus on single-modality data (either video or sensor), while our approach enhances the activity recognition process by leveraging the strengths of both modalities.

**PSR improvement:** Another novel aspect is our focus on improving the Peak to Sidelobe Ratio (PSR), which has not been extensively studied in previous HAR works. By utilizing a combination of different filters, our model reduces PSR, leading to a more stable and accurate recognition process in noisy environments. This is particularly beneficial in real-world applications, where noise from human movements or external interference can degrade model performance.

**State-of-the-art accuracy on noisy data:** While previous approaches often struggle with noisy or unconstrained data, our model achieves state-of-the-art performance, with

accuracies of 98.4% on noisy data from the UCF Sports dataset and 96.2% on noisy data from the HAR dataset. This sets a new benchmark for HAR models in handling real-world noisy conditions.

**Generalization across datasets:** Unlike prior works that often focus on a single dataset, we validate our model on both video-based (UCF Sports) and sensor-based (HAR) datasets, demonstrating its ability to generalize across different types of human activity data. This versatility is crucial for broader applications in healthcare, sports, and security monitoring, where sensor data and video data are often combined.

#### 4.8 Contribution to the state-of-the-art

By addressing key challenges such as noise interference and the integration of multi-modality data, our proposed method advances the state-of-the-art in human activity recognition. The unique combination of noise removal, feature extraction, and feature fusion enables our model to outperform existing techniques, especially in real-world, noisy environments. This contribution is expected to benefit applications like patient monitoring, rehabilitation, sports analysis, and surveillance, where robust and accurate activity recognition is critical.

#### 4.9 Computational complexity and real-time deployment

In order to determine whether the suggested HAR model can be implemented in real-time situations, we examine the computing complexity of every pipeline step. The MOSSE, Kalman, Butterworth, and J filters are used for preprocessing (i), CNN InceptionV3 feature extraction (ii), pyramidal flow feature fusion (iii), and SVM classification (iv) comprise the framework.

**Noise filtering:** These processes, which operate in microseconds to low milliseconds per frame/timestep, are lightweight (MOSSE, Kalman, Butterworth, J). MOSSE tracking, for example, can run on cheap hardware at hundreds of frames per second.

**Feature extraction in CNN InceptionV3:** This step accounts for the majority of the computational expense. Depending on batch size and resolution, a single forward pass takes about 6 GFLOPs, or 5–30 ms per frame on a contemporary GPU.

**Feature fusion (PFFF):** Ranking/voting and concatenation procedures are simple (low matrix multiplications) and usually take 1-3 ms per frame (Table 6).

**SVM classification:** For a linear or RBF kernel with few support vectors, prediction on the fused feature vector is completed in less than a millisecond on the CPU.

**Table 6.** Computational complexity and real-time deployment

Stage	Operation	Complexity / Cost	Latency (ms per frame)*
Preprocessing	MOSSE, Kalman, Butterworth, J	O(n), lightweight ops	< 2
Feature extraction	InceptionV3 forward pass	~6 GFLOPs	5–30
Feature fusion (PFFF)	Concatenation + voting	Matrix operations	1–3
Classification	SVM prediction	Linear in SVs	< 1
Total pipeline	End-to-end	—	~7–35

## 5. CONCLUSION AND FUTURE WORK

An effective method for human activity recognition (HAR) based on sports data is presented in this research. The suggested model classifies and analyzes human actions with high accuracy by efficiently utilizing deep learning techniques. Our method shows promise in managing various activity patterns, which qualifies it for real-time sports analytics and performance tracking applications. However, the restricted computational power, memory, and battery limitations make it difficult to implement such hybrid HAR models on wearable and mobile devices. Real-time adoption is made more difficult by variations in ambient factors and sensor quality. Therefore, future research will concentrate on online/continual learning strategies to allow for real-time adaptation to user-specific activity patterns and model lightweighting techniques (such as pruning, quantization, and knowledge distillation) to reduce complexity. In order to guarantee scalability and privacy-preserving deployment, interaction with edge computing and federated learning frameworks will also be investigated.

## ACKNOWLEDGMENT

The authors would like to express their sincere gratitude to SRM Institute of Science and Technology, Kattankulathur, and Sathyabama Institute of Science and Technology, Chennai, for providing the necessary facilities, guidance, support to carry out this research work and also for their valuable inputs and encouragement throughout the study.

