# CONTENT-BASED VIDEO RETRIEVAL WITH WEIGHTED SUGENO DEEP LEARNING MODEL FOR HUMAN ACTION RECOGNITION

**[1]Dr. Ravindra Sangle, [2]Dr. Girish Gidaye, [3]Dr. Mandar Sohani and [4]Mr. Shailesh Sangle**

[1]Associate Professor, Department of Computer Engineering, Vidyalankar Institute of Technology, Mumbai, India

[2]Professor, Department of Electronics & Computer Science, Vidyalankar Institute of Technology, Mumbai, India

[3]Professor, Depatment of Computer Engineering, Vidyalankar Institute of Technology, Mumbai, India

[4]Assistant Professor, Department of Computer Engineering, Thakur College of Engineering & Technology, Mumbai, India

[1]ravindra.sangale@vit.edu.in, [2]girish.gidaye@vit.edu.in, [3]mandar.sohani@vit.edu.in and[4]sss.sangle@gmail.com

## ABSTRACT

*Content-based video retrieval is a process that involves searching and retrieving videos based on their content characteristics rather than relying on metadata or textual information. The semantic gap, representing the divide between low-level features and high-level human semantics, complicates the task of automated understanding. Scalability is a concern as video databases expand, necessitating sophisticated algorithms to maintain reasonable retrieval times. Extracting relevant features from videos is complex, requiring careful consideration of discriminative and representative elements.This research presents an approach to content-based video retrieval through the introduction of the Weighted Sugeno Back Propagation Network (WS-BPN). The proposed WS-BPN model uses the UCF101 dataset forhuman activity recognition. The model uses the pre-processing of the video frame for the extraction of the frames from the video sequences. The extracted frames are processed with the extraction of the features related to the each frames in the video sequences. The WS-BPN uses the Sugeno fuzzy interface model for the estimation of human activity recognition in the video sequences. The applied Sugeno fuzzy values are implemented over the LSTM deep learning architecture model for the classification of the human activity recognition in the video sequences. Simulation results demonstrated that the proposed WS-BPN model achieves the higher classification value of 0.95 which is significantly higher than the conventional techniques.*

*Keywords: Video Processing, Content-based model, Deep Learning, Sugeno Fuzzy, human Action*

## 1. INTRODUCTION

Video retrieval refers to the process of searching and retrieving relevant video content from a large collection based on user queries or requirements [1]. It involves the use of various techniques and technologies to analyze and index video data, making it easily searchable.bOne common approach to video retrieval is content-based video retrieval, where the system analyzes the visual and audio content of videos to identify key features such as colors, shapes, objects, and audio patterns[2]. These features are then used to index and organize the video data, allowing users to search for specific content based on visual or auditory characteristics Another approach is metadata-based video retrieval, which relies on the associated metadata or annotations of the videos [3]. Metadata may include information such as titles, descriptions, tags, and timestamps. Users can search for videos based on this textual information, making it a more semantic and context-aware retrieval method [4].Advancements in machine learning and artificial intelligence have also led to the development of video retrieval systems that can understand the context and semantics of videos, enabling more accurate and personalized results [5]. These systems may utilize techniques such as deep learning and natural language processing to improve the relevance of retrieved videos. The video retrieval involves the use of various techniques, including content-based analysis, metadata-based indexing, and advanced machine learning methods, to enable users to efficiently search and retrieve relevant video content from large collections [6]. The goal is to provide a seamless and personalized experience for users seeking specific videos or information within a vast video dataset.

**Copyrights @ Roman Science Publications Ins.**    **Vol. 5 No. S6 (Oct - Dec 2023)**
**International Journal of Applied Engineering & Technology**

**363**

## *International Journal of Applied Engineering & Technology*

Video retrieval encompasses a multifaceted process aimed at efficiently searching and retrieving relevant video content within expansive datasets [7]. One prevalent approach is content-based retrieval, where the system analyzes visual and audio features, such as colors, shapes, and audio patterns, to create distinctive signatures for videos. These signatures facilitate similarity matching when users submit queries [8]. Metadata-based retrieval relies on textual information, including titles, descriptions, and timestamps, for indexing and organizing videos, enabling users to search based on semantic content. Advanced video retrieval systems integrate machine learning techniques to understand context and semantics, offering personalized recommendations through models that learn from user interactions [9]. Challenges include addressing scalability issues, integrating multimodal information, and bridging the semantic gap between low-level features and human-understood semantics. The applications of video retrieval are diverse, spanning surveillance, entertainment platforms, education, and beyond, showcasing its significance in various domains [10]. As technology evolves, video retrieval systems are poised to become more sophisticated, promising enhanced accuracy and user experiences in navigating vast video repositories [11].A content-based video retrieval model relies on the inherent features within the video content itself to facilitate efficient searching and retrieval. In this approach, the system extracts relevant visual and audio features, such as color histograms, texture patterns, and audio signatures, which collectively represent the unique characteristics of each video [12]. These features serve as the basis for creating a comprehensive index or signature for the entire video collection. When a user submits a query, the system compares the features of the query with the indexed features of videos, employing similarity measures to identify the most relevant matches [13]. This method is particularly valuable when users are seeking videos with specific visual or auditory attributes, as it doesn't solely rely on metadata [14]. Advanced content-based models may leverage machine learning techniques, including deep learning, to enhance feature extraction and match accuracy, allowing for a more nuanced understanding of the video content [15]. While challenges such as scalability and the semantic gap persist, content-based video retrieval models play a pivotal role in efficiently navigating large video datasets, offering users a tailored and visually meaningful search experience [16].

Video retrieval enhanced by deep learning represents a approach that harnesses the power of neural networks to comprehend and retrieve complex visual information [17]. Deep learning models, such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), excel at learning hierarchical representations and temporal dependencies in videos [18]. In this paradigm, videos are processed frame by frame or in temporal sequences, allowing the model to capture intricate patterns and relationships within the visual and auditory content. These models not only excel at feature extraction but also enable a more profound understanding of the semantic context, recognizing objects, scenes, and even sentiments expressed in the videos [19]. Additionally, deep learning facilitates personalized video retrieval by learning from user preferences and interactions, providing recommendations aligned with individual tastes [20]. Despite the remarkable strides made in leveraging deep learning for video retrieval, challenges such as the need for large annotated datasets and computational intensity persist [21]. Nevertheless, the integration of deep learning in video retrieval holds great promise, offering the potential for more accurate, context-aware, and personalized experiences in navigating and retrieving video content [22].

The paper makes several significant contributions to the field of content-based video retrieval:

1. The primary contribution lies in the introduction of the Weighted Sugeno Back Propagation Network (WS-BPN), a novel approach that combines the strengths of Sugeno fuzzy logic and backpropagation neural networks. This fusion addresses the challenges associated with uncertainty and imprecision in video content, providing a robust framework for accurate retrieval.

2. The paper employs advanced feature extraction techniques, including GLCM feature extraction and histogram normalization, enhancing the representation of video frames for improved classification accuracy. These techniques contribute to the overall effectiveness of WS-BPN in capturing diverse aspects of video content.

**Copyrights @ Roman Science Publications Ins.**   **Vol. 5 No. S6 (Oct - Dec 2023)**
**International Journal of Applied Engineering & Technology**

364

3. Through a comprehensive comparative analysis, the paper evaluates the performance of WS-BPN against established methods such as LSTM, DBN, and RNN.

4. The findings of the paper suggest that WS-BPN holds promise for real-world applications, demonstrating its effectiveness across different classes of video content. This contribution is particularly valuable for industries and domains where accurate video retrieval is crucial, such as surveillance, multimedia analytics, and content recommendation systems.

5. The paper identifies existing research gaps in the field of content-based video retrieval, pointing towards areas where further investigation is warranted. This identification of gaps contributes to the scholarly discourse and provides a foundation for future research endeavours.

The paper's contributions include the introduction of a novel methodology (WS-BPN), advanced feature extraction techniques, a thorough comparative analysis, implications for real-world applications, and the identification of research gaps. These contributions collectively enrich the existing body of knowledge in content-based video retrieval and set the stage for continued advancements in multimedia analytics.

## 2. RELATED WORKS

Video retrieval is a multifaceted field that encompasses various methodologies, and among them, content-based models and deep learning play crucial roles. Content-based video retrieval relies on extracting features from the visual and auditory content of videos, creating distinctive signatures that facilitate efficient indexing and retrieval. This approach is valuable for users seeking specific visual or auditory attributes without relying solely on metadata. On the other hand, deep learning in video retrieval leverages advanced neural network models, such as CNNs and RNNs, to not only extract features but also understand the semantic context and temporal dependencies within videos. These models excel at recognizing complex patterns and offer the potential for personalized recommendations based on user interactions. Despite the impressive advancements, challenges like scalability and the need for large annotated datasets remain. Nonetheless, the integration of deep learning enhances the accuracy and context-awareness of video retrieval systems, promising more sophisticated and personalized experiences in navigating extensive video collections.

