

# Human Action Recognition using Long Short-Term Memory and Convolutional Neural Network Model

Shreyas Pagare, Rakesh Kumar



**Abstract:** Human Action Recognition (HAR) is the difficulty of quickly identifying strenuous exercise performed by people. It is feasible to sample some measures of a body's tangential acceleration and speed using inertial sensors and exercise them only to learn model skills of incorrectly categorizing behavior into the relevant categories. In detecting human activities, the use of detectors in personal and portable devices has increased to better understand and anticipate human behavior. Many specialists are working toward developing a classification that can distinguish between a user's behavior and uncooked data while utilizing as few nerves as possible. A Long-term Recurrent Convolutional Network (LRCN) is proposed as a comprehensive human action recognition system based on deep neural networks in this paper.

**Keywords:** Human Action Recognition, Convolutional Neural Network, Long Short-Term Memory, Long Short-Term Memory

## I. INTRODUCTION

Human action recognition (HAR) is an imperative aspect of today's generation since it can analyse novel information about human actions from raw data. As a result of the rise of interpersonal communication applications, HAR innovation appears to have been a key study field both locally and internationally. Individuals can categorize any type of urban transportation and gather the knowledge that the body requires to function effectively to convey by extracting information from ordinary items, setting the framework for future applications [1].

Observing human behavior is essential in interpersonal interactions and intimate communication. It's tough to retrieve because this includes details about something like a specific individual, attitude and mental factors. In the research domains of computational intelligence, one of the really important objects of analysis is indeed the individual ability to perceive another person's behavior. A fundamental element of several consumer items was action recognition [2]. For example, game systems like the Nintendo Wii and Microsoft Kinect rely on gesture acknowledgment or even full-body arrangements to radically modify the game experience. While originally created for the entertainment sector, these systems have found new applications in fields like as personal fitness training and rehabilitation, as well as driving new action recognition research.

Manuscript received on 13 July 2023 | Revised Manuscript received on 10 May 2024 | Manuscript Accepted on 15 May 2024 | Manuscript published on 30 June 2024.

\*Correspondence Author(s)

**Shreyas Pagare\***, Research Scholar, Department of Computer Science & Engineering, RNTU University, Bhopal (M.P), India. E-mail: shreyas\_au211443@aisectuniversity.ac.in, ORCID ID: 0009-0001-8147-2389

**Dr. Rakesh Kumar**, Research Guide, Department of Computer Science & Engineering, RNTU University, Bhopal (M.P), India. E-mail: rakeshmittan@gmail.com

© The Authors. Published by Blue Eyes Intelligence Engineering and Sciences Publication (BEIESP). This is an open access article under the CC-BY-NC-ND license <http://creativecommons.org/licenses/by-nc-nd/4.0/>

Finally, some sporting goods, such as the Philips Direct Life or Nike+ running shoes, contain motion sensors and provide feedback to both amateur and professional athletes.

All of these instances demonstrate the importance of recognizing human activities in both academia and industry. Developing HAR systems that match application and user needs remains a difficult issue, despite significant progress in inferring activities from on-body encoders and prototyping and deploying action recognition systems. This is true even if HAR approaches that worked well for one recognition challenges are applied to a different problem domain.

### A. Machine Learning

A Machine Learning algorithm is a sequence of instructions and empirical methods for retrieving relevant insights from respondents and internal restructuring from it. It is the overarching theme of a Machine Learning model. A model is the most essential subset of machine learning. A Machine Learning Algorithm is employed as a model. In order to achieve the appropriate result, an approach integrates every one of the assumptions that a model is poised to deliver relying on the supplied information. A Machine Learning procedure occurs with massive quantities of data being supplied to the device and then used to train an algorithm to discover hidden actionable insight leveraging this fact. These observations then are utilized to make a Training Data that addresses the problems using a method [3].

### B. Deep Learning

Deep Learning is a machine learning subfield that understands to reflect real life as a recursive hierarchical order of principles, with each separate dimension in connection to smaller units and many more conceptual representations deduced in aspects of less intangible ones. It achieved excellent scalability and performance by learning to reflect reality as a recursive hierarchical system of notions. Deep learning networks learn by detecting complicated patterns from the data they receive. The systems can establish numerous areas of granularity to make sense of the data by establishing simulation tools that are composed of many convolutional filters [4].

