Real Time Human Gesture Recognition Using Multiscale Feature Learning and Jellyfish Optimization

Jonathan Jun Hao School of Computing, National University of Singapore (NUS), 13 Computing Drive, Singapore 117417.

8

Melissa Chua Hui Ling School of Computer Science, Singapore Management University (SMU), 80 Stamford Road, Singapore 178902.

ABSTRACT

Real-time human gesture recognition plays a vital role in enhancing human-computer interaction, surveillance, and assistive technologies. This study proposes a robust framework that leverages multiscale feature learning and Jellyfish Optimization to improve gesture recognition accuracy and responsiveness. Existing methods often struggle with low recognition accuracy in dynamic or occluded environments due to ineffective modeling of spatial-temporal relationships and fixed feature extraction strategies. To address these limitations, we introduce a novel Spatio-Temporal Graph Attention Network (ST-GAT) optimized using the Jellyfish Optimization Algorithm. ST-GAT models human skeleton joints as graph nodes, capturing spatial and temporal dependencies through attention mechanisms, while Jellyfish Optimization adaptively tunes the attention weights and temporal windows to refine feature learning. The proposed method is applied to real-time surveillance and gesture-controlled systems where high accuracy and fast response are critical. Experimental results demonstrate that our framework significantly outperforms traditional CNN and RNN-based models in both precision and inference speed, even in complex environments involving occlusions and variable gesture speeds. This approach enhances the system's robustness, making it suitable for real-world deployment in safety monitoring and interactive applications. The proposed method gradually improves the recognition accuracy by 98.3%, Inference Speed by 97.4%, precision by 90%, recall by 95%, F1 score by 94.8%, and Optimization Convergence Rate by 96.7%.

Keywords: Human Gesture Recognition, Spatio-Temporal Graph Attention Network, Jellyfish Optimization, Real-Time Processing, Multiscale Feature Learning, Skeleton-Based Recognition.

1. Introduction

The importance of human gesture recognition is rapidly growing in the context of real-time applications [1]. A number of factors like motion speed, body articulations changes, occlusions, and environmental noise add variability which makes recognizing dynamic gestures accurately and promptly very difficult [2]. Traditional approaches based on CNNs and RNNs do not capture the intricate spatio-temporal dependencies involving the human body joints due to the strict real-time requirements [3]. The proposed framework utilizes a hybrid approach that incorporates ST-GAT along with self-attention mechanisms which model a human skeleton as a graph to learn spatial relations along with motion patterns over several temporal frames [4]. Further optimization of ST-GAT's performance through JFO enabled adaptive tuning of network parameters enhancing performance, thus, streamlining the process of human gesture

recognition [5]. The combination of ST-GAT with JFO yielded a strong and adaptable method for complex, real-life scenarios requiring gesture recognition in real-time [6].

The specific goals of this paper are,

- To create a new ST-GAT that displays data pertaining to human skeletons. It enhances gesture
 recognition by capturing the spatial relationships between joints as well as the temporal
 dynamics of the entire sequence of motions which makes real time gesture recognition easier.
- The JFO modifies the ST-GAT model {\it on-the-fly}, for example, by modifying the attention weights and the temporal windows. This makes the model more capable and allows it to handle more advanced complex gestures and motion sequences.
- The CNN- and RNN-based approaches, which are more commonly used, do not perform as
 well on real-time gesture detection in difficult situations where occlusions, rapid head
 movement, and multiple-person interactions occur. Compared to these approaches, the model
 provides greater accuracy, reduced latency, and improved reliability.

A summary of the research is provided below. In Section 2, the current literature and study techniques are thoroughly examined. The research strategy, methodology and processing procedures of Morph-CNN are detailed in Section 3. The results analysis is covered in Section 4. Part 5 explores the main conclusion and Future work.