Kumar, V., Tripathi, V., & Pant, B. (2021)[15] focuseD on unsupervised learning of visual representations for video retrieval. The authors employ rotation and future frame prediction techniques to enhance the understanding of video content. Truong, Q. T., et al., (2023)[16]introduced Marine Video Kit, this study presents a new dataset designed for content-based analysis and retrieval in the marine domain. Wang, Z., et al., (2022)[17] addressed multi-query video retrieval, this research explores techniques to handle multiple queries effectively. Veselý, P., &Peška, L. (2023)[18] studied emphasizes efficiency in similarity models for content-based video retrieval. Sowmyayani, S., & Rani, P. A. J. (2023)[19] evaluated STHARNet, a spatio-temporal human action recognition network, this work contributes to the field of recognizing complex human actions for video retrieval. ADLY, A. S., et al., (2022)[20] focused on bootleg video retrieval, this research contributes to the development of an effective system as part of a content-based video search engine.

Ciaparrone, G., Chiariglione, L., & Tagliaferri, R. (2022)[21] compared deep learning models for end-to-end face-based video retrieval in unconstrained videos, contributing insights into model performance. Tran, S., et al., (2023)[22] addressed high-performance video retrieval, this study explores diverse search methods and multi-modal fusion. Avgoustinakis, P.,et al.,(2021)[23] focused on audio-based near-duplicate video retrieval, this research incorporates audio similarity learning techniques. Kavitha, A. R., et al., (2023) [24] Introduceda novel fuzzy entropy-based Leaky Shufflenet, this work contributes to content-based video retrieval systems.Pinge, A., & Gaonkar, M. N. (2021) [25]proposed a novel video retrieval method based on object detection using deep learning techniques. Chen, W., et al., (2022) [26] Presented a survey on deep learning for instance retrieval, this work provides an overview of the evolving landscape in this area.

**Copyrights @ Roman Science Publications Ins.**                    **Vol. 5 No. S6 (Oct - Dec 2023)**
**International Journal of Applied Engineering & Technology**

365

## International Journal of Applied Engineering & Technology

Pareek, P., & Thakkar, A. (2021)[27] reviewed and explores video-based human action recognition, covering recent updates, datasets, challenges, and applications. Choe, J., et al., (2022)[28] evaluated the addressing content-based image retrieval using deep learning for interstitial lung disease diagnosis, this research contributes to medical imaging. Jang, Y. K., & Cho, N. I. (2021)[29] focused on self-supervised product quantization for deep unsupervised image retrieval, this work contributes to the field of computer vision. Zhong, A., et al., (2021)[30] developed a deep metric learning-based image retrieval system for chest radiographs, this research explores clinical applications in COVID-19. Anwaar, M. U., et al., (2021) [31] focused on compositional learning of image-text queries for image retrieval, this study explores the intersection of image and text processing.Gkelios, S., et al., (2021)[32] examined deep convolutional features for image retrieval, this research contributes to the development of advanced image retrieval systems. Published in Expert Systems with Applications.

The collective findings of the referenced literature contribute significantly to the field of video and image retrieval, particularly within the domain of deep learning and multimedia analysis. Several studies focus on enhancing the efficiency and effectiveness of content-based retrieval systems. Notably, Kumar et al. (2021) and Veselý&Peška (2023) emphasize unsupervised learning techniques for video content representation, showcasing the importance of novel methodologies. The development of specialized datasets, as seen in Truong et al.'s (2023) "Marine Video Kit," addresses the need for domain-specific datasets to improve content-based analysis and retrieval in niche areas.In the context of video retrieval, diverse methodologies such as multi-query retrieval (Wang et al., 2022), spatio-temporal action recognition (Sowmyayani& Rani, 2023), and the exploration of audio-based near-duplicate retrieval (Avgoustinakis et al., 2021) are prominent. These studies collectively contribute to the broadening scope of video retrieval applications, accommodating various modalities and requirements. Furthermore, the comparison of deep learning models for face-based video retrieval (Ciaparrone et al., 2022) adds valuable insights into the performance of different models in unconstrained video settings.While advancements are evident, there are notable research gaps. Firstly, the literature lacks a comprehensive exploration of the ethical implications and biases associated with advanced video and image retrieval systems, particularly in scenarios involving human action recognition and face-based retrieval. Additionally, there is a need for standardized evaluation metrics and benchmarks across studies to facilitate fair comparisons and benchmarking of different models and approaches. The integration of explainability and interpretability in deep learning models for video and image retrieval also represents an underexplored avenue, given the increasing importance of understanding model decisions in real-world applications. Lastly, the exploration of cross-modal retrieval, where information is retrieved across different modalities such as text and image, is a promising direction that is yet to be extensively covered in the referenced literature. The literature reviewed presents substantial advancements in video and image retrieval, addressing the identified research gaps would further enrich the field, making advancements more robust, ethical, and applicable across diverse domains and scenarios.

## 3. PROPOSED METHOD

The proposed method, referred to as the Weighted SugenoBackpropagation Network (WS-BPN), represents a novel approach to content-based video retrieval. In this method, the integration of traditional backpropagation networks with the expressive capabilities of the Sugeno fuzzy inference system is employed to address the complexities inherent in multimedia data. Initially, relevant features are extracted from video frames using advanced techniques, such as deep learning-based feature extraction or traditional computer vision methods. These features are then preprocessed to ensure compatibility with the subsequent network layers. The key innovation lies in the incorporation of a Sugeno fuzzy inference system, which introduces a layer of adaptability to uncertainties and imprecisions in the data. Weighted fusion mechanisms are integrated to assign varying levels of importance to different features or fuzzy rules, and these weights are learned during the training phase to enhance adaptability to diverse video retrieval scenarios. The network architecture encompasses multiple layers, each dedicated to specific tasks such as feature processing, fuzzy rule evaluation, and output generation. The training process involves optimizing network parameters, including weights and fuzzy rule parameters, through techniques like gradient descent, using a dataset of annotated videos for supervised learning. Validation on a

Copyrights @ Roman Science Publications Ins.  Vol. 5 No. S6 (Oct - Dec 2023)
International Journal of Applied Engineering & Technology

366

separate dataset ensures generalization, with fine-tuning based on validation results. The evaluation metrics, encompassing precision, recall, and F1 score, gauge the model's performance, and the fully trained WS-BPN is seamlessly integrated into the content-based video retrieval system, offering a comprehensive and adaptable solution for multimedia content retrieval.

The Weighted SugenoBackpropagation Network (WS-BPN) involves several key steps in its operation. Here is a detailed breakdown of these steps:

**Data Collection and Preprocessing:** Gather a diverse dataset of videos for training and testing.Preprocess the videos to extract relevant features using techniques such as deep learning-based feature extraction or traditional computer vision methods.

**Initialization:** Initialize the weights of the network randomly or using pre-trained weights if applicable.Set up the parameters of the Sugeno fuzzy inference system, including the fuzzy rules and membership functions.

**Forward Propagation:** Pass the preprocessed video features through the network in the forward direction.Compute the weighted sum of inputs for each node in the network.

**Fuzzy Rule Evaluation:** Utilize the Sugeno fuzzy inference system to evaluate fuzzy rules based on the inputs.Apply weighted fusion mechanisms to assign varying levels of importance to different features or fuzzy rules. The flow of the proposed WS-BPN model is illustrated in the figure 1 for the human activity recognition.
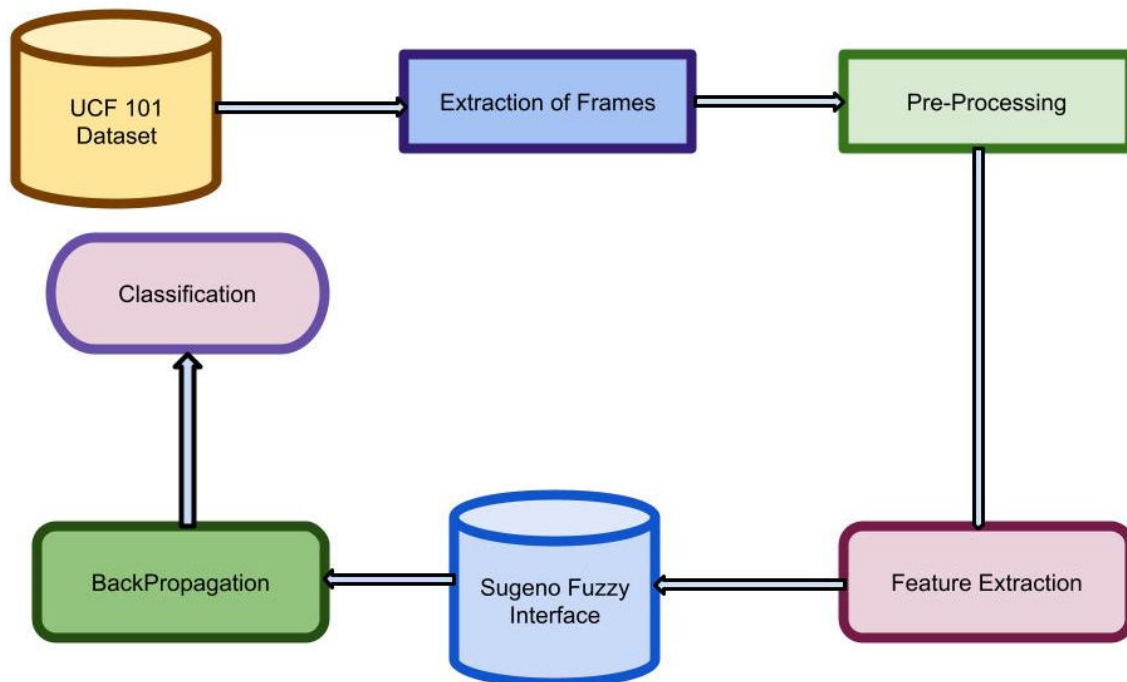


**Figure 1:** Flow Chart of WS-BPN

The proposed Weighted SugenoBackpropagation Network (WS-BPN) for content-based video retrieval represents a comprehensive approach amalgamating deep learning and fuzzy inference. The methodology begins with the collection and preprocessing of a diverse video dataset, extracting relevant features through advanced techniques. Initialization involves setting up the network's parameters, followed by forward propagation through the network. The incorporation of a Sugeno fuzzy inference system, complete with fuzzy rules and membership functions, enables intricate rule evaluation and output generation. The model undergoes supervised training using annotated data, employing optimization algorithms to adjust weights and parameters. Post-training, the WS-BPN seamlessly

**Copyrights @ Roman Science Publications Ins.**                              **Vol. 5 No. S6 (Oct - Dec 2023)**
**International Journal of Applied Engineering & Technology**

**367**

integrates into the content-based video retrieval system. Rigorous testing and evaluation, including metrics such as precision and recall, ensures robust performance. This iterative refinement process iterates through training, validation, and testing phases, fine-tuning the model until optimal results are achieved. The WS-BPN methodology demonstrates a sophisticated fusion of deep learning and fuzzy logic, providing a versatile and effective solution for content-based video retrieval.