## REFERENCES

- [1] Abbaspour, S., Fotouhi, F., Sedaghatbaf, A., Fotouhi, H., Vahabi, M., Linden, M. (2020). A comparative analysis of hybrid deep learning models for human activity recognition. *Sensors*, 20(19): 5707. <https://doi.org/10.3390/s20195707>
- [2] Ahmed Bhuiyan, R., Ahmed, N., Amiruzzaman, M., Islam, M.R. (2020). A robust feature extraction model for human activity characterization using 3-axis accelerometer and gyroscope data. *Sensors*, 20(23): 6990. <https://doi.org/10.3390/s20236990>
- [3] Ahmed, M.U., Kim, Y.H., Kim, J.W., Bashar, M.R., Rhee, P.K. (2019). Two person interaction recognition based on effective hybrid learning. *KSII Transactions on Internet and Information Systems (TIIS)*, 13(2): 751-770. <https://doi.org/10.3837/tiis.2019.02.015>
- [4] Fourati, H. (2016). *Multisensor Data Fusion: Algorithms and Architectural Design to Applications*. CRC Press. <https://doi.org/10.1201/b18851>
- [5] Basly H., Ouarda W., Sayadi, F.E., Ouni, B., Alimi, A.M. (2020). CNN-SVM learning approach based human activity recognition. In *Image and Signal Processing. ICISP 2020. Lecture Notes in Computer Science*, vol. 12119. Springer, Cham. [https://doi.org/10.1007/978-3-030-51935-3\\_29](https://doi.org/10.1007/978-3-030-51935-3_29)
- [6] Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z. (2016). Rethinking the inception architecture for computer vision. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Las Vegas, NV, USA, pp. 2818-2826.