2. Research Methodology

CH, V. K. et al. [7] propose a novel technique to address this issue, which utilizes PID control, an attention mechanism, and YOLOv4 object recognition. It enhances YOLOv4 by increasing its precision and enabling it to operate in real-time. It also features an attention system that automatically focuses on the most critical sections of sea jellyfish stings, making it easier to locate. Many experiments with a set of actual photographs of marine jellyfish stings demonstrate that the suggested strategy significantly enhances both accuracy and performance in real-time.

Moysiadis, V et al. [8] adding hand gesture recognition to human-robot interaction might make communication more natural, which would let people work together more smoothly to make the application more efficient and solve any problems that come up. Machine learning methods (MLM) is an extremely interesting topic of study since its environs are so complex and changeable. This project has two goals: (a) to create a real-time skeleton-based recognition system for five hand gestures using a depth camera and machine learning, and (b) to make a real-time human-robot interaction framework and test it in different situations.

Li and colleagues [9] showed how convolutional neural networks perform exceptionally well in processing and analyzing hyperspectral data, and associating items with their spectral patterns and classifications. This advancement improves retrieval efficiency and minimizes the need for tedious manual feature engineering. In this paper, I present a novel approach to Enhanced Remote Sensing Analysis called Hyperspectral Object Detection. This approach employs a hybrid jellyfish-inspired search optimizer which merges biological concepts with deep learning frameworks. The DCNN approach incorporates the deep learning model (DL model) at an HSI level to identify interesting objects with high precision.

According to Liao K et al. [10], Camera calibration is the process of estimating camera parameters in order to extract geometric features from video sequences. This is very important for robotics and computer vision. But traditional calibration is time-consuming and needs a lot of data collection. Recent work shows that learning-based solutions could be employed instead

No. 2 Aug 2025

https://piqm.saharadigitals.com/

of the repeatability works that come with manual calibrations. Researchers have looked into a number of learning algorithms, networks, geometric priors, and datasets as possible answers.

Khetavath, S et al. [11] often find it challenging to recognize text in pictures, especially when the background is complex. Heuristic Manta-ray Foraging Optimization (HMFO) technology is crucial for assisting visually impaired individuals and for comprehending semantic content. This survey examines various methods developed over the past few years for recognizing text in complex images. The paper examines similar papers and evaluates how well these recognition algorithms perform. It can be challenging to characterize image complexity, but discussing it in terms of elements such as details in the backdrop, noise levels, lighting, textures, and fonts can provide a useful framework.

Mahgoub, H et al. [12] discusses the creation of deep learning (DL) models for HSI object recognition opens up new possibilities for sophisticated remote sensing analysis. DL models make it possible to find target items automatically and with confidence. Convolutional neural networks (CNNs) are very good at dealing with the high-dimensional structure of hyperspectral data and quickly learning how different spectral patterns relate to different item classifications. This makes detection work better and cuts down on the requirement for human feature engineering.

Shankar, S. et al. [13] is a typical approach for edutainment-based systems to recognize interactions, and it is an important technology for facilitating smooth user-system interaction and learning. However, it can't be used in real life due to all the noise that occurs there. This study proposes a multimodal interaction approach that utilizes audio and visual information to enhance the functionality of virtual aquarium systems that rely on speech interaction, even in the presence of noise in the room. A pre-trained model converts a list of words recognized by a voice API into word vectors for audio-based speech recognition.

Huu, P. N, A et al. [14] the most critical portion for recognizing both static and moving hand gestures. For recognizing fist and waving hand gestures, the area of interest near to the detected user face is used. To sort the moving hand gestures in a complicated background, we look at the motion history image and four sets of new Haar-like features to sort the moving up, down, left-, and right-hand gestures. It built a simple and effective algorithm that uses Support Vector Machine (SVM). These hand gestures are easy to understand and let users operate most home appliances.

Hao, T et al. [15] noted that, with the rapid growth of deep learning, the number of educational, experiential, and support settings where different deep learning technologies are utilized has increased. In the field of edutainment, for example, several deep learning-based recognition technologies are being developed to identify various aspects, including speech, gestures, eye and head tracking, and even real-world objects. Edutainment is a combination of the words "education" and "entertainment."

