Journal Pre-proof

Design of a Computational Model to Detect Hybrid Emotion through Facial Expressions in Videos Using CNN-LSTM

Sahaya Sugirtha Cindrella S and Jayashree R DOI: 10.53759/7669/jmc202505155 Reference: JMC202505155 Journal: Journal of Machine and Computing.

Received 12 May 2025 Revised from 28 June 2025 Accepted 09 July 2025



Please cite this article as: Sahaya Sugirtha Cindrella S and Jayashree R, "Design of a Computational Model to Detect Hybrid Emotion through Facial Expressions in Videos Using CNN-LSTM", Journal of Machine and Computing. (2025). Doi: https:// doi.org/10.53759/7669/jmc202505155.

This PDF file contains an article that has undergone certain improvements after acceptance. These enhancements include the addition of a cover page, metadata, and formatting changes aimed at enhancing readability. However, it is important to note that this version is not considered the final authoritative version of the article.

Prior to its official publication, this version will undergo further stages of refinement, such as copyediting, typesetting, and comprehensive review. These processes are implemented to ensure the article's final form is of the highest quality. The purpose of sharing this version is to offer early visibility of the article's content to readers.

Please be aware that throughout the production process, it is possible that errors or discrepancies may be identified, which could impact the content. Additionally, all legal disclaimers applicable to the journal remain in effect.

© 2025 Published by AnaPub Publications.



Design of a Computational Model to Detect Hybrid Emotion through Facial Expressions in Videos using CNN-LSTM

¹S Sahaya Sugirtha Cindrella, ²R Jayashree^{*}

^{1,2} Department of Computer Applications, Faculty of Science and Humanities, SRM Institute of Science and Technology, Kattankulathur, Chengalpattu, Tamilnadu-603202, India

sc1905@srmist.edu.in, jayashrr@srmist.edu.in

*Corresponding Author: R Jayashree (jayashrr@srmist.edu.in)

Abstract

In many applications of human-computer interaction, emotion prediction is To enhance emotion categorization, we present a hybrid deep learning m dy that blends his el h convolutional neural networks (CNN) with long short-term memory (LST) h networks. The preprocessing step refines the input data using Q-based score normaliz to ensure ideal feature scale and distribution. Emotional states are robustly classified when CN is employed to extract spatial data, and LSTM captures temporal relationships. Our models ability to identify intricate emotion patterns is demonstrated through training and eval atto on a benchmark emotion dataset. According to experimental results, our suggested CNN LaTM model performs exceptionally well on the test dataset, attaining 100% accuracy, precision, received of F1-score. These exceptional results highlight the power of combining Q Λ as LS I in handling emotion prediction's spatial and continuous aspects. Q-based score formalization further enhances the model's performance by ensuring a well-distributed feature space withmately improving classification consistency. This study underscores the potential of hybrid app learning architectures in improving emotion recognition applications. Our findings can be applied in diverse domains such as emotional computing, mental analytics, an computer interaction. hup

Keywords: Hybrid emotion brediction CNN, LSTM, Q based score, normalization, precision, recall, fl score and accuracy.

1. Introduction

interface between humans and computers more natural and efficient, kn. In ad ter intraction (HCI) aims to promote online learning, aid in human growth, and huma esthe cally pleasing designs and user experiences. Emotions have naturally emerged as produc a si vifica t component in creating HCI-based applications as they are essential to human ips. Numerous technology methods may be used to record and evaluate emotions, elation Voice, physiological signs, and facial expressions. Emotions expressed through signals luding be accurately identified and processed to provide a more intuitive and natural humanshou puter connection. Numerous machine learning methods have been created and continuously enhanced throughout the past 20 years of research on automatic emotion perception. Automated systems that can comprehend and interpret human emotions have the potential to revolutionize a society that is becoming more and more digital. Accurately interpreting human emotion might drastically alter these environments, given the rise of internet platforms for social engagement, education, healthcare, and remote work. For instance, automated facial expression detection might enhance telehealth services by providing physicians with real-time information on patients' emotional reactions [1].

Human emotions convey information through various channels, including speech patterns, body language, and facial expressions. Faces are important for communication and relationshift dynamics because they show how individuals react to different situations. Face expressions act as transmitters, giving others subtle clues. In daily interactions, these expressions—frequently unplanned and instinctive responses to stimuli—are essential to nonverbal communication that transmits social cues and reveals an individual's emotional condition. The late 20th-century vork by Ekman and Friesen, which categorized six universal facial expressions—placeus is surprise, disgust, sadness, and rage—solidified Darwin's idea [2].