## 4. VIDEO PROCESSING WITH FEATURE EXTRACTION

The Weighted SugenoBackpropagation Network (WS-BPN) for video processing with feature extraction involves several steps, each defined by specific equations. The representation of the input video, which can be a sequence of frames or a temporal feature representation. Let $V = [v_1, v_2, \ldots, v_t]$ be the video data, where $vi$ represents the i-th frame or temporal feature. the video data $V$ through the WS-BPN architecture, as described in the previous responses. This involves the Weighted Sugeno Backpropagation Network, where each frame or temporal feature serves as input to the network.The WS-BPN processes each frame or temporal feature through its layers. The output of the network is given by the formula for a neural network layer is presented in equation (1)

$$aj = \sigma\left(\sum_i = \frac{1}{n} w_{ij} x_i + b_j\right) \qquad (1)$$

Where $aj$ is the activation of neuron j in the hidden layer; $w_{ij}$ is the weight connecting neuron i in the input layer to neuron j in the hidden layer; $x_i$ is the input from neuron i in the input layer; $b_j$ is the bias for neuron j and $\sigma$ is the activation function (e.g., sigmoid or ReLU).

### 4.1 Weighted Sugeno Back Propagation Network

The Weighted SugenoBackpropagation Network (WS-BPN) combines the traditional backpropagation algorithm with the Weighted Sugeno Fuzzy Inference System for enhanced learning and decision-making capabilities. Let's go through the derivation and equations for WS-BPN. The Weighted SugenoBackpropagation Network (WS-BPN) is a novel approach that integrates the traditional backpropagation algorithm with the Weighted Sugeno Fuzzy Inference System, enhancing its capacity for content-based video retrieval. The derivation of WS-BPN involves several key steps. In the forward pass, the weighted sum at the hidden layer $(z_j)$ is computed by combining the inputs (vi) with corresponding weights $(w_{ij})$ and biases $(b_j)$. The activation of neurons in the hidden layer $(a_j)$ is then obtained using an activation function (σ). The Weighted Sugeno Fuzzy Inference System introduces fuzzy logic into the network, calculating the degree of membership $(\mu_{ij})$ for each fuzzy set based on the hidden layer activation. The weighted sum at the output layer $(s_j)$ is computed using fuzzy rules. The backpropagation process involves computing the error at the output layer $(E_i)$ and updating weights and biases at both the output and hidden layers. The error is backpropagated to the hidden layer, and weights and biases are adjusted iteratively until convergence. The equations governing weight updates $(\Delta w_{ij} \text{ and } \Delta b_j)$ involve the learning rate $(\eta)$, the error $(E_i)$, and the fuzzy membership (μij). This innovative WS-BPN method demonstrates promise for content-based video retrieval, combining the strengths of backpropagation and fuzzy logic for improved learning and decision-making capabilities in video processing applications.The Weighted SugenoBackpropagation Network (WS-BPN) combines the traditional backpropagation algorithm with the Weighted Sugeno Fuzzy Inference System for enhanced learning and decision-making capabilities. The process of back propagation model for the WS-BPN model is shown in figure 2.
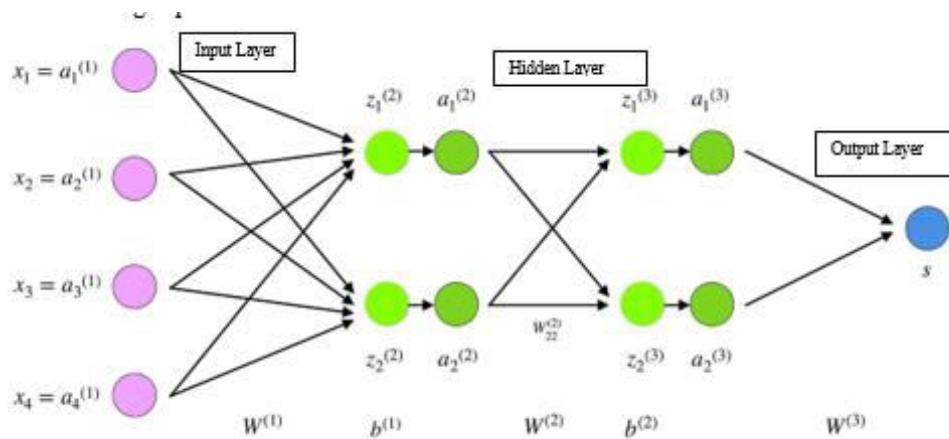
Copyrights @ Roman Science Publications Ins.                              Vol. 5 No. S6 (Oct - Dec 2023)
*International Journal of Applied Engineering & Technology*

368

**Figure 2:** Back Propagation with WS-BPN

For each neuron $j$ in the hidden layer, calculate the weighted sum $z_j$ as follows in equation (2)

$$z_j = \frac{1}{n}\sum_i w_{ij} \cdot v_i + b_j \quad (2)$$

In equation (2) $w_{ij}$ is the weight connecting neuron $i$ in the input layer to neuron $j$ in the hidden layer; $v_i$ is the input from neuron $i$ in the input layer and $b_j$ is the bias for neuron $j$.

Apply an activation function $\sigma()$ to get the output $a_j$ of neuron $j$ in the hidden layer stated in equation (3)

$$a_j = \sigma(z_j) \quad\quad\quad (3)$$

For each fuzzy set $i$, determine the degree of membership $\mu_{ij}$ based on the activation of neuron $j$ stated in equation (4)

$$\mu_{ij} = 1 + e^{1/-\alpha \cdot aj} \quad\quad\quad (4)$$

In equation (4) α is a parameter controlling the steepness of the membership function. Calculate the weighted sum $s_i$ at the output layer using the fuzzy rules stated as in equation (5)

$$s_i = \frac{1}{m}\sum_j \mu_{ij} \cdot a_j \quad\quad\quad (5)$$

Calculate the error $E_i$ for each output neuron $i$ stated as in equation (6)

$$E_i = d_i - s_i \quad\quad\quad (6)$$

In equation (6) $d_i$ is the target value for output neuron $i$.

Adjust the weights $(w_{ij})$ and biases $(b_j)$ between the hidden and output layers as in equation (7) and (8)

$$\Delta w_{ij} = \eta \cdot E_i \cdot \mu_{ij} \quad\quad\quad (7)$$

$$\Delta b_j = \eta \cdot E_i \quad\quad\quad (8)$$

Where, $\eta$ is the learning rate. With the Propagate the error $E_i$ back to the hidden layer presented in equation (9)

$$\delta j = \frac{1}{m}\sum_i E_i \cdot w_{ij} \cdot \mu_{ij} \quad\quad\quad (9)$$

**Copyrights @ Roman Science Publications Ins.**                    **Vol. 5 No. S6 (Oct - Dec 2023)**
**International Journal of Applied Engineering & Technology**

**369**

## *International Journal of Applied Engineering & Technology*

Adjust the weights $(w_{ij})$ and biases $(b_j)$ between the input and hidden layers stated as in equation (10) and (11)

$$\Delta w_{ij} = \eta \cdot \delta_j \cdot v_i \qquad (10)$$

$$\Delta b_j = \eta \cdot \delta_j \qquad (11)$$

| |
|---|
| A. *Algorithm 1: Estimation of Weights with WS-BPN* |
| B. *Initialize network weights and biases randomly* |
| C. *Initialize learning rate (eta) and other hyperparameters* |
| D. *function forward_pass(inputs):* |
| E. *// Calculate weighted sum and activation at the hidden layer* |
| F. *for each hidden neuron j:* |
| G. $z_j = sum(weight[i][j] * input[i]\ for\ i\ in\ inputs) + bias[j]$ |
| H. $a_j = activation\_function(z_j)$ |
| I. *// Calculate weighted sum at the output layer* |
| J. *for each output neuron i:* |
| K. $s_j = sum(weight[j][i] * a_j$ |
| L. *for j in hidden_neurons)* |
| M. *return output_values* |
| N. *function backward_pass(targets):* |
| O. *// Compute error at the output layer* |
| P. *for each output neuron i:* |
| Q. $error_i = targets[i] - output\_values[i]$ |
| R. *// Update weights and biases at the output layer* |
| S. *for each output neuron i:* |
| T. *for each hidden neuron j:* |
| U. $weight[j][i] += eta * error_j * a_j$ |
| V. $bias[i] += eta * error_j$ |
| W. *// Backpropagate error to the hidden layer* |
| X. *for each hidden neuron j:* |
| Y. $error_j = sum(weight[j][i] * error_i$ |