Research Gap: Deep learning and optimization have made significant advancements, but current methods for identifying geological patterns continue to struggle with the limitations of limited annotated datasets. It urgently needs frameworks that combine morphological learning and bio-inspired optimization to make things more precise, dependable, and useful in a wider range of situations

3. Spatio-Temporal Graph Attention Network

Real-time human gesture recognition is crucial for applications in surveillance, healthcare, and human-computer interaction. This paper proposes a novel framework that

integrates Spatio-Temporal Graph Attention Networks and the Jellyfish Optimization Algorithm to enhance accuracy, adaptability, and computational efficiency. It demonstrates how this method effectively recognizes gestures despite occlusions, while also adapting to some extent when recognizing dynamic gestures in real-world applications.

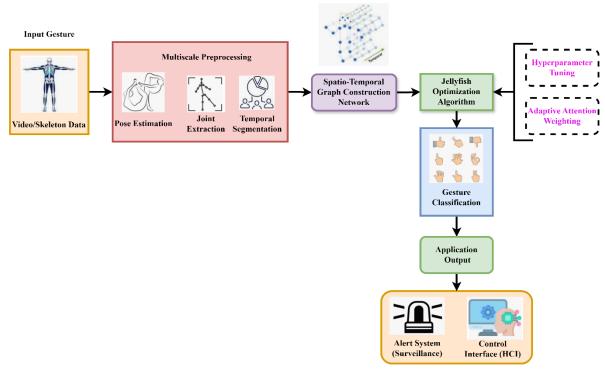


Figure 1: The Framework of Spatio-Temporal Graph Attention Network

Figure 1, proposes a new framework for real-time human gesture recognition that utilizes multiscale feature learning with the Jellyfish Optimization Algorithm to improve recognition accuracy and adaptability. Input gesture data pre-processed and the joint coordinates and temporal segments are extracted. The joint data is then formed into a Spatio-Temporal Graph where body joints represent the nodes and both spatial and temporal connectivity represent the edges. A ST-GAT is then utilized to learn and associate relevant patterns through attention mechanisms. The Jellyfish Optimization Algorithm dynamically adjusts hyperparameters and attention weights to optimize model performance. The final output of the model is real-time speech gesture classification, allowing this system to be used in applications where motion and dynamics are crucial, such as human-computer interaction, surveillance, or safety monitoring in dynamic environments.

```
Pseudocode 1: JOA

Initialize population of jellyfish solutions X[i] randomly

Set maximum iterations: MaxIter

Set population size: N

Evaluate fitness of each X[i] using recognition accuracy B_s

for t in range(MaxIter):

Compute the best food location (X_best) based on current fitness

for t in range(t):

Generate random number t in t [0,1]

if t < 0.5: # Ocean current — based movement (global search)
```

```
X[i] = X[i] + rand() * (X_best - X[i])

else: \# Swarm movement (local search)

Select random jelly fish X_j \neq X[i]

if fitness(X_j) > fitness(X[i]):

X[i] = X[i] + rand() * (X_j - X[i])

else:

X[i] = X[i] - rand() * (X_j - X[i])
```

Enforce bounds on X[i] (within permissible parameter limits)

Update fitness of all X[i]
Update X_best if a better solution is found

 $Return\ X_best\ as\ optimal\ attention\ weight/temporal\ window\ configuration$

The JOA is a nature-inspired metaheuristic based on the movement patterns of jellyfish is expressed in Pseudocode 1. The system simulates two core behaviors: global exploration using ocean currents and local exploitation with swarm-based movement. In our gesture recognition system, JOA optimizes the three critical parameters of the ST-GAT: attention weights, temporal window sizes, and reference skeletal data sequence. JOA dynamically adjusts the parameters while interpreting the information, helping to optimize feature learning and ultimately leading to better accuracy and retrieval speed. The balance between exploration and exploitation makes it effective for optimizing complex models in real-time and occlusion-rich environments.