A deep neural network, a category of deep learning model, could be learned using a substantial array (in the millions) of imagery, including those depicting living beings. Essentially, this form of neuron adapts from the cells in the visuals it obtains. It can distinguish and label collections of data that resemble distinct qualities, including such hand, head, and ears, which confirm the existence of a human in a vision.

## II. LITERATURE REVIEW

The method of identifying body gestures or motions and then predicting or defining stages of action or behavior is known as human action recognition (HAR) [3]. Its many applications, including Army care system, athletic rehabilitation from a disability or damage, and medical trials abnormality therapy, are stimulate industrial and academic interest for future research and development. With the pervasive use of handheld personal digital devices like smartphones, smartwatches, and multi-media stations that create a wide range of different forms of chronic data, such as recording devices, photograph streams, and geographic timbers There is going to be a considerable Personalization based on human action recognition, especially in national healthcare [2] which elucidate a deep learning-based HAR strategy Researchers use an accelerometer pulse to build a power spectrum visual, which they then inject into a convent. The feature extraction step is effectively replaced by the spectrogram synthesis phase, which adds some initial cost to the network training [3]. Taking pure accelerometer sensor as that of the feed to a convnet, perform a 1-D inversion on every enough. The spatial linkages between distinct sensor components may be lost as a result of this method. They concentrate on datasets that are freely available and originate mostly from embedded sensors (such as phones) or monitoring devices. [5] adopt a comparable tactics, Researchers have used the same existing data, and have used two-dimensional inversion to express dynamic impulses in a specific cable during their investigation [6] This uses an inter multilayer neural which incorporates linear accelerating and rotational mobility signals to categorize usual tasks using a labeled database comprising topmost motions for implementation of CNNs for the occupation classification task. Because the sorting assignment participants undertake is highly individual, data indicators collected out of each respondent are being used to train various learning models [8][24] so as to build a monitoring system out of six inertial measurement devices. Following network analysis, to create a functionality, the researchers listed a number of performance parameters that survived the statistical method and then they employed the random forest (RF) classifier to categorize the events. Finally, the correctness level was 84.6 %. The study described a haptic feedback inertial sensors action recognition system and its submission in healthcare diagnosis [6]. For feature selection, Reprieve and consecutive advance drifting scan (RCADS) were coupled. Finally, for action classification and comparison, the approaches of Nave Bayesian and k-nearest neighbor (KNN) have been used. ML algorithms could depend heavily on heuristics subjective extracting features mostly in people's daily movement detection tasks. The knowledge base possessed by humans is usually a major hurdle [9]. Deep learning techniques, that can automatically generate key features from telemetry data even during preprocessing step and substantially lower original descriptors with elevated conceptual patterns, have been used to resolve this difficulty. Adapting deep neural networks to the scope of human wearable sensors is indeed a recent research issue in analytical thinking [19], considering their effectiveness in classification, speech synthesis, computational linguistics, as well as other areas. The Zens et

al have suggested translating three-axis sensor readings into a "pictures" representation, then employing CNN with three convolutions and one fully-connected network to represent human movements. Consultation of additional literature available for researcher [18][23] has been done so as to improve perception of the broad field.

## III. RESEARCH METHODOLOGY

Deep Learning is used in this research to construct a system that recognizes human action. For the implantation of human action recognition on videos, the suggested system uses a convolutional neural network paired with a long short-term memory Network. In Tensor Flow, we'll need to use 2 distinct structures and ways. Finally, we'll adapt the successful system to YouTube video estimations [4].

### A. Image Classification

Image classification is where an image can be examined by a computer to see which 'class' it belongs to. (Or there's a chance it's a 'class.')

A class is simply a phrase, such as 'vehicle,' 'animal', and so on. Pass an image into a filter (whether it's a learned deep neural network (convolutional neural network or multi-layer perceptron) or a conventional classifier) to get class prediction. Let's say you want to insert a human photo. The image is analyzed by the computer, and it is determined to be a male or female. (Or the possibility that it's a female.) [16]

A video is simply a series of static images (known as frames) that are updated very quickly to provide the impression of motion. Take a look at the footage of a cat jumping on a bookcase below. (Converted to.gif format), which is simply a compilation of 15 separate still images that are updated one after the other [19].

### B. Video Classification

Video classification is the way of applying a tag to a central aspect on its own pixels. A competent YouTube clip predictor does far more than just provide precise panel tags; so, it reflects the real clip that uses the elements and tags of the consecutive frames [7]. For example, a clip might show a shrub inside one shot however the major championship (e.g., "hiking") could be something wholly distinct. The objective establishes the level of precision for the identifiers that must be used to characterize the pieces and clip. Common tasks include providing one or even more worldwide tags to a clip, and also one or maybe more descriptions for each panel from the inside of the clip [10].