The micro-expressions were highlighted by closely examining the video of a modally ill patient. The patient briefly displayed a concealed melancholy emotion, laking only two frames (1/12 second) while always seeming happy. Studies on micro expressions highlight their importance for spotting dishonesty, especially in high-stakes scenario take police interrogations when lying might have deadly consequences [3].

People who are struggling with mental health disorder in the find it challenging to correctly interpret and react to social signals, such as tone and facul expressions. Reading these cues and rehending those of others, and reacting correctly is essential for controlling our environment fostering relationships. Building adaptive set in ractions requires understanding how our social activities affect other people's views. To av id bad ffects, angry people, for example, may xp ssions. The subject of computer vision, known as decide not to show it verbally or by facial Facial Expression Recognition (FER), offers variety of methods for identifying human emotions from facial expressions. Understanding an individual's emotions helps improve behavioral science, therapeutic practice, and human-machine interaction, which is why researchers are interested in FER. More F. R systems can be created thanks to technological egorized methods. Applications for FER systems can be found in developments and picture rketing, healthcare systems, school counseling programs, lie the music business. detection, and targe ising [2]. d adve

age edges and represent objects according to their size, shape, and In or e detection techniques prior to machine learning. However, computer vision texture ised et , Fr has be in technique used by FER in recent years. Applications for FER are many and the identifying dishonesty and mental illnesses to identifying human emotions. FER rang fro. are based on facial photographs, which frequently need the identification of Regions of system OIs) because processing whole images would be computationally expensive. The erest tions of FER extend to key areas of computer vision, such as data-driven animation, ram active gaming, robotics, neuromarketing, and human-computer interface (HCI) [4].

In this instance, the use of FER for lie detection in police investigation is very complex and necessitates a deep understanding of body language and psychology. Reducing error rates and validating lie detection systems require the involvement of experts in these domains. Similarly, the judicial sector frequently uses complex FER approaches that produce hybrid CNN-RNN algorithms for lie detection. Lower complexity FER is useful for measuring client satisfaction in

industries like the fashion sector. This entails differentiating between good feelings that indicate interest or a desire to buy and negative emotions that indicate indifference to a product. This procedure is simpler than the difficult work of identifying falsehoods in police investigations [5].

FER's primary problem is handling large amounts of data variability that are impacted be personality, ethnicity, and expression intensity. Furthermore, factors like head posture and lighting greatly influence how well facial expressions can convey human emotions. Another difficulty is categorizing facial expressions into a small number of sorts since various ethnic groups have quite varied characteristics. Therefore, creating a strong identification model for FER is crucia. The primary technique for creating the FER model in recent years has been deep learning, cont Convolutional Neural Networks (CNN) [6].

Contributions of the study

- First, we introduce a novel hybrid CNN-LSTM framework that affectively combines spatial and temporal feature extraction, leading to superior chosic ation performance.
- Second, we propose Q-based score normalization as a prepressing technique that improves data consistency and overall model robustness.
- Third, our test findings show that the suggested method is resilient in emotion identification tasks, achieving 100% accuracy precision secall, and F1-score.
- Lastly, our study opens the door for future developments in real-time emotion analysis by offering insights into the useful uses of deep learning in emotional computing, mental health monitoring, and human-computer interaction.

Organization of the paper

The hybrid emotion prediction strappeviews the contributions of various authors in Section 2 (Related Works). Section 2 presents the proposed algorithm within the methodology. Section 4 focuses on dataset analysis, result prediction, and performance evaluation. Finally, Section 5 concludes the research why key indings and discussions.

2. Related work

Audio-usual entition recognition is essential for intelligent human-machine interactions, although training deep learning techniques often requires large amounts of data. To solve this problem, anultimodal Conditional Generative Adversarial Network (GAN) is proposed. Utilizing Hirschold-Kebelin-René (HGR) maximal correlations, the GAN simulates the strong reliance is tween ne visual and hearing senses utilizing category information as input. The created data is then usended to the data manifold to account for class imbalance [7]. The model performs better networks of accuracy, precision, and memory efficiency than pre-trained models like VGG16 and RESNET 50. This concept might be useful for cyber security, social media content regulation, and identity verification. The model further incorporates Channel-Wise Attention Mechanisms during feature extraction to improve RESNET50's overall performance [8]. By balancing contextual and granular information, this study suggests a hierarchical feature fusion network (DHF-Net) to detect emotions during conversations. DHF-Net analyzes joint attention, action/intention, and discourse

emotion to comprehend cross-influence. It employs a hybrid convolutional neural network (CNN) grouping technique with a bidirectional gate recurrent unit (Bi-GRU) [9].

Emotions are a multifaceted domain that impacts cognition, planning, and decision-making. Utilized in e-learning, marketing, and human-robot interaction, automated human emotio recognition (AHER) is a crucial area of study in computer science [10]. Combined text and visual semiotic systems with online content using deep learning networks and machine learning. The model consists of four parts: segmentation, text analysis, picture analysis, and the decision ock. The model outperforms individual text and picture modules with an accuracy of about 91% 11]. To lower model training error, the TK-GAN model is introduced, and the GAN methy Soft attention was established to enhance the learning capacity of stock elated lata characteristics. BERT improves the model's applicability to certain financial area branch. The MAE and MSE values in the experimental findings s 0.01949 and lov 0.00091[12].

th GANs that combines To estimate stock market prices, this article creates a hybrid mode. natural language processing, machine learning, sentiment analysis and statistics. The study's conclusions help investors make knowledgeable decisions on whether to buy, sell, or keep shares [13]. For characterizing and merging EEG data, EEGE seN is a hybrid unsupervised deep d mç convolutional recurrent generative adversarial network-l el. Deep EEG characteristics lly descroed. The model can recognize with temporal and spatial dynamics are automatic emotions and is dependable, robust, and simple to the. It is also effective at describing and combining dynamic EEG data [14]. The model uses relation-based transfer learning to analyze low-volume sentiment on SemEval public det. The extracted features are injected into an implicit ich has synergistic effects due to its ability to neural network (INN) in the target domain, process continuous and intermitter tata. The model achieves an accuracy of 88% and a 3% loss rate on SemEval data [15].

and experience, multimodal emotion recognition is essential. In affective computing, However, assessing e om several modalities is a challenge for traditional systems. GM ons functions apply dyr ghts to emotions in a hybrid multimodal emotion recognition (Hlmic w MMER) framework. Vith a average accuracy of 98.19%, the framework can accurately predict four distinct motional states [16]. Proposed hybrid deep learning model that uses eye movement ction with MesoNet4 and ResNet101 architectures to identify Deepfakes in realdata | con time. Th mode uses MesoNet4 to edit face images and ResNet101 to extract complex visual data. wratings of 0.9873 on FaceForensics++, 0.9689 on CelebV1, and 0.9790 on With rccur elebV. the hybrid model shows promise for video forensics and content integrity verification. If y, the model performs exceptionally well on several datasets, such as FaceForensics++, litior Selebv1, and CelebV2 [17]. Restaurant operations have been greatly influenced by online media, raising customer reviews. Although it struggles with imbalanced datasets, sentiment analysis (SA) aids in predicting review sentiments. This study suggests a hybrid strategy that blends PSO, oversampling, and SVM techniques. Compared to previous classification methods, this method increases accuracy, F-measure, G-mean, and AUC, benefiting the restaurant industry's success [18]. A pre-trained Spatial Transformer Network and bi-LSTM with an attention mechanism are used by the face emotion recognizers. Using the RAVDESS dataset, the program identified eight

emotions with an accuracy of 80.08%. Combining these modalities enhances system performance in several applications, such as road safety and healthcare systems [19]. A temporal-detection pipeline for microscopic-typo comparison of video frames. The program classifies authentic and fake visual input using 512 face landmarks and a Recurrent Neural Network (RNN) pipeline. The suggested algorithm and network showed competitive performance on any fake-generated imag or video and established a new standard for visual counterfeit detection [20]. To enhance emotion identification in Brain-Computer Interaction (BCI), this study suggests feature extraction and augmentation strategies. To produce more EEG characteristics, the technique uses Cond onal Wasserstein GAN and deep generative models. The method's efficacy is assessed using the D dataset. According to experimental results, EEG-based emotion identification mode are improved by adding enhanced data, with mean accuracy rising by 3.0% for aro nd 6.59 for valence [21]. Multimedia content has been mined using attention-based des neura (DNNs). Attention mechanisms have recently been included to high ght e otio ly significant information. This study looks at how attention processes affect performance and evaluates current advancements in SER. A benchmark database called IEMOCAP hused to compare system accuracies [22]. To generate a feature vector for emotional recognition the model combines outputs from N adaptive neuro-fuzzy inference system classifier using input vectors. On the y ratings of 73.49% and 95.97%, DEAP and Feeling Emotions datasets, the model received ccur respectively [23].

The study assessed a hybrid deep neur HDNN) for emotion identification using EEG data. The activation function impressed the odel's scuracy by accounting for the intricacy of the input and output data. The DEAP , which contained physiological and EEG inputs, વજ showed that the model's performance was encoded by the ELU function [24]. The model refines characteristics and predicts relevance using me-grained emotional similarity and feature canonical correlation analysis (CCA). After that, XGBoost is used to calculate similarity while considering emotional comp astan e and fine-grained affective semantic distance [25]. By th CNI bidirectional RNN with CNN (BRDC) model improves integrating fusion features the performance of b retation and classification. With 99.90% accuracy, 98.41% F1, inte 97.96% precision, % recall during training, it outperforms existing models [26]. It nd 99. TM hybrid model that combines an improved CIM optimization presents an IChOA NN-L using, long short-term memory networks for continuous data processing, and method 1 ature al feature extraction. With an impressive 97.8% accuracy rate, our model CNN ns cut ont techniques [27] outper lo

It of its a morough analysis of EEG emotion identification, emphasizing deep learning methods is chas ReN, CNN, and DBN. It looks at benchmark data sets, discusses innovative applications, assert the promise and issues, and offers suggestions for more study in this challenging field of [28]. To enhance the identification of emotional states from EEG data by combining deep features from wavelet CNNs with multiclass support vector machines (MSVM). In the process, EEG data is preprocessed and optimized using popular CNNs, and the best feature layer is chosen for MSVM classification. The approach produced enhanced accuracy rates of 77.47% and 87.45%, respectively, when tested on the DEAP and MAHNOB-HCI datasets [29]. Artificial intelligence, medical technology, and online education all depend on the ability to recognize emotions. While existing approaches frequently offer multi-category single-label predictions, deep learning

techniques increase the accuracy of EEG emotion identification. A novel method builds emotion categorization labels using DBSCAN and linguistic resources. With an average emotion classification accuracy of 92.98%, the DEAP dataset experiments have potential applications in social media, education, and mental health care [30].

3. Proposed methodology

The proposed hybrid emotion prediction model is based on convolutional neural network (CNNs) and long short-term memory (LSTM) networks. During the preprocessing age, normalization is applied to improve data consistency and model performance. The CNN ex spatial information from video frames to identify important patterns required for en tion recognition. These collected characteristics are then sent into an LSTM network ch loo sequential relationships to improve classification accuracy. The model's lictive is are assessed using key performance measures, with precision, recall, curac and 1 score. This hybrid method improves emotion identification from video data ng the advantages of utili ations, as seen in Figure both CNN and LSTM, making it appropriate for many real-world a 1.



Figure 1. Overall flow diagram of proposed hybrid prediction

Dataset

Real-time image acquisition is performed, and the captured videos are saved in a folder. Each video represents a distinct emotion. For example, one video may show a person smiling to express happiness, while another may show a person laughing to express joy. These videos serve as the primary source of emotional data, providing valuable information for analysis and integration into the emotion prediction process.

Preprocessing: Q-based score Normalization

One method to obtain the threshold as a function of Q is to break the issue down into many simple independent problems. Finding the proper boundaries θi for each neighborhood is simple. EQ may be separated into K-linked and unconnected neighborhoods, Ni, if EQ = N1 U··UNK Ni \cap Nj = 0/ \forall i = j. This formula may be readily modified if more than one quality indiction is offered for a single verification technique. Neighborhoods are defined by their borders. $\mu = (li, li+1)$, in the straightforward scenario when only one quality metric is given.

An acceptable method for creating these intervals must consider that the and it of prifection attempts within an interval determines how reliable any threshold estimated is. Q is, therefore, divided into intervals with about equal numbers of verification attervals. Assume that he sequence $T = \{(s1,q1,c1),...,(sNT,qNT,cNT)\}$ is arranged