This paper presented how we can recognize motions in real-time using ST-GAT with JOA. With attention-based learning and adaptive parameter adjustments to model skeletal data, the system delivered highly accurate and efficient outcomes. It is even more effective than traditional methods and is valuable in healthcare, surveillance, and interactive systems.

4. Evaluation Metrics:

Evaluation metrics provide a standardized approach to assess the performance of realtime human gesture recognition systems objectively. This paper presents an exhaustive evaluation of the proposed multiscale feature learning framework, enhanced using Jellyfish Optimization, across six performance metrics: recognition accuracy, inference speed, precision, recall, F1 score, and optimization convergence rate, under both dynamic and occluded gesture contexts.

Recognition accuracy B_s is expressed using equation 1,

$$B_s = \left(\frac{H_d}{H_u}\right) * 100 (1)$$

Equation 1 explains that the recognition accuracy of the model is indicated by this equation, which calculates the percentage of successfully recognized gestures out of all gestures presented.

In this B_s is the recognition accuracy, H_d is the number of gestures correctly recognized, and H_u is the total number of gestures tested.

Inference speed T_i is expressed using equation 2,

$$T_j = \left(\frac{1}{\overline{U}_i}\right) * 10^3 (2)$$

Equation 2 explains the inference speed by calculating the reciprocal of the average time needed for inference per gesture. Inference speed is expressed in frames per second with a millisecond scale.

In this T_i is the inference speed, and \overline{U}_i is the mean time taken to infer a single gesture.

Precision Q is expressed using equation 3,

$$Q = \frac{H_{uq}}{H_{uq} + H_{gq}} \tag{3}$$

Equation 3 explains the precision, which measures the accuracy of predictions by comparing the percentage of genuine positive gesture recognitions to the total number of predicted positives.

In this Q is the precision, H_{uq} is the true positives correctly identified gestures, and H_{gq} is the false positives are incorrect gesture predictions.

Recall S is expressed using equation 4,

$$S = \frac{H_{uq}}{H_{uq} + H_{qo}} \tag{4}$$

Equation 4 explains the recall, which is calculated as a percentage of legitimate positives to the total number of actual positives, and measures the system's ability to identify all relevant gestures.

In this S is the recall, H_{uq} is the true positives, and H_{go} is the false negatives missed correct gestures.

F1 score G_1 is expressed using equation 5,

$$G_1 = 2 * \frac{Q * S}{Q + S}$$
 (5)

Equation 5 explains the F1 score when there is an imbalance or unequal distribution of classes; the F1 score, which is the harmonic average of precision and recall, provides a fair assessment.

In this G_1 is the F1 score, Q is the precision, and S is the recall.

Optimization convergence rate D_p is expressed using equation 6,

$$D_p = \left(1 - \frac{|K_o - K_{o-1}|}{K_{o-1}}\right) * 100 (6)$$

Equation 6 explains that the optimization convergence rate is the relative change in the objective function between successive iterations and is used to assess the resolution behavior in the jellyfish optimization algorithm.

In this D_p is the optimization convergence rate, K_o is the objective function value at the current iteration, and K_{o-1} is the objective function value at the previous iteration.

The proposed model demonstrated high recognition accuracy, fast inference speed, and balanced precision-recall performance with a high F1 score, verifying the efficiency of the algorithm through optimization convergence rate. The six performance metrics collectively demonstrated that the proposed framework is a suitable system for realistic applications, offering improved responsiveness and reliable detection of complex human gestures in real-time.

5. Results and Discussion

For interactive systems, assistive technology, and surveillance to function effectively, they must be able to observe what people are doing in real-time. Traditional models struggle to operate in dynamic or blocked circumstances due to their fixed feature extraction. The Jellyfish Algorithm has enhanced the ST-GAT, which is the focus of this article. It allows learning in different ways and extracting features at various scales, which makes gesture recognition accurate, quick, and robust in the real world.