### C. Convolutional Neural Network

A convolutional neural network (CNN or ConvNet) is a multi-layer perceptron that has been required in conducting to operate using visual data and specializes in vision estimation and forecasting. It performs with seeds (called filters) that are over the portrait and simply create a saliency map (which personifies how well a particular component is visible at a precise point of the image or not), and so it actually creates a comparatively tiny number of parameters initially.

However, as we proceed with greater depth into the intranet, this same number of nodes develops as well as the magnitude of graphs gets smaller without having to lose pertinent data utilizing accumulating processes [18].

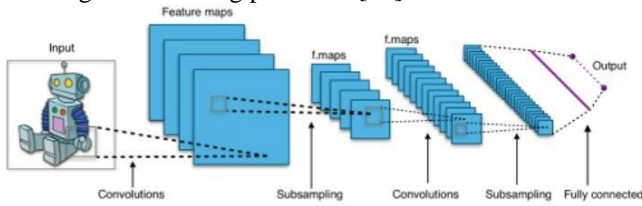


Figure 1: Convolutional Neural Network

The layers of a ConvNet learn features of increasing complexity, such as detecting shapes during the first layer and gaining a prominent place in diverse postures in the last layer.

**D. Long Short-Term Memory (LSTM)**

Because it considers all various sources in order when generating output, a long short - term memory system was developed specifically to accommodate a data series. Although LSTMs are a subtype of RNN (Recurrent neural network), RNN has been demonstrated to be inadequate in dealing with long-standing dependency in input sequences induced by the Disappearing slope problem. To overcome the disappearing slope and allow an LSTM unit to recall the context of long input data, LSTMs were developed [11].

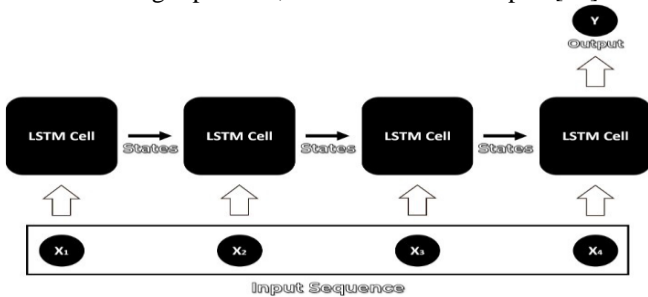


Figure 2: Long Short-Term Memory

This improves an LSTM's capacity to deal with sequential data problems including time series forecasting, speech synthesis, dictionary lookup, and orchestral composition. But for now, let's focus on how LSTMs can help us build more powerful action recognition models. We'll now look at how Action Recognizer will be put into action. we will utilize a deep convolution network (DCN) + long short-term memory (LSTM) Network to undertake Action Recognition when leveraging the geographic element of the clips [19].

**E. Recurrent Neural Network (RNN)**

Recurrent neural networks identify consecutive features of the information and anticipate the next likely situation using patterns. RNNs have been used in deep learning and the creation of algorithms that replicate neuronal behavior in the nervous system. They're extremely beneficial in cases where comprehension was essential to predicting a response. They differ from those other neural networks in that they use responses to examine a set of memories that influences the accurate results. Such natural cycles and good documentation to endure. The recall is a common term for this phenomenon [12][20][21][22].

Prominent Recurrent neural network usage instances comprise word embeddings wherein the subsequent word in

a term and the next expression in a statement are reliant on the facts that came before everything. Another innovative experiment was using a Recurrent neural network that received training on Literary plays and created Shakespeare-like prose. Computational inventiveness is a sort of RNN-based literature. This reproduction of artistic expression is facilitated by the AI's linguistic knowledge and interpretation obtained from its training phase.

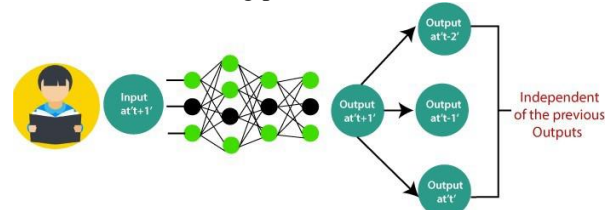


Figure 3: Recurrent Neural Networks [18]

**F. LSTM vs RNN**

The long short-term memory network is a type of recurrent neural network including both special and conventional subunits. There really is a 'memory block' in LSTM units that could store information for long durations. This memory cell allows them to learn longer- term dependence. The main distinction between RNN and LSTM is it retains information in memory for a lengthy period of time. In this scenario, LSTM outperforms RNN because it can keep information in memory for a longer period of time [13].