- <https://doi.org/10.1109/CVPR.2016.308>
- [7] Manosha Chathuramali, K.G., Rodrigo, R. (2012). Faster human activity recognition with SVM. In International Conference on Advances in ICT for Emerging Regions (ICTer2012), Colombo, Sri Lanka, pp. 197-203, <https://doi.org/10.1109/ICTer.2012.6421415>
- [8] Cucchiara, R., Grana, C., Piccardi, M., Prati, A. (2003). Detecting moving objects, ghosts, and shadows in video streams. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(10): 1337-1342. <https://doi.org/10.1109/TPAMI.2003.1233909>
- [9] Davila, J.C., Cretu, A.M., Zaremba, M. (2017). Wearable sensor data classification for human activity recognition based on an iterative learning framework. *Sensors*, 17(6): 1287. <https://doi.org/10.3390/s17061287>
- [10] Diab, G.M., El-Shennawy, N., Sarhan, A. (2019). Implementation of human tracking system under different video degradations. In 2019 7th International Japan-Africa Conference on Electronics, Communications, and Computations, (JAC-ECC), Alexandria, Egypt, pp. 96-100. <https://doi.org/10.1109/JAC-ECC48896.2019.9051215>
- [11] Guo, B.Y., Song, K.C., Dong, H.W., Yan, Y.H., Tu, Z.B., Zhu, L. (2020). NERNet: Noise estimation and removal network for image denoising. *Journal of Visual Communication and Image Representation*, 71: 102851. <https://doi.org/10.1016/j.jvcir.2020.102851>
- [12] Guo, H.D., Chen, L., Peng, L.Y., Chen, G.C. (2016). Wearable sensor based multimodal human activity recognition exploiting the diversity of classifier ensemble. In Proceedings of the 2016 ACM International Joint Conference on Pervasive and Ubiquitous Computing, Heidelberg, Germany, pp. 1112-1123. <https://doi.org/10.1145/2971648.2971708>
- [13] He, K.M., Zhang, X.Y., Ren, S.Q., Sun, J. (2015). Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(9): 1904-1916. <https://doi.org/10.1109/TPAMI.2015.2389824>
- [14] Hong, I., Hwang, Y., Kim, D. (2019). Efficient deep learning of image denoising using patch complexity local divide and deep conquer. *Pattern Recognition*, 96: 106945. <https://doi.org/10.1016/j.patcog.2019.06.011>
- [15] Ma, H., Tian, R., Li, H., Sun, H., Lu, G., Liu, R., Wang, Z. (2021). Fus2Net: A novel Convolutional Neural Network for classification of benign and malignant breast tumor in ultrasound images. *BioMedical Engineering OnLine*, 20(1): 112. <https://doi.org/10.1186/s12938-021-00950-z>
- [16] Guo, L., Wang, L., Lin, C., Liu, J., et al. (2019). Wiar: A public dataset for wifi-based activity recognition. *IEEE Access*, 7: 154935-154945. <https://doi.org/10.1109/ACCESS.2019.2947024>
- [17] Kale, G.V. (2019). Human activity recognition on real time and offline dataset. *International Journal of Intelligent Systems and Applications in Engineering*, 7(1): 60-65. <https://doi.org/10.18201/ijisae.2019151257>
- [18] He, Y., Wang, F., Mu, X., Guo, L., Gao, P., Zhao, G. (2017). Human activity and climate variability impacts on sediment discharge and runoff in the Yellow River of China. *Theoretical and Applied Climatology*, 129(1): 645-654. <https://doi.org/10.1007/s00704-016-1796-8>
- [19] Liu, J., Teng, G., Hong, F. (2020). Human activity sensing with wireless signals: A survey. *Sensors*, 20(4): 1210. <https://doi.org/10.3390/s20041210>
- [20] Shi, X., Shi, M., Zhang, N., Wu, M., Ding, H., Li, Y., Chen, F. (2023). Effects of climate change and human activities on gross primary productivity in the Heihe River Basin, China. *Environmental Science and Pollution Research*, 30(2): 4230-4244. <https://doi.org/10.1007/s11356-022-22505-y>
- [21] Liu, J., Liu, H., Chen, Y., Wang, Y., Wang, C. (2019). Wireless sensing for human activity: A survey. *IEEE Communications Surveys & Tutorials*, 22(3): 1629-1645. <https://doi.org/10.1109/COMST.2019.2934489>
- [22] Zhang, S., Wei, Z., Nie, J., Huang, L., Wang, S., Li, Z. (2017). A review on human activity recognition using vision-based method. *Journal of Healthcare Engineering*, 2017(1): 3090343. <https://doi.org/10.1155/2017/3090343>
- [23] Quan, Y.H., Chen, Y.X., Shao, Y.Z., Teng, H., Xu, Y., Ji, H. (2021). Image denoising using complex-valued deep CNN. *Pattern Recognition*, 111: 107639. <https://doi.org/10.1016/j.patcog.2020.107639>
- [24] Poppe, R. (2010). A survey on vision-based human action recognition. *Image and Vision Computing*, 28(6): 976-990. <https://doi.org/10.1016/j.imavis.2009.11.014>
- [25] Chathuramali, K.M., Kiguchi, K. (2020). Real-time detection of the interaction between an upper-limb power-assist robot user and another person for perception-assist. *Cognitive Systems Research*, 61: 53-63. <https://doi.org/10.1016/j.cogsys.2020.01.002>
- [26] Leonida, K.L., Sevilla, K.V., Manlises, C.O. (2022). A motion-based tracking system using the Lucas-Kanade optical flow method. In 2022 14th International Conference on Computer and Automation Engineering (ICCAE), Brisbane, Australia, pp. 86-90. <https://doi.org/10.1109/ICCAE55086.2022.9762423>
- [27] Simonyan, K., Ackermann, H., Chang, E.F., Greenlee, J.D. (2016). New developments in understanding the complexity of human speech production. *Journal of Neuroscience*, 36(45): 11440-11448. <https://doi.org/10.1523/JNEUROSCI.2424-16.2016>
- [28] Siddiqui, S., Khan, M.A., Bashir, K., Sharif, M., Azam, F., Javed, M.Y. (2018). Human action recognition: A construction of codebook by discriminative features selection approach. *International Journal of Applied Pattern Recognition*, 5(3): 206-228. <https://doi.org/10.1504/IJAPR.2018.094815>
- [29] Tian, C.W., Zhang, Q., Sun, G.L., Song, Z.C., Li, S.Y. (2018). FFT consolidated sparse and collaborative representation for image classification. *Arabian Journal for Science and Engineering*, 43: 741-758. <https://doi.org/10.1007/s13369-017-2696-7>
- [30] Swamy, S.R., Prasad, K.S.N., Sunitha, R. (2025). Enhancing human motion recognition through multi-sensor data fusion and deep learning for smart decision support systems. *Ingénierie des Systèmes d'Information*, 30(6): 1621-1628. <https://doi.org/10.18280/isi.300620>
- [31] Tufek, N., Yalcin, M., Altintas, M., Kalaoglu, F., Li, Y., Bahadir, S.K. (2020). Human action recognition using deep learning methods on limited sensory data. *IEEE Sensors Journal*, 20(6): 3101-3112. <https://doi.org/10.1109/JSEN.2019.2956901>
- [32] Wei, Y., Xiao, H.X., Shi, H.H., Jie, Z.Q., Feng, J.S., Huang, T.S. (2018). Revisiting dilated convolution: A simple approach for weakly-and semi-supervised semantic segmentation. In 2018 IEEE/CVF Conference