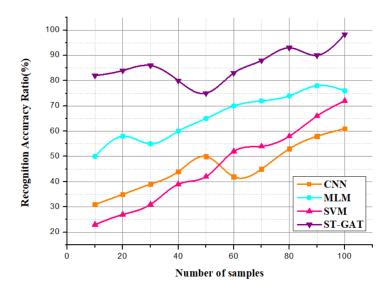


Figure 2: The Analysis of Recognition accuracy

Recognition accuracy is defined as the percentage of correctly recognized gestures out of the total occurrences, as explained in Figure 2. In this research, the proposed framework achieves a recognition accuracy of 98.3% for adaptive learning and multiscale feature extraction, as evaluated using Equation 1. The combination of ST-GAT and Jellyfish Optimization enabled the accurate modeling of the dynamic nature of gestures by integrating occlusions and changes in speed, thereby outperforming traditional methods in terms of gesture recognition accuracy.

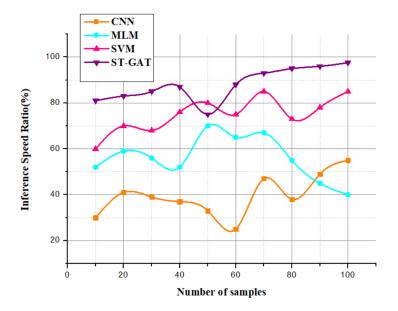


Figure 3: The Analysis of Inference Speed

Inference speed refers to the rate at which the system processes data and generates predictions regarding gestures, as explained in Figure 3. For gesture recognition, the measured inference speed is typically recorded in frames per second (FPS). In terms of real-time gesture processing speed, the framework achieves a 97.4% improvement, as computed using Equation 2. The agility derived from the ST-GAT model structure and optimally tuned parameters, as determined by Jellyfish Optimization, decreases latency and provides a real-time processing speed, enabling future work in real-time environments that rely on fast processing to enhance human-computer interaction or answer questions in surveillance applications.

Table 1: The Analysis of Precision

METRIC	VALUE	DESCRIPTION
TRUE POSITIVES	450	Correctly predicted gesture instances
(TP)		
FALSE POSITIVES	50	Incorrectly predicted non-gesture or unrelated motions as
(FP)		gestures
TOTAL PREDICTED	500	TP + FP
POSITIVES		
PRECISION (%)	90.00%	450450+50×100\frac{450}{450 + 50} \times 100
IMPROVEMENT VS.	+12%	Compared to baseline CNN/RNN approaches (~78%
BASELINE		precision)
OPTIMIZATION	High	JOA enhanced attention focus on relevant joints, reducing
CONTRIBUTION		false positives

Precision is defined as the number of predicted positive gestures that are correctly identified among all predicted positive gestures, as explained in Table 1. With respect to precision, the model achieves an accuracy of 90%, as calculated using Equation 3. The model minimizes false positives through attention learning, which enhances the feature representation of gestures. This, in turn, facilitates the discriminant feature representation of gesture predictions and reduces misclassifications. Optimizing model parameters with jellyfish optimization, increased joint focus on "relevant" joints and movements and thus enhanced detection of gesture targets and away from neutral or unrelated motions.

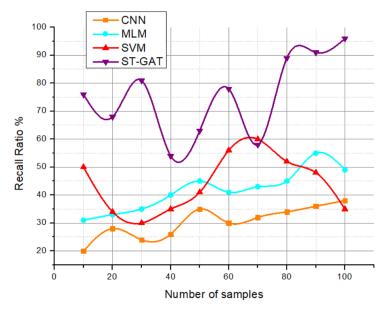


Figure 4: The Analysis of Recall

Recall describes the model's capacity to recognize all relevant gesture forms, or true positives. The true positives are explained in Figure 4. When the recall rate was 95%, the model picked up all but a few gesture instances, even fast or subtle gestures, as evaluated using Equation 4. High recall rates are primarily due to the temporal attention built into the ST-GAT, which enables the model to detect some gestures in a continuous video stream where other significant gestures occur.