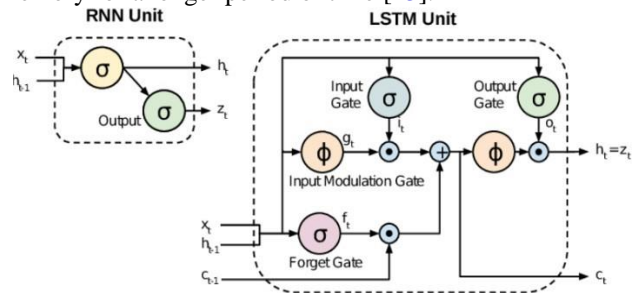


Figure 4: LSTM vs RNN [16]

**IV. PROPOSED MODEL**

**A. CNN + LSTM**

A CNN would be used to retrieve distance measures at a certain sampling interval inside the original signal (clip), and an LSTM would be used to discover and predict future performance among panels.

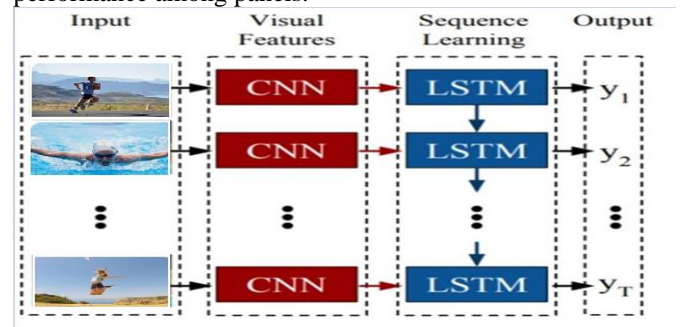


Figure 4: CNN and LSTM [14]



These two structures we are using to combine CNN and LSTM are as follows:

- 1 ConvLSTM
- 2 LRCN

TensorFlow can be used for each of these approaches.

## B. Conv LSTM

ConvLSTM is a recurrent neural network featuring multilayer properties within both insight and state-to-state stages which anticipates spatial patterns. The ConvLSTM forecasting condition of the unit inside the grid uses supplies and previous conditions of its local neighbors. A simple way to achieve this is by using a transformation function in the state-to-state and insight conversions. The following are indeed the ConvLSTM fundamental equations, where '\*' represents the convolution operator and 'o' indicates the Hadamard product:

$$it = \sigma(Wxi * Xt + Whi * Ht-1 + Wci \circ Ct-1 + bi)$$

$$ft = \sigma(Wxf * Xt + Whf * Ht-1 + Wcf \circ Ct-1 + bf)$$

$$Ct = ft \circ Ct-1 + it \circ \tanh(Wxc * Xt + Whc * Ht-1 + bc)$$

$$ot = \sigma(Wxo * Xt + Who * Ht-1 + Wco \circ Ct + bo)$$

$$Ht = ot \circ \tanh(Ct)$$

We should anticipate a ConvLSTM with a larger transitioning unit to detect greater oscillations and a lesser unit to identify delayed actions if we interpret its transitions to be disguised descriptions of moving objects.

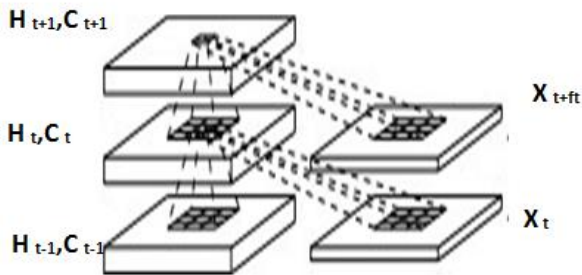


Figure 6: Inner Structure of Conv LSTM

Before executing the convolutional operation, padding is essential to ensure that results have the same number of data points as the supplies. Determining the aspect of the outside globe could be equivalent to padding the disguised conditions on the border crossing points. Before the first entry, we generally set all the variables of the LSTM to zeroes, which essentially translates to "profound ignorance" of the long term.

## C. LRCN

A model for activities that involve sequence information (inputs or outputs), perhaps optical, verbal, or other, that blends a fully convolutional graphical extracted features (such as a CNN) with such a system that really can be taught to identify and manufacture time evolution. The core of our approach is depicted in Figure 7.