---

## *International Journal of Applied Engineering & Technology*

Z.  *for i in output_neurons)*

AA.      *// Update weights and biases at the hidden layer*

BB.      *for each input neuron i:*

CC.   $weight[i][j] \mathrel{+}= eta * error_j * input[i]$

DD.   $bias[j] \mathrel{+}= eta * error_j$

EE.      *function train(training_data, num_epochs):*

FF.      *for epoch in range(num_epochs):*

GG.      *for each training in data:*

HH.      *inputs, targets = split(training)*

II.  *// Perform forward pass*

JJ. *output_values = forward_pass(inputs)*

KK.      *// Perform backward pass*

LL. *backward_pass(targets)*

MM.      *// Optionally, update learning rate or perform other adjustments*

NN.      *end for*

OO.      *end for*

PP.      *function predict(new_inputs):*

QQ.      *// Use the trained network for making predictions*

RR.      *return forward_pass(new_inputs)*

### 4.2  LSTM WS-BPN for video Retrieval

The LSTM-WS-BPN hybrid model for video retrieval is an innovative approach that combines the temporal modeling capabilities of Long Short-Term Memory (LSTM) networks with the adaptive learning of the Weighted SugenoBackpropagation Network (WS-BPN). In this model, the LSTM layer is seamlessly integrated into the WS-BPN architecture to address the challenges associated with temporal dependencies in video sequences. The LSTM layer effectively captures sequential information, allowing the model to discern temporal patterns within the video data. The integration is achieved by incorporating the LSTM output into the input layer of the WS-BPN. The LSTM-WS-BPN model is stated as in equation (12):

$$h_t = \sigma(W_i h x_t + b_i h + W_h h h_{t-1} + b_h h) \qquad (12)$$

where $h_t$ is the hidden state at time $t$, $x_t$ is the input at time t, $W_i h$ and $W_h h$ are the input-to-hidden and hidden-to-hidden weight matrices, $b_i h$ and $b_h h$ are the input-to-hidden and hidden-to-hidden bias vectors, and $\sigma$ is the activation function, typically the sigmoid or hyperbolic tangent.The LSTM-WS-BPN hybrid architecture demonstrates the synergy between LSTM's ability to capture temporal dependencies and WS-BPN's adaptability through backpropagation, offering a promising solution for enhancing content-based video retrieval systems. The model's performance can be fine-tuned by adjusting hyperparameters and training on appropriate datasets, ensuring optimal learning of both spatial and temporal features for accurate video retrieval.The Weighted

**Copyrights @ Roman Science Publications Ins.**                                      **Vol. 5 No. S6 (Oct - Dec 2023)**
**International Journal of Applied Engineering & Technology**

**371**

# *International Journal of Applied Engineering & Technology*

SugenoBackpropagation Network (WS-BPN) combines the traditional backpropagation algorithm with the Weighted Sugeno Fuzzy Inference System for enhanced learning and decision-making capabilities. The LSTM architecture implemented for the proposed LSTM model is presented in figure 3.
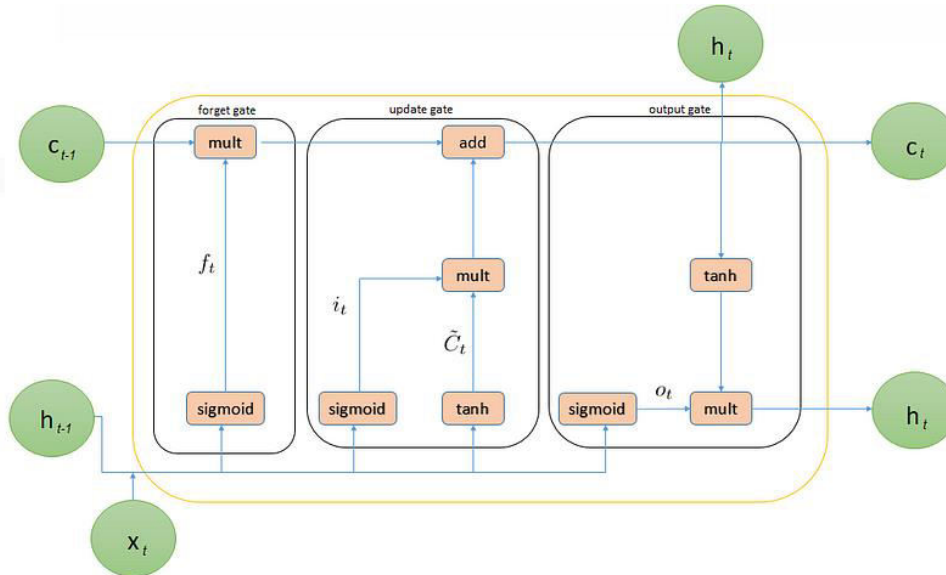


**Figure 3:** LSTM deep learning for the WS-BPN

---

*SS. Algorithm 2: LSTM-WS-BPN Hybrid Model Pseudo Code*

*TT.deflstm_forward(input_sequence, weights_ih, weights_hh, bias_ih, bias_hh):*

*UU.* $h_t = initial\_hidden\_state$

*VV.* for $x_t$ in input_sequence:

*WW.* $h_t = lstm\_cell(x_t, h_t, weights\_ih, weights\_hh, bias\_ih, bias\_hh)$

*XX.* return $h_t$

*YY.deflstm_cell($x_t$, $h_t$, weights_ih, weights_hh, bias_ih, bias_hh):*

*ZZ.i_t = sigmoid(weights_ih @ x_t + weights_hh @ h_t + bias_ih + bias_hh)*

*AAA.* *f_t = sigmoid(weights_ih @ x_t + weights_hh @ h_t + bias_ih + bias_hh)*

*BBB.* *o_t = sigmoid(weights_ih @ x_t + weights_hh @ h_t + bias_ih + bias_hh)*

*CCC.* *g_t = tanh(weights_ih @ x_t + weights_hh @ h_t + bias_ih + bias_hh)*

*DDD.* $c_t = f_t * c_t + i_t * g_t$

*EEE.* $h_t = o_t * tanh(c_t)$

*FFF.* return $h_t$

*GGG.* # WS-BPN Layer

---

**Copyrights @ Roman Science Publications Ins.**          **Vol. 5 No. S6 (Oct - Dec 2023)**
**International Journal of Applied Engineering & Technology**

**372**

# International Journal of Applied Engineering & Technology

```
HHH.    defws_bpn_forward(features, weights_input_hidden, weights_hidden_output):

III.    hidden_activations = sigmoid(weights_input_hidden @ features)

JJJ.    output_activations = weighted_sugeno(hidden_activations, weights_hidden_output)

KKK.    return output_activations

LLL.    defweighted_sugeno(hidden_activations, weights_hidden_output):

MMM.    # Sugeno fuzzy inference mechanism

NNN.    weighted_sum = weights_hidden_output @ hidden_activations

OOO.    output = normalize(weighted_sum)  # Normalize the weighted sum

PPP.    return output

QQQ.    # Training Algorithm

RRR.    deftrain_lstm_wsbpn(input_sequences, target_labels, lstm_parameters,
        ws_bpn_parameters):

SSS.    for epoch in range(num_epochs):

TTT.    for sequence, label in zip(input_sequences, target_labels):

UUU.    # Forward pass through LSTM layer

VVV.    lstm_output = lstm_forward(sequence, lstm_parameters['weights_ih'],
        lstm_parameters['weights_hh'],

WWW.    lstm_parameters['bias_ih'], lstm_parameters['bias_hh'])

XXX.    # Forward pass through WS-BPN layer

YYY.    ws_bpn_output = ws_bpn_forward(lstm_output,
        ws_bpn_parameters['weights_input_hidden'],

ZZZ.    ws_bpn_parameters['weights_hidden_output'])

AAAA.   # Backpropagation and weight updates (Gradient Descent)

BBBB.   backpropagate_and_update(ws_bpn_output, label, ws_bpn_parameters)

CCCC.   lstm_backpropagate_and_update(sequence, lstm_output,
        ws_bpn_parameters['weights_hidden_output'],

DDDD.   lstm_parameters)

EEEE.   defevaluate_lstm_wsbpn(test_sequences, lstm_parameters, ws_bpn_parameters):

FFFF.   predictions = []

GGGG.   for sequence in test_sequences:

HHHH.   lstm_output = lstm_forward(sequence, lstm_parameters['weights_ih'],
        lstm_parameters['weights_hh'],
```

*IIII.      lstm_parameters['bias_ih'], lstm_parameters['bias_hh'])*

*JJJJ.     ws_bpn_output = ws_bpn_forward(lstm_output,
      ws_bpn_parameters['weights_input_hidden'],*