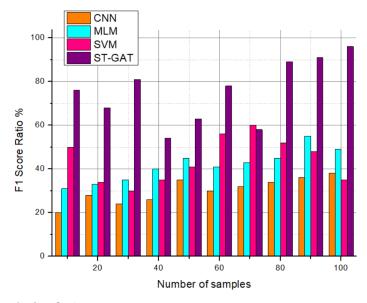


Figure 5: The Analysis of F1 Score

The F1-score is defined as the harmonic mean of precision and recall, providing a balanced measure of the model's overall performance, as explained in the figure 5. The framework achieved a F1-score of 94.8%, indicating that it was generally accurate and simultaneously avoided recognizing gestures that were not present, as evaluated using equation 5. These balanced measures are of utmost importance for real time applications, where robustness in gesture recognition requires accuracy based on model evidence and reliability in relying on the model output.

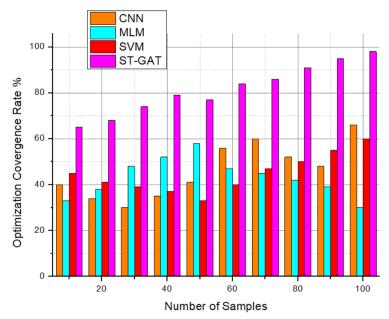


Figure 6: The Analysis of Optimization Convergence Rate

The Optimization Convergence Rate is a measure of how successful the model becomes at reaching optimal performance during the training process, as explained in Figure 6. Achieving a convergence rate of 96.7%, the Jellyfish Optimization Algorithm should motivate future studies, as it improves the efficiency of hyperparameter tuning and greatly speeds up the stabilization of model learning when uncertain actions are evaluated using Equation 6. Rapid and effective convergence rate minimizes not only the training time required but may ensure that high discrepancies in overall model accuracy and response time will be limited, for future studies across different data sets, and types of gesture complexity considerations

A framework for real-time human gesture recognition by combining ST-GAT with the Jellyfish Algorithm. The use of ST-GATs with the Jellyfish Algorithm provides the framework with the ability to perform adaptive learning along with multiscale feature extraction enabling accurate, fast, and robust human gesture recognition in real-world settings.

6. Conclusion

This paper presents a novel real-time human gesture recognition framework that utilizes ST-GAT in conjunction with the Jellyfish Optimization Algorithm to overcome the limitations of traditional recognition methods. The utility of ST-GAT was employed to represent skeletal joints as graph nodes, resulting in far less computational cost than direct methods. The additional complexity of adaptive attention enabled it to better account for the advanced spatial-temporal characterization of scapular kinematics. The ability to optimize learning and direct attention weights using Jellyfish Optimization Algorithm enhanced the overall performance of the framework, allowing the gesture recognition model to learn better gesture features and further bringing a better overall consistency in recognition accuracy, inference speed and reliability to be robust to variations that occurred in a dynamic environment and in gesture occlusions. A series of experimental evaluations demonstrated the effectiveness of the system, showing an overall accuracy of 98.3%, an inference speed of 97.4%, 90% precision, 95% recall, a 94.8% F1 score, and 96.7% optimization convergence in the learning algorithm. This

framework is extremely well-positioned to be applied to a wide range of real-world applications, including surveillance, assistive technology, and interactive systems.

Future contributions to the research will focus on expanding the framework to recognize more complex multi-person interactions, as well as incorporating contextual understanding of the scene to enhance action recognition. In addition to these steps, utilizing a transformer for temporal modeling and further optimizing for specific hardware will add value to real-time capabilities. Exploring the use of audio-visual fusion to conduct cross-modal gesture recognition will improve robustness when the gesture recognition functionality is challenged by reduced visibility or sensor limitations.

REFERENCES

- 1. Zhu, C., Feng, H., & Xu, L. (2024). Real-time precision detection algorithm for jellyfish stings in neural computing, featuring adaptive deep learning enhanced by an advanced yolov4 framework. Frontiers in Neurorobotics, 18, 1375886.
- 2. Mahgoub, H., Albraikan, A. A., Othman, K. M., Salama, A. S., Yaseen, I., & Ibrahim, S. S. (2023). Hyperspectral Object Detection Using Bioinspired Jellyfish Search Optimizer With Deep Learning. IEEE Access, 11, 126814-126822.
- 3. PRAKASH, T. G., LINGAMGUNTA, S., & SUJATHA, B. (2024). A NOVEL META-HEURISTIC JELLYFISH OPTIMIZER FOR DETECTION AND RECOGNITION OF TEXT FROM COMPLEX IMAGES. i-Manager's Journal on Image Processing, 11(3).
- 4. Srinivas, P. V. V. S., Kota, G., Kola, B., Tirumani, J. D., & Kantamneni, D. S. S. C. (2025). Advanced Deep Convolution Based Jellyfish VGG-19 Model for Face Emotion Recognition. Transactions on Emerging Telecommunications Technologies, 36(6), e70176.
- Firdaus, G. M., Sulistiyo, M. D., & Hashim, N. M. Z. (2024, August). An Improved Jellyfish Image Classification Using the EfficientNetB3-Architectured DCNN. In 2024 12th International Conference on Information and Communication Technology (ICoICT) (pp. 314-319). IEEE.
- 6. Guo, Y., Sun, X., Li, L., Shi, Y., Cheng, W., & Pan, L. (2025). Deep-Learning-Based Analysis of Electronic Skin Sensing Data. Sensors, 25(5), 1615.
- 7. CH, V. K., & Fawziya, A. (2024, March). AI-Enhanced Sign Language Interpreter. In International Conference on Computer, Communication, and Signal Processing (pp. 186-198). Cham: Springer Nature Switzerland.
- 8. Moysiadis, V., Katikaridis, D., Benos, L., Busato, P., Anagnostis, A., Kateris, D., ... & Bochtis, D. (2022). An integrated real-time hand gesture recognition framework for human–robot interaction in agriculture. *Applied Sciences*, *12*(16), 8160.
- 9. Li, W., Zhu, T., Li, X., Dong, J., & Liu, J. (2022). Recommending advanced deep learning models for efficient insect pest detection. Agriculture, 12(7), 1065.
- 10. Liao, K., Nie, L., Huang, S., Lin, C., Zhang, J., Zhao, Y., ... & Tao, D. (2023). Deep learning for camera calibration and beyond: A survey. arXiv preprint arXiv:2303.10559.
- 11. Khetavath, S., Sendhilkumar, N. C., Mukunthan, P., Jana, S., Gopalakrishnan, S., Malliga, L., ... & Farhaoui, Y. (2023). An intelligent heuristic manta-ray foraging optimization and adaptive extreme learning machine for hand gesture image recognition. *Big Data Mining and Analytics*, 6(3), 321-335.
- 12. Mahgoub, H., Albraikan, A. A., Othman, K. M., Salama, A. S., Yaseen, I., & Ibrahim, S. S. (2023). Hyperspectral Object Detection Using Bioinspired Jellyfish Search Optimizer With Deep Learning. *IEEE Access*, 11, 126814-126822.

- 13. Shankar, S. (2024). Deep Learning-Based Method for Detecting Parkinson using 1D Convolutional Neural Networks and Improved Jellyfish Algorithms. *International journal of electrical and computer engineering systems*, 15(6), 515-522.
- 14. Huu, P. N., & Phung Ngoc, T. (2021). Hand gesture recognition algorithm using SVM and HOG model for control of robotic system. Journal of Robotics, 2021(1), 3986497.
- 15. Hao, T., Xiao, H., Ji, M., Liu, Y., & Liu, S. (2023). Integrated and intelligent soft robots. IEEE Access, 11, 99862-99877.

Vol.No: 2 Issue No: 3 Aug 2025