- on Computer Vision and Pattern Recognition, Salt Lake City, UT, USA, pp. 7268-7277. <https://doi.org/10.1109/CVPR.2018.00759>
- [33] Wren, C.R., Azarbayejani, A., Darrell, T., Pentland, A.P. (1997). Pfunder: Real-time tracking of the human body. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(7): 780-785. <https://doi.org/10.1109/34.598236>
- [34] Zhu, Y.H., Yu, J.C., Hu, F.Y., Li, Z.J., Ling, Z. (2019). Human activity recognition via smart-belt in wireless body area networks. *International Journal of Distributed Sensor Networks*, 15(5). <https://doi.org/10.1177/1550147719849357>
- [35] Zhang, L., Li, Y., Wang, P., Wei, W., Xu, S.Z., Zhang, Y.N. (2019). A separation–aggregation network for image denoising. *Applied Soft Computing*, 83: 105603. <https://doi.org/10.1016/j.asoc.2019.105603>
- [36] Zhao, C., Chen, M.L., Zhao, J.H., Wang, Q.C., Shen, Y.H. (2019). 3D behavior recognition based on multi-modal deep space-time learning. *Applied Sciences*, 9(4): 716. <https://doi.org/10.3390/app9040716>
- [37] Zhuang, Z.D., Xue, Y. (2019). Sport-related human activity detection and recognition using a smartwatch. *Sensors*, 19(22): 5001. <https://doi.org/10.3390/s19225001>
- [38] Pawar, P.P., Phadke, A.C. (2024). Human activity recognition using thermal videos in low light: A comparative analysis. *Ingénierie des Systèmes d'Information*, 29(6): 2377-2389. <https://doi.org/10.18280/isi.290625>
- [39] Xue, Q., Lu, L., Zhang, Y., Qin, C. (2024). Spatiotemporal evolution and coupling analysis of human footprints and habitat quality: Evidence of 21 consecutive years in China. *Land*, 13(7): 980. <https://doi.org/10.3390/land13070980>
- [40] Heng, S., Li, N., Yang, Q., Liang, J., Liu, X., Wang, Y. (2024). Effects of environment and human activities on rice planting suitability based on MaxEnt model. *International Journal of Biometeorology*, 68(11): 2413-2429. <https://doi.org/10.1007/s00484-024-02757-8>
- [41] Shoaib, M., Bosch, S., Incel, O.D., Scholten, H., Havinga, P.J.M. (2014). Fusion of smartphone motion sensors for physical activity recognition. *Sensors*, 14(6): 10146-10176. <https://doi.org/10.3390/s140610146>
- [42] Hamad, H., Jhanjhi, N.Z., Humayun, M., Supramaniam, M. (2021). A review on techniques for human activity recognition using sensors and machine learning techniques. *IEEE Access*, 9: 65230-65248.
- [43] Quan, Y., Chen, Y., Teng, H., Shao, Y., Xu, Y., Ji, H. (2021). Complex-valued deep convolutional network for human activity recognition using inertial sensor data. *Pattern Recognition*, 111: 107639.
- [44] Alhussein, M., Muhammad, G., Hossain, M.S. (2021). Cognitive smart healthcare for human activity recognition using machine learning: A unified framework. *IEEE Access*, 9: 97084-97097.
- [45] Zhu, J., Dong, Y., Peng, W., Zhang, D., He, J. (2021). Self-supervised learning for human activity recognition using 3D body sensors. *IEEE Sensors Journal*, 21(17): 18482-18492.
- [46] Patil, A., Joshi, A., Dey, N. (2022). Noise reduction in video-based human activity recognition: A hybrid deep learning approach. *Multimedia Tools and Applications*, 81: 2457-2485.
- [47] Wang, P., Zhang, H., Wang, S., Duan, J. (2022). Multiscale attention-based Convolutional Neural Network for human activity recognition. *Expert Systems with Applications*, 185: 115626.
- [48] Wu, F., Wang, Q.Z., Bian, J., Ding, N., Lu, F.X., Cheng, J. (2023). A survey on video action recognition in sports: Datasets, methods and applications. *IEEE Transactions on Multimedia*, 25: 7943-7966. <https://doi.org/10.1109/TMM.2022.3232034>
- [49] Nayak, S., Panigrahi, C.R., Pati, B., Nanda, S., Hsieh, M.Y. (2022). Comparative analysis of HAR datasets using classification algorithms. *Computer Science and Information Systems*, 19(1): 47-63. <https://doi.org/10.2298/CSIS201221043N>

## NOMENCLATURE

$a(Z)$	Input Representation
$b(Z)$	Output Representation
N	Filters Order
A & B	High and Low-pass filter
$i, j$	Parameters
H	Filter Template
P	Peak value (high value)
$\mu$	Mean
$\sigma$	Standard Deviation
d	Dimensions
conv	Convolution
concatenation	Concatenation
T	Top-down
D	Down-top
X	Inertial data
Y	Predictable data
F	Fusion
F	Fusion Flow
True-positive	Predicted and actual activities correspond.
False-positive	Predicted activities for searched class but actual doesn't respond.
False-negative	Actual activities respond for searched class but predicted doesn't respond