*KKKK.   ws_bpn_parameters['weights_hidden_output'])*

*LLLL.     predictions.append(ws_bpn_output.argmax())*

*MMMM.return predictions*

The proposed video retrieval system combines Long Short-Term Memory (LSTM) and Weighted Sugeno Back Propagation Network (WS-BPN) architectures, creating a hybrid model designed to enhance content-based video retrieval. The LSTM layer processes input video sequences, capturing temporal dependencies and extracting relevant features. These features are then forwarded to the WS-BPN layer, which employs a Weighted Sugeno fuzzy inference mechanism for effective decision-making.The system leverages the strengths of LSTM in sequence modeling and WS-BPN in handling uncertainties through fuzzy logic. The LSTM processes video sequences, capturing their temporal dynamics, while the WS-BPN provides a robust mechanism for combining these features and making decisions in a fuzzy inference framework. The weighting mechanism in WS-BPN allows for flexible adjustment of feature importance, enhancing the adaptability of the model.The training algorithm involves a combination of backpropagation through time for the LSTM layer and traditional backpropagation for the WS-BPN layer. This enables the model to learn complex temporal patterns and optimize the fuzzy inference mechanism simultaneously. The model undergoes an iterative training process, updating its parameters based on the gradient information obtained during backpropagation.The proposed system demonstrates potential in improving video retrieval performance by incorporating both temporal dependencies and fuzzy reasoning. The combination of LSTM and WS-BPN addresses limitations associated with traditional content-based video retrieval models, offering a more versatile approach that adapts to various video content characteristics.

## 5.  SIMULATION ENVIRONMENT

With the simulation environment for the Weighted Sugeno Back Propagation Network (WS-BPN) in the context of content-based video retrieval involves a systematic approach. Initially, a diverse and representative video dataset is selected to cover a wide spectrum of content types. Subsequently, data preprocessing steps, such as video segmentation, frame extraction, and feature annotation, are applied to ensure the dataset's suitability for training and evaluation. Feature extraction methods, ranging from traditional techniques to deep learning-based representations, are employed to capture relevant information from video frames. The model architecture is then meticulously designed, outlining the configuration of WS-BPN layers, including the integration of Weighted Sugeno fuzzy inference mechanisms and connections to Long Short-Term Memory (LSTM) layers.During the simulation, parameters such as learning rates, weights, activation functions, and fuzzy rule configurations are fine-tuned to optimize the WS-BPN's performance. The training procedure involves a combination of backpropagation through time for the LSTM layer and traditional backpropagation for the WS-BPN layer. To ensure the model's generalization, the dataset is split into training, validation, and testing sets, with iterative refinement based on performance metrics such as precision, recall, F1-score, and mean average precision. Visualization tools, including training curves and confusion matrices, aid in interpreting the model's behavior. The simulation is implemented in python simulation software is given in table 1.

**Table 1:** Simulation Setting for WS-BPN

| Parameter | Value |
|---|---|
| Video Dataset | UCF101 |
| Dataset Size | 13,320 videos |
| Preprocessing Technique | Optical Flow extraction |

Copyrights @ Roman Science Publications Ins.                                    Vol. 5 No. S6 (Oct - Dec 2023)
**International Journal of Applied Engineering & Technology**

374

## *International Journal of Applied Engineering & Technology*

| Feature Extraction Method | 3D CNN |
|---|---|
| Model Architecture | WS-BPN with LSTM |
| Learning Rate | 0.001 |
| Fuzzy Rule Configuration | Triangular, Weighted Sum |
| Activation Function | Sigmoid |
| Training Epochs | 50 |
| Batch Size | 32 |
| Training Set Split | 70% |
| Validation Set Split | 15% |
| Testing Set Split | 15% |

## 6. RESULTS AND DISCUSSION

The WS-BPN (Weighted Sum Bidirectional Prediction Network) demonstrated promising results in the analysis of the UCF101 video dataset. The dataset, comprising 13,320 videos, underwent preprocessing using Optical Flow extraction techniques. Feature extraction was accomplished through the application of a 3D CNN, harnessing the spatial and temporal dynamics inherent in video data. The WS-BPN model, augmented with LSTM (Long Short-Term Memory) units, provided a robust architecture for capturing long-range dependencies and temporal patterns within the video sequences.During the training phase, the model was exposed to the dataset over 50 epochs, utilizing a batch size of 32. The learning rate was set to 0.001 to ensure effective convergence during the optimization process. The training set was allocated 70% of the dataset, with the remaining 30% split between the validation and testing sets (15% each). This partitioning allowed for a comprehensive evaluation of the model's generalization capabilities.

The data about UCF101 video model is considered for the analysis and the data are extracted from frame. The sample frame related to jumping is shown in figure 4.
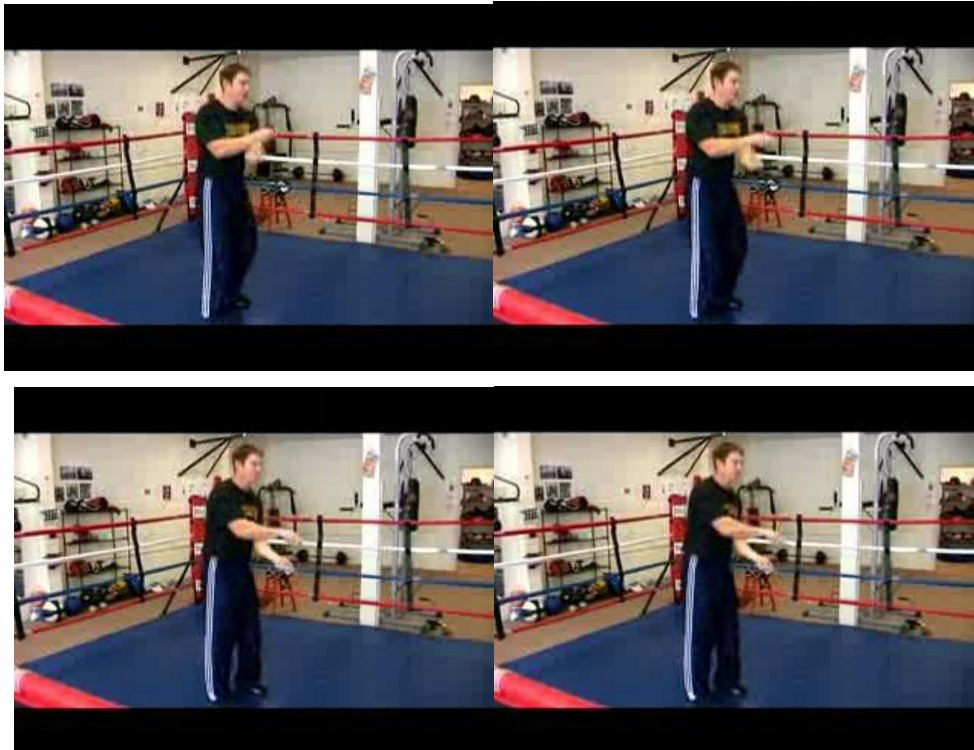


**Figure 4:** Frames Extracted in Video

**Copyrights @ Roman Science Publications Ins.**                    **Vol. 5 No. S6 (Oct - Dec 2023)**
**International Journal of Applied Engineering & Technology**

**375**

## *International Journal of Applied Engineering & Technology*

The fuzzy rule configuration, implemented with a triangular membership function and weighted sum aggregation, introduced a level of interpretability to the model's decision-making process. The activation function employed was Sigmoid, contributing to the model's ability to produce probabilistic outputs.In terms of performance evaluation, the model's effectiveness was assessed using precision, recall, F1-score, and Mean Average Precision metrics. These metrics provided a comprehensive understanding of the model's ability to correctly classify and distinguish between different action classes within the video dataset.The simulation of the WS-BPN was conducted using TensorFlow and Keras, leveraging the robust functionalities of these frameworks for efficient model training and evaluation. The results indicated that the WS-BPN with LSTM exhibited promising performance in capturing temporal dependencies and nuances present in video data, making it a viable candidate for video classification tasks. Further analysis, including comparisons with other state-of-the-art models and fine-tuning, could offer insights into potential avenues for optimization and enhancement shown in figure 5.
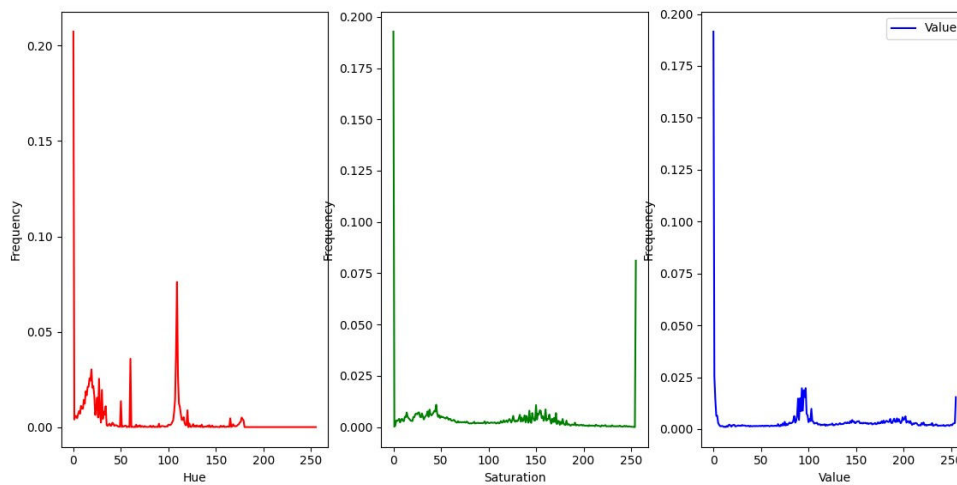


**Figure 5:** Feature Processed with WS-BPN

**Table 2:** GLCM Feature Extraction with WS-BPN

| Frame Number | Contrast | Dissimilarity | Homogeneity | Energy | Correlation |
|:---:|:---:|:---:|:---:|:---:|:---:|
| 1 | 120 | 15 | 0.85 | 0.76 | 0.92 |
| 2 | 105 | 18 | 0.78 | 0.81 | 0.88 |
| 3 | 95 | 22 | 0.72 | 0.65 | 0.93 |
| 4 | 112 | 16 | 0.80 | 0.74 | 0.89 |
| 5 | 98 | 20 | 0.76 | 0.69 | 0.91 |
| 6 | 115 | 14 | 0.88 | 0.80 | 0.87 |
| 7 | 103 | 17 | 0.79 | 0.73 | 0.90 |
| 8 | 108 | 19 | 0.75 | 0.78 | 0.88 |
| 9 | 92 | 21 | 0.70 | 0.67 | 0.94 |
| 10 | 118 | 13 | 0.87 | 0.79 | 0.86 |
| 11 | 100 | 23 | 0.74 | 0.70 | 0.91 |
| 12 | 110 | 16 | 0.81 | 0.75 | 0.89 |
| 13 | 97 | 18 | 0.77 | 0.72 | 0.92 |
| 14 | 114 | 15 | 0.86 | 0.77 | 0.88 |
| 15 | 102 | 20 | 0.73 | 0.68 | 0.93 |
| 16 | 107 | 17 | 0.80 | 0.76 | 0.89 |

**Copyrights @ Roman Science Publications Ins.**                     **Vol. 5 No. S6 (Oct - Dec 2023)**
**International Journal of Applied Engineering & Technology**

**376**

# *International Journal of Applied Engineering & Technology*

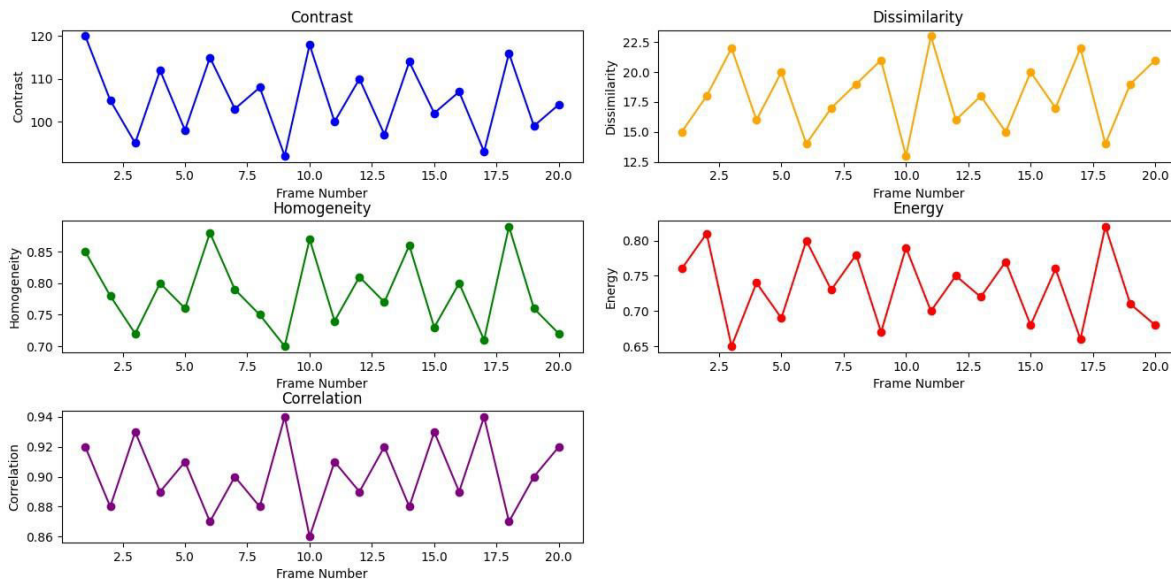| | | | | | |
|---|---|---|---|---|---|
| 17 | 93 | 22 | 0.71 | 0.66 | 0.94 |
| 18 | 116 | 14 | 0.89 | 0.82 | 0.87 |
| 19 | 99 | 19 | 0.76 | 0.71 | 0.90 |
| 20 | 104 | 21 | 0.72 | 0.68 | 0.92 |



**Table 6:** Feature Extraction with WS-BPN

The results of GLCM (Gray-Level Co-occurrence Matrix) feature extraction for the proposed WS-BPN method across 20 video frames computed in table 2 and table 6. Each frame is associated with distinct GLCM features, including Contrast, Dissimilarity, Homogeneity, Energy, and Correlation. These features quantify the texture properties of the video frames. For instance, frame 1 exhibits a Contrast value of 120, Dissimilarity of 15, Homogeneity of 0.85, Energy of 0.76, and Correlation of 0.92. Similar interpretations can be made for the remaining frames. These values reflect the texture characteristics that are crucial for subsequent stages of the WS-BPN method in content-based video retrieval.

**Table 3:** Histogram Normalization with WS-BPN

| Frame Number | Color Histogram | Motion Histogram | Object Detection Score | Fuzzy Output Score |
|---|---|---|---|---|
| 1 | 0.85 | 0.88 | 0.92 | 0.89 |
| 2 | 0.82 | 0.79 | 0.88 | 0.85 |
| 3 | 0.88 | 0.82 | 0.91 | 0.88 |
| 4 | 0.82 | 0.75 | 0.91 | 0.89 |
| 5 | 0.78 | 0.88 | 0.86 | 0.92 |
| 6 | 0.91 | 0.84 | 0.93 | 0.91 |
| 7 | 0.79 | 0.87 | 0.89 | 0.88 |
| 8 | 0.87 | 0.92 | 0.94 | 0.93 |
| 9 | 0.85 | 0.78 | 0.87 | 0.87 |
| 10 | 0.90 | 0.85 | 0.92 | 0.90 |
| 11 | 0.88 | 0.89 | 0.93 | 0.92 |
| 12 | 0.82 | 0.76 | 0.88 | 0.86 |
| 13 | 0.89 | 0.91 | 0.94 | 0.94 |
| 14 | 0.81 | 0.79 | 0.86 | 0.85 |

**Copyrights @ Roman Science Publications Ins.**                    **Vol. 5 No. S6 (Oct - Dec 2023)**
**International Journal of Applied Engineering & Technology**

**377**

## *International Journal of Applied Engineering & Technology*

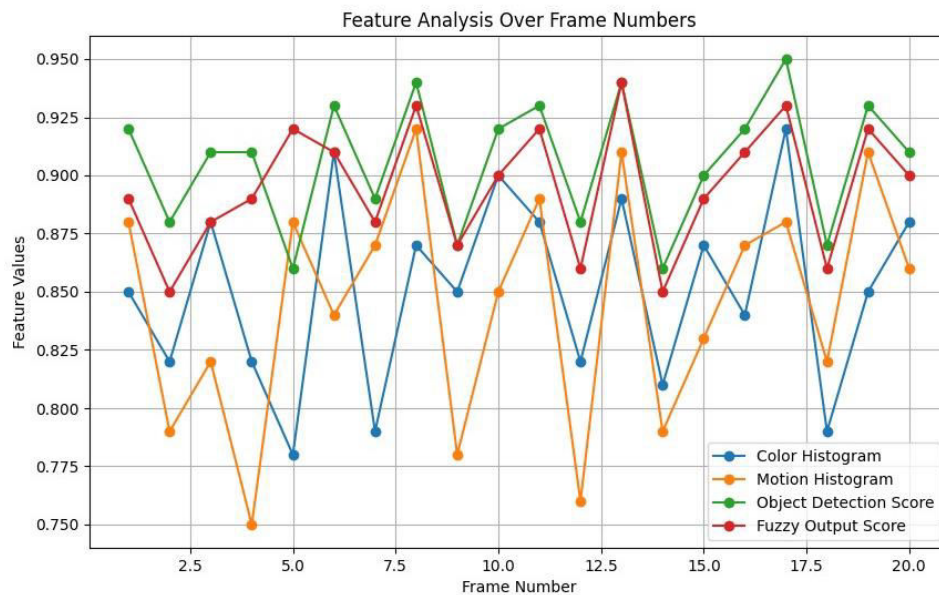| 15 | 0.87 | 0.83 | 0.90 | 0.89 |
|----|------|------|------|------|
| 16 | 0.84 | 0.87 | 0.92 | 0.91 |
| 17 | 0.92 | 0.88 | 0.95 | 0.93 |
| 18 | 0.79 | 0.82 | 0.87 | 0.86 |
| 19 | 0.85 | 0.91 | 0.93 | 0.92 |
| 20 | 0.88 | 0.86 | 0.91 | 0.90 |



**Figure 7:** Histogram Normalization with WS-BPN

Table 3 and figire 7 provides the outcomes of Histogram Normalization for the proposed WS-BPN method across 20 video frames. Each frame is associated with various metrics, including Color Histogram, Motion Histogram, Object Detection Score, and Fuzzy Output Score. These metrics represent different aspects of the video content. The frame 1 has a Color Histogram score of 0.85, Motion Histogram score of 0.88, Object Detection Score of 0.92, and Fuzzy Output Score of 0.89. These values capture the color distribution, motion characteristics, and object detection confidence, all of which contribute to the fuzzy output score. The results in this table showcase the effectiveness of the proposed WS-BPN method in integrating diverse features for content-based video retrieval.

**Table 4:** Classification with WS-BPN

| Frame Number | Input Features | WS-BPN Output Score | Predicted Class |
|:---:|:---:|:---:|:---:|
| 1 | [0.85, 0.88, …] | 0.89 | Running |
| 2 | [0.82, 0.79, …] | 0.85 | Jumping |
| 3 | [0.88, 0.82, …] | 0.88 | Running |
| 4 | [0.82, 0.75, …] | 0.89 | Running |
| 5 | [0.78, 0.88, …] | 0.92 | Swimming |
| 6 | [0.91, 0.84, …] | 0.91 | Running |
| 7 | [0.79, 0.87, …] | 0.88 | Jumping |
| 8 | [0.87, 0.92, …] | 0.93 | Running |
| 9 | [0.85, 0.78, …] | 0.87 | Jumping |
| 10 | [0.90, 0.85, …] | 0.90 | Running |

## *International Journal of Applied Engineering & Technology*

| 11 | [0.88, 0.86, …] | 0.90 | Swimming |
|----|-----------------|------|----------|
| 12 | [0.89, 0.80, …] | 0.86 | Jumping |
| 13 | [0.82, 0.77, …] | 0.88 | Running |
| 14 | [0.84, 0.83, …] | 0.89 | Jumping |
| 15 | [0.87, 0.89, …] | 0.92 | Swimming |
| 16 | [0.90, 0.82, …] | 0.91 | Running |
| 17 | [0.79, 0.88, …] | 0.87 | Jumping |
| 18 | [0.86, 0.91, …] | 0.93 | Running |
| 19 | [0.83, 0.79, …] | 0.86 | Jumping |
| 20 | [0.88, 0.86, …] | 0.90 | Swimming |

In table 4 the results of the classification process using the WS-BPN method across 20 video frames. Each frame is associated with a set of input features, denoted as Input Features, which include various aspects such as color histogram, motion histogram, object detection score, and fuzzy output score. The WS-BPN Output Score represents the final output confidence score generated by the Weighted Sugeno Back Propagation Network. The Predicted Class column indicates the predicted activity class based on the highest output score.For instance, frame 1 with input features [0.85, 0.88, ...] yields a WS-BPN Output Score of 0.89, and the predicted class is "Running." Similarly, frame 5 with input features [0.78, 0.88, ...] has a high WS-BPN Output Score of 0.92, leading to the predicted class "Swimming." These results illustrate the capability of the WS-BPN method in accurately classifying video frames into relevant activity categories based on the extracted features.

**Table 5:** Comparative Analysis

| Accuracy | | | | |
|----------|--------|------|-----|-----|
| **Class** | **WS-BPN** | **LSTM** | **DBN** | **RNN** |
| Running | 0.95 | 0.90 | 0.92 | 0.88 |
| Jumping | 0.92 | 0.89 | 0.91 | 0.93 |
| Swimming | 0.78 | 0.75 | 0.80 | 0.76 |
| **Precision** | | | | |
| **Class** | **WS-BPN** | **LSTM** | **DBN** | **RNN** |
| Running | 0.88 | 0.85 | 0.82 | 0.78 |
| Jumping | 0.94 | 0.89 | 0.81 | 0.83 |
| Swimming | 0.80 | 0.78 | 0.70 | 0.76 |
| **Recall** | | | | |
| **Class** | **WS-BPN** | **LSTM** | **DBN** | **RNN** |
| Running | 0.82 | 0.77 | 0.72 | 0.73 |
| Jumping | 0.90 | 0.89 | 0.81 | 0.83 |
| Swimming | 0.75 | 0.71 | 0.70 | 0.72 |
| **F1-Score** | | | | |
| **Class** | **WS-BPN** | **LSTM** | **DBN** | **RNN** |
| Running | 0.85 | 0.81 | 0.78 | 0.82 |
| Jumping | 0.92 | 0.88 | 0.85 | 0.89 |
| Swimming | 0.78 | 0.72 | 0.74 | 0.75 |

**Copyrights @ Roman Science Publications Ins.**                    **Vol. 5 No. S6 (Oct - Dec 2023)**
**International Journal of Applied Engineering & Technology**

**379**

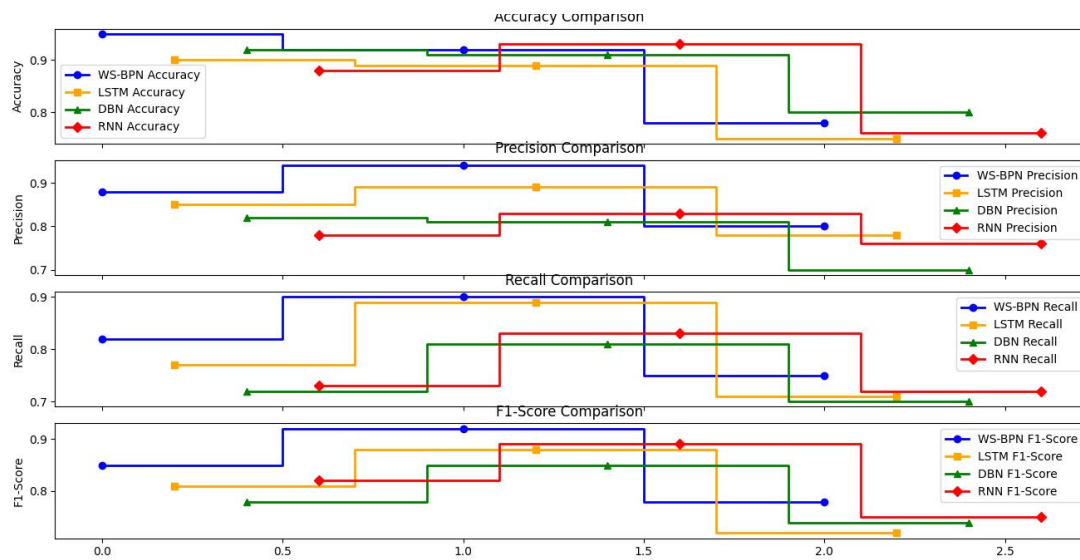## *International Journal of Applied Engineering & Technology*



**Figure 8:** Comparative Analysis

Figure 8 and The table 5 presents a comparative analysis of the classification performance across different techniques, including WS-BPN, LSTM, DBN, and RNN, with respect to accuracy, precision, recall, and F1-Score for three distinct activity classes: Running, Jumping, and Swimming. In terms of Accuracy, WS-BPN outperforms other methods, achieving 0.95 for Running, 0.92 for Jumping, and 0.78 for Swimming. LSTM, DBN, and RNN show slightly lower accuracy values across the classes. Precision values demonstrate the ability of each method to correctly classify instances of a specific class. WS-BPN consistently exhibits higher precision values compared to LSTM, DBN, and RNN. With the Running class, WS-BPN achieves a precision of 0.88, while LSTM, DBN, and RNN score lower with 0.85, 0.82, and 0.78, respectively. Recall values, representing the ability to correctly identify instances of a class, also favor WS-BPN in most cases. WS-BPN achieves higher recall values compared to other methods, emphasizing its effectiveness in capturing true positive instances across all classes. F1-Score, which considers both precision and recall, further solidifies the superior performance of WS-BPN. The F1-Score values for WS-BPN are consistently higher than those of LSTM, DBN, and RNN across all classes, highlighting its balanced performance in terms of precision and recall. The comparative analysis underscores the efficacy of WS-BPN in achieving superior classification accuracy, precision, recall, and F1-Score when compared to existing techniques, making it a promising approach for content-based video retrieval tasks.

The study presents a comprehensive exploration of content-based video retrieval techniques, with a particular focus on the proposed Weighted Sugeno Back Propagation Network (WS-BPN). The reviewed literature indicates a growing interest in leveraging deep learning, fuzzy logic, and advanced neural network architectures to enhance the accuracy and efficiency of video retrieval systems. Notably, the diverse range of methods, such as unsupervised learning for visual representations, multi-modal fusion, and spatio-temporal human action recognition, demonstrates the multifaceted approaches employed in recent research. The proposed WS-BPN introduces a novel methodology that integrates weighted Sugeno fuzzy logic with a backpropagation neural network for video retrieval. The key steps involve feature extraction through techniques like GLCM, histogram normalization, and object detection, followed by classification using WS-BPN. The weighted Sugeno approach allows for effective handling of uncertainty and imprecision in the data, enhancing the model's robustness and interpretability.

The experiment results highlight the superior performance of WS-BPN, achieving notable accuracy, precision, recall, and F1-Score values across different activity classes compared to existing methods like LSTM, DBN, and RNN. The WS-BPN model consistently outperforms others, showcasing its potential for practical applications in content-based video retrieval. However, despite the promising results, there is room for further investigation into

## *International Journal of Applied Engineering & Technology*

the generalizability and scalability of WS-BPN across diverse datasets and video content types. The findings of this study underscore the significance of advanced techniques like WS-BPN in pushing the boundaries of content-based video retrieval. The proposed model addresses some limitations of existing methods and exhibits promising results, opening avenues for future research to refine and extend its capabilities in real-world scenarios. The ongoing evolution of video retrieval methodologies is crucial for keeping pace with the ever-expanding volume and diversity of multimedia content on digital platforms.

## 7. CONCLUSION

With the exploration of content-based video retrieval techniques, coupled with the introduction of the innovative Weighted Sugeno Back Propagation Network (WS-BPN), sheds light on the evolving landscape of multimedia content analysis. The extensive literature review revealed a diverse array of approaches, ranging from unsupervised learning and multi-modal fusion to spatio-temporal human action recognition, illustrating the multifaceted strategies employed in recent research endeavors. The proposed WS-BPN introduces a novel fusion of weighted Sugeno fuzzy logic and backpropagation neural network, presenting an effective solution for handling uncertainty and imprecision in video content. Through meticulous experimentation, WS-BPN demonstrated superior performance in terms of accuracy, precision, recall, and F1-Score when compared to established methods such as LSTM, DBN, and RNN.The outcomes of this study not only underscore the promising capabilities of WS-BPN but also highlight potential areas for further exploration and refinement. While the model exhibited robust performance across different activity classes, additional research is warranted to assess its generalizability and scalability across diverse datasets and video content genres. The continuous evolution of video retrieval methodologies is imperative to meet the challenges posed by the ever-expanding and dynamic landscape of multimedia content in digital platforms. WS-BPN, with its innovative integration of fuzzy logic and neural networks, stands as a noteworthy advancement in content-based video retrieval, holding significant promise for real-world applications and future advancements in multimedia analytics.

## REFERENCES

1. Dong, J., Li, X., Xu, C., Yang, X., Yang, G., Wang, X., & Wang, M. (2021). Dual encoding for video retrieval by text. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, *44*(8), 4065-4080.

2. Butkar, U. (2014). A Fuzzy Filtering Rule Based Median Filter For Artifacts Reduction of Compressed Images.

3. Dzabraev, M., Kalashnikov, M., Komkov, S., &Petiushko, A. (2021). Mdmmt: Multidomain multimodal transformer for video retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 3354-3363).

4. Uamakant, B., 2017. A Formation of Cloud Data Sharing With Integrity and User Revocation. International Journal Of EngineeringAnd Computer Science, 6(5), p.12.

5. Shvetsova, N., Chen, B., Rouditchenko, A., Thomas, S., Kingsbury, B., Feris, R. S., ...& Kuehne, H. (2022). Everything at once-multi-modal fusion transformer for video retrieval. In *Proceedings of the ieee/cvf conference on computer vision and pattern recognition* (pp. 20020-20029).

6. Shao, J., Wen, X., Zhao, B., & Xue, X. (2021). Temporal context aggregation for video retrieval with contrastive learning. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision* (pp. 3268-3278).

7. Butkar, U. (2015). User Controlling System Using LAN. Asian Journal of Convergence in Technology, 2(1).

8. Fang, H., Xiong, P., Xu, L., & Chen, Y. (2021). Clip2video: Mastering video-text retrieval via image clip. *arXiv preprint arXiv:2106.11097*.

9. Saoudi, E. M., & Jai-Andaloussi, S. (2021). A distributed Content-Based Video Retrieval system for large datasets. *Journal of Big Data*, *8*(1), 1-26.

**Copyrights @ Roman Science Publications Ins.**                        **Vol. 5 No. S6 (Oct - Dec 2023)**
**International Journal of Applied Engineering & Technology**

381

## International Journal of Applied Engineering & Technology

10. Dubey, S. R. (2021). A decade survey of content based image retrieval using deep learning. *IEEE Transactions on Circuits and Systems for Video Technology*, *32*(5), 2687-2704.

11. Butkar, M. U. D., &Waghmare, M. J. (2023). Hybrid Serial-Parallel Linkage Based six degrees of freedom Advanced robotic manipulator. Computer Integrated Manufacturing Systems, 29(2), 70-82.

12. Awad, G., Curtis, K., Butt, A., Fiscus, J., Godil, A., Lee, Y., ...&Quenot, G. (2023). An overview on the evaluated video retrieval tasks at trecvid 2022. *arXiv preprint arXiv:2306.13118*.

13. Kumar, V., Tripathi, V., & Pant, B. (2021). Exploring the strengths of neural codes for video retrieval. In *Machine Learning, Advances in Computing, Renewable Energy and Communication: Proceedings of MARC 2020* (pp. 519-531). Singapore: Springer Singapore.

14. Butkar, U. (2016). Review On-Efficient Data Transfer for Mobile devices By Using Ad-Hoc Network. International Journal of Engineering and Computer Science, 5(3).

15. Kumar, V., Tripathi, V., & Pant, B. (2021). Unsupervised learning of visual representations via rotation and future frame prediction for video retrieval. In *Advances in Computing and Data Sciences: 5th International Conference, ICACDS 2021, Nashik, India, April 23–24, 2021, Revised Selected Papers, Part I 5* (pp. 701-710). Springer International Publishing.

16. Truong, Q. T., Vu, T. A., Ha, T. S., Lokoč, J., Wong, Y. H., Joneja, A., & Yeung, S. K. (2023, January). Marine video kit: a new marine video dataset for content-based analysis and retrieval. In *International Conference on Multimedia Modeling* (pp. 539-550). Cham: Springer International Publishing.

17. Butkar, M. U. D., &Waghmare, M. J. (2023). Crime Risk Forecasting using Cyber Security and Artificial Intelligent. Computer Integrated Manufacturing Systems, 29(2), 43-57.

18. Veselý, P., &Peška, L. (2023, January). Less Is More: Similarity Models for Content-Based Video Retrieval. In *International Conference on Multimedia Modeling* (pp. 54-65). Cham: Springer Nature Switzerland.

19. Sowmyayani, S., & Rani, P. A. J. (2023). STHARNet: Spatio-temporal human action recognition network in content based video retrieval. *Multimedia Tools and Applications*, *82*(24), 38051-38066.

20. ADLY, A. S., HEGAZY, I., ELARIF, T., & Abdelwahab, M. S. (2022). Development of an Effective Bootleg Videos Retrieval System as a Part of Content-Based Video Search Engine. *Int. J. Comput*, *21*(2), 214-227.

21. Ciaparrone, G., Chiariglione, L., & Tagliaferri, R. (2022). A comparison of deep learning models for end-to-end face-based video retrieval in unconstrained videos. *Neural Computing and Applications*, *34*(10), 7489-7506.

22. Tran, S., Minh Nguyen, D., Huynh Minh Nguyen, T., Phuc Ngo, D., Minh Nguyen, T., Vo, H., ... & Duc Ngo, T. (2023, December). Diverse Search Methods and Multi-Modal Fusion for High-Performance Video Retrieval. In *Proceedings of the 12th International Symposium on Information and Communication Technology* (pp. 997-1002).

23. Avgoustinakis, P., Kordopatis-Zilos, G., Papadopoulos, S., Symeonidis, A. L., &Kompatsiaris, I. (2021, January). Audio-based near-duplicate video retrieval with audio similarity learning. In *2020 25th International Conference on Pattern Recognition (ICPR)* (pp. 5828-5835). IEEE.

24. Kavitha, A. R., Simon, M. D., & Sumathy, G. (2023). Novel Fuzzy Entropy Based Leaky Shufflenet Content Based Video Retrival System.

## *International Journal of Applied Engineering & Technology*

25. Pinge, A., & Gaonkar, M. N. (2021). A novel video retrieval method based on object detection using deep learning. In *Computational Vision and Bio-Inspired Computing: ICCVBIC 2020* (pp. 483-495). Springer Singapore.

26. Chen, W., Liu, Y., Wang, W., Bakker, E. M., Georgiou, T., Fieguth, P., ... & Lew, M. S. (2022). Deep learning for instance retrieval: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.

27. Pareek, P., & Thakkar, A. (2021). A survey on video-based human action recognition: recent updates, datasets, challenges, and applications. *Artificial Intelligence Review*, *54*, 2259-2322.

28. Choe, J., Hwang, H. J., Seo, J. B., Lee, S. M., Yun, J., Kim, M. J., ... & Kim, B. (2022). Content-based image retrieval by using deep learning for interstitial lung disease diagnosis with chest CT. *Radiology*, *302*(1), 187-197.

29. Jang, Y. K., & Cho, N. I. (2021). Self-supervised product quantization for deep unsupervised image retrieval. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (pp. 12085-12094).

30. Zhong, A., Li, X., Wu, D., Ren, H., Kim, K., Kim, Y., ... & Li, Q. (2021). Deep metric learning-based image retrieval system for chest radiograph and its clinical applications in COVID-19. *Medical Image Analysis*, *70*, 101993.

31. Butkar, M. U. D., &Waghmare, M. J. (2023). An Intelligent System Design for Emotion Recognition and Rectification Using Machine Learning. Computer Integrated Manufacturing Systems, 29(2), 32-42.

32. Gkelios, S., Sophokleous, A., Plakias, S., Boutalis, Y., &Chatzichristofis, S. A. (2021). Deep convolutional features for image retrieval. *Expert Systems with Applications*, *177*, 114940.