The addition of new features  $\phi V(\cdot)$  with variables  $V$ , almost always a CNN, is applied to every realizing  $x_t$  (a single picture or video frame) to produce a repaired generative model  $\phi V(x_t)$ . The results of  $V$  are therefore already in modules that learn recurring sequences.

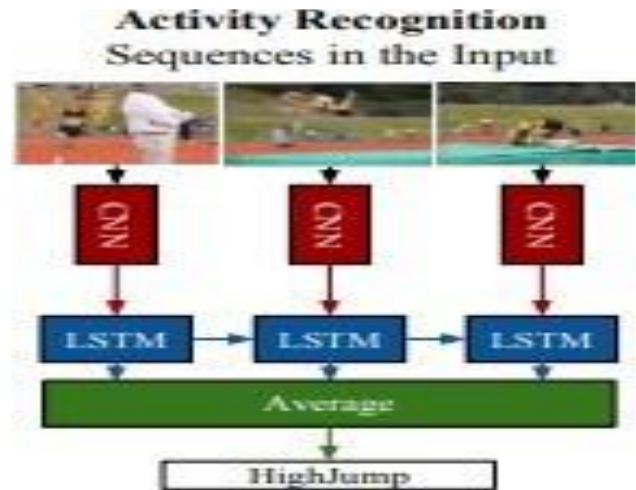


Figure 7: LRCN Model for Action Recognition

## D. Data Collection

UCF50 is a data set of 50 action categories for action recognition culled from real-life action videos on YouTube. Such set of data is a supplement to the Vimeo Action set of data (UCF11), which contains 11 different action categories. The common of the action recognition numbers sets available are unrealistic and staged by actors. Our primary objective in gathering data is to provide the computer image community with a source of action recognition data drawn from realistic YouTube videos. Our data set is highly demanding due to large changes in camera activities, object look and posture, object ruler, perspective, disorderly backdrop, lighting conditions, and other parameters.

## V. MODEL CREATION

### A. Conv LSTM Approach

We combine ConvLSTM cells to execute the first approach. Inside the system, the ConvLSTM unit is a variant of the LSTM unit which includes convolution procedures. It was an LSTM having convolution built-in, enabling it all to distinguish spatial properties as information even while considering predicting future performance into account. Regarding video identification, this method numerically captures the relationship between individual frames in terms of both space and time. While a regular LSTM can only handle one dimension of input, due to its convolution structure, the ConvLSTM can handle three dimensions (width, height, and number of channels). As such, it is not appropriate for modeling spatiotemporal data on its own [17].

### B. Model of ConvLSTM

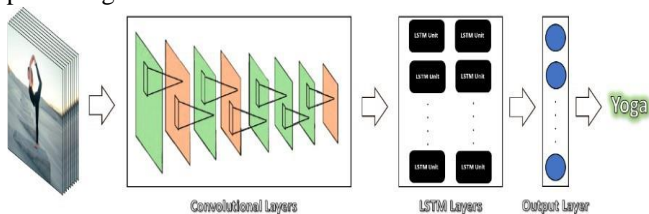
Recurrent layers from Keras ConvLSTM2D will be used to construct the model. Additionally, the filter and kernel sizes required for multilayer processes are considered by the ConvLSTM2D layer. Finally, a deep network with a linear kernel is fed the flattened layer output, which provides a frequency for every action category. We'll also utilize Dropout layers to avoid overfitting the model on the data and MaxPooling3D layer to reduce the size of the frames and avoid wasteful calculations.

The design is straightforward, with only a few trainable parameters. This is due to the fact that we are only dealing with a small portion of the information, which does not need the adoption of a large-scale model [18].

**C. LRCN Approach**

To apply the LRCN Approach, we'll merge the Pooling layer and LSTM layers into a single component in this stage. In a similar way, a Convolution layer and a Hidden layer prototype instructed individually could be used. A pre-trained system that could be perfectly all right for the purpose could be used to retrieve spatial features from video sequences to use the Network model. The LSTM model could then use the Feature representation to forecast the action taken inside the clip.

However, the Long-term Recurrent Convolutional Network (LRCN), which integrates CNN and LSTM layers into a single model, will be used in this case. The LSTM layer(s) is(are) tasked with temporal sequence modeling after the Convolutional layers have extracted spatial properties from the frames at each time step. The network can acquire spatiotemporal features through end-to-end training, producing a reliable model.



**Figure 9: LRCN Approach**

A Time Distributed container component would also be used, enabling us to deploy the same surface toward each pixel of the clip independently. As just a corollary, if indeed the gradient intake structure were (breadth, elevation, number of streams), it can additionally receive information of pattern (number of images, spacing, size, number of channels), which is tremendously useful as it enables users to enter the entire movie into the method in one go [19].

**D. LRCN Model**

We will use time-circulated Conv2D layers, MaxPooling2D, and Dropout strands to develop the LRCN layout. The feature gathered from the Conv2D layers will be flattened and sent to an LSTM layer using a flattening surface. The output from the LSTM layer would then be used by the Thick layer, which uses hidden neurons, to forecast the operation that would be carried out.

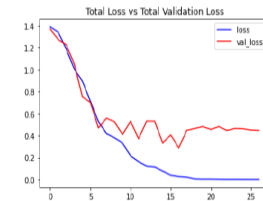
**E. LRCN Model Evaluation**

We'll analyze the model on testing data after it has been trained, and the accuracy rate is 92 %, with a loss rate of 22 %.

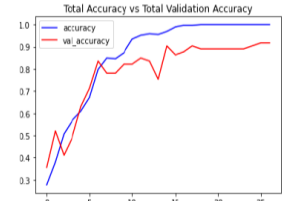
**VI. RESULTS**

In this paper we have used two different approaches i.e., is a grouping of convolution neural networks and long-short term memory. The First approach we have used is ConvLSTM and the second approach we have used is LRCN. Here is the visualization of the ConvLSTM

**Loss of ConvLSTM model**

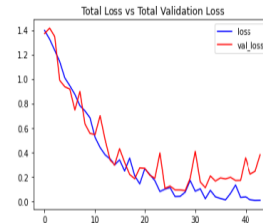


**Accuracy of ConvLSTM model**

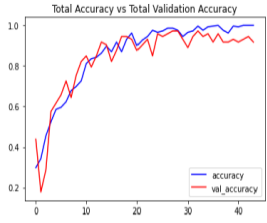


Here is the visualization of the LRCN Approach

**Loss of LRCN model**



**Accuracy of LRCN model**



Approach	Loss	Accuracy
ConvLSTM	89%	80%
LRCN	22%	92%

Conv LSTM model accuracy rate is 80 % and LRCN model accuracy rate is 92%. As we can see that LRCN give a better accuracy as compared to Conv LSTM model. So, we will use LRCN Model to perform prediction with it on YouTube videos.

**VII. CONCLUSION**

Human action recognition does have a variety of applications because of its effect on people's pleasure. It has become a major tool in personalized medicine including preventing weight and elderly care. In this paper, we have used a convolutional neural network collective with long short-term memory. we tried two different approaches. First, we create a model using ConvLSTM, and the second model using LRCN and TensorFlow was used for each of these approaches. We have used the UCF50 dataset for model creation. UCF50 is a data set of 50 action categories for action recognition culled from real-life action videos on YouTube. ConvLSTM approach produces 80% accuracy with 89 % of loss, and the LRCN approach produces 92% accuracy with 22% loss. After seeing the result of both the approaches, we found that the LRCN approach produced a better result as compared to ConvLSTM. So, LRCN Model was used to perform prediction with it on YouTube videos.

**DECLARATION STATEMENT**

Funding	No, I did not receive it.
Conflicts of Interest	No conflicts of interest to the best of our knowledge.
Ethical Approval and Consent to Participate	No, the article does not require ethical approval and consent to participate with evidence.
Availability of Data and Material	Not relevant.
Authors Contributions	I am only the sole author of the article.

## REFERENCES

1. M. G. Morshed, T. Sultana, A. Alam, and Y. K. Lee, "Human Action Recognition: A Taxonomy-Based Survey, Updates, and Opportunities," *Sensors*, vol. 23, no. 4. MDPI, Feb. 01, 2023. doi: 10.3390/s23042182. <https://doi.org/10.3390/s23042182>
2. T. O. Araoye, E. C. Ashigwuike, A. C. Adeyemi, S. V. Egoigwe, N. G. Ajah, and E. Eronu, "Reduction and control of harmonic on three-phase squirrel cage induction motors with voltage source inverter (VSI) using ANN-grasshopper optimization shunt active filters (ANN-GOSAF)," *Sci Afr*, vol. 21, Sep. 2023, doi: 10.1016/j.sciaf.2023.e 01785. <https://doi.org/10.1016/j.sciaf.2023.e01785>
3. K. D. Addo, F. Davis, Y. A. K. Fiagbe, and A. Andrews, "Machine learning for predictive modelling of the performance of automobile engine operating without coolant thermostat," *Sci Afr*, vol. 21, Sep. 2023, doi: 10.1016/j.sciaf.2023.e 01802. <https://doi.org/10.1016/j.sciaf.2023.e01802>
4. P. Lara-Benítez, M. Carranza-García, and J. C. Riquelme, "An Experimental Review on Deep Learning Architectures for Time Series Forecasting," *Int J Neural Syst*, vol. 31, no. 3, Mar. 2021, doi: 10.1142/S0129065721300011. <https://doi.org/10.1142/S0129065721300011>
5. E. C. Nwosu, G. N. Nwaji, C. Ononogbo, I. Ofong, N. V. Ogueke, and E. E. Anyanwu, "Effects of water thickness and glazing slope on the performance of a double-effect solar still," *Sci Afr*, vol. 21, Sep. 2023, doi: 10.1016/j.sciaf.2023.e 01777. <https://doi.org/10.1016/j.sciaf.2023.e01777>
6. S. Chae, J. Shin, S. Kwon, S. Lee, S. Kang, and D. Lee, "PM10 and PM2.5 real-time prediction models using an interpolated convolutional neural network," *Sci Rep*, vol. 11, no. 1, Dec. 2021, doi: 10.1038/s41598-021-91253-9. <https://doi.org/10.1038/s41598-021-91253-9>
7. H. W. Hounkpatin, H. E. V. Donnou, K. V. Chegnimonhan, G. H. Hounguè, and B. B. Kounouhewa, "Thermal characterisation of insulation panels based on vegetable typha domengensis and starch," *Sci Afr*, vol. 21, Sep. 2023, doi: 10.1016/j.sciaf.2023.e 01786. <https://doi.org/10.1016/j.sciaf.2023.e01786>
8. S. Ghimire, Z. M. Yaseen, A. A. Farooque, R. C. Deo, J. Zhang, and X. Tao, "Streamflow prediction using an integrated methodology based on convolutional neural network and long short-term memory networks," *Sci Rep*, vol. 11, no. 1, Dec. 2021, doi: 10.1038/s41598-021-96751-4. <https://doi.org/10.1038/s41598-021-96751-4>
9. H. A. Mupambwa et al., "A 2-year study on the spatio-temporal changes in trace metal concentrations in sediment, water and plants within the Walvis Bay Lagoon, Namibia," *Sci Afr*, vol. 21, Sep. 2023, doi: 10.1016/j.sciaf.2023.e01787. <https://doi.org/10.1016/j.sciaf.2023.e01787>
10. J. Donahue et al., "Long-term Recurrent Convolutional Networks for Visual Recognition and Description," Nov. 2014, [Online]. Available: <http://arxiv.org/abs/1411.4389> <https://doi.org/10.21236/ADA623249>
11. O. G. Famutimi, I. O. Adewale, and K. R. Adegoke, "Inhibition characteristics of peptide extracts of four medicinal plants on activities of bovine trypsin," *Sci Afr*, vol. 21, Sep. 2023, doi: 10.1016/j.sciaf.2023.e01795. <https://doi.org/10.1016/j.sciaf.2023.e01795>
12. M. Shreyas Pagare and D. Rakesh Kumar, "Survey on Human Action Recognition System, Challenges and Applications," 2023.
13. M. O. Mario, "Human activity recognition based on single sensor square HV acceleration images and convolutional neural networks," *IEEE Sens J*, vol. 19, no. 4, pp. 1487–1498, Feb. 2019, doi: 10.1109/JSEN.2018.2882943. <https://doi.org/10.1109/JSEN.2018.2882943>
14. G. Theses and J. Pang, "Scholar Commons Human Activity Recognition Based on Transfer Learning," 2018. [Online]. Available: <https://scholarcommons.usf.edu/etd>
15. T. N. Sainath, O. Vinyals, A. Senior, and H. H. Sak, "Convolutional, Long Short-Term Memory, Fully Connected Deep Neural Networks."
16. S. Gopali, F. Abri, S. Stami-Namini, and A. S. Namin, "A Comparative Study of Detecting Anomalies in Time Series Data Using LSTM and TCN Models," Dec. 2021, [Online]. Available: <http://arxiv.org/abs/2112.09293>
17. Y. Kim, Y. Jernite, D. Sontag, and A. M. Rush, "Character-Aware Neural Language Models," Aug. 2015, [Online]. Available: <http://arxiv.org/abs/1508.06615>
18. O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, "Show and Tell: A Neural Image Caption Generator," Nov. 2014, [Online]. Available: <http://arxiv.org/abs/1411.4555>
19. X. Shi, Z. Chen, H. Wang, D.-Y. Yeung, W. Wong, and W. Woo, "Convolutional LSTM Network: A Machine Learning Approach for Precipitation Nowcasting," Jun. 2015, [Online]. Available: <http://arxiv.org/abs/1506.04214>
20. Aralimarad, M., S. M., Dr. M., & Mallapur, Dr. J. D. (2020). A Comprehensive Survey on Human Action Recognition. In *International Journal of Recent Technology and Engineering (IJRTE)* (Vol. 9, Issue 2, pp. 902–908). <https://doi.org/10.35940/ijrte.b3933.079220>
21. Jahagirdar, A., & Nagmode, M. (2019). Human Action Recognition using Scaled Convolutional Neural Network. In *International Journal of Engineering and Advanced Technology* (Vol. 9, Issue 1, pp. 1820–1826). <https://doi.org/10.35940/ijeat.a1440.109119>
22. Magapu, H., Krishna Sai, M. R., & Goteti, B. (2024). Human Deep Neural Networks with Artificial Intelligence and Mathematical Formulas. In *International Journal of Emerging Science and Engineering* (Vol. 12, Issue 4, pp. 1–2). <https://doi.org/10.35940/ijese.c9803.12040324>
23. Panicker, M. J., Upadhyay, V., Sethi, G., & Mathur, V. (2021). Image Caption Generator. In *International Journal of Innovative Technology and Exploring Engineering* (Vol. 10, Issue 3, pp. 87–92). <https://doi.org/10.35940/ijitee.c8383.0110321>
24. Das, S., S. S., M. A., & Jayaram, S. (2021). Deep Learning Convolutional Neural Network for Defect Identification and Classification in Woven Fabric. In *Indian Journal of Artificial Intelligence and Neural Networking* (Vol. 1, Issue 2, pp. 9–13). <https://doi.org/10.54105/ijainn.b1011.041221>

## AUTHORS PROFILE



**Shreyas Pagare** is a Ph.D. scholar in the Department of Computer Science Engineering from Rabindranath Tagore University Bhopal (M.P.). He completed his Master of Engineering (Computer Science & Engineering) from IET DAVV, Indore in 2012. His interests include Database Technologies, Data Structures, Data Mining, Data Warehousing, Artificial Intelligence, and Machine Learning. He has experience is about 15 years of experience in teaching and Research at the college and University level. He has organized many Trainings and Workshops. He has delivered many Expert Lectures in various organizations. He is working on various funded projects. He has published many Research papers in Reputed Journals (i.e., SCOPUS and UGC Approved) and presented many papers in various National and International Conferences. He is the author of Book Chapters and Patents.



**Dr. Rakesh Kumar** is an Associate Professor, Department of Computer Science & Engineering and Head, Centre for IoT & Advance Computing. He is member of Institute Innovation Council (IIC) MHRD and the Innovation Ambassador at University. He completed his Ph.D. in Computer Science & Engineering. His area interests includes IoT, IoE, Software Engineering, Operating System, Database Technologies, Data Structures, Data Mining, Data Warehousing, Robotics, Drone Technology, Artificial Intelligence and Computer Architecture. He has experience is about 17 years of experience of teaching and Research at college and University level. He has organized many Trainings and Workshops. He has delivered many Expert Lectures in various organizations. He has organized AICTE funded FDP. He has guided various Projects at State and National Level competitions. He is working on various funded projects. He has published many Research papers in Reputed Journals (i.e., SCI, SCOPUS and UGC Approved) and presented many papers in various National and International Conferences. He is the author of Book Chapters and Patents. He has appointed as an Editorial board member and reviewer in various Journals, International conferences. He was recognize by "Quarterly Franklin Membership" (Membership ID#VY34049), London Journal of Engineering Research (LJER), London Journals Press (UK) for research paper "Inside Agile Family Software Development Methodologies. He is the founder member of Centre for IoT & Advance Computing.

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of the Blue Eyes Intelligence Engineering and Sciences Publication (BEIESP)/ journal and/or the editor(s). The Blue Eyes Intelligence Engineering and Sciences Publication (BEIESP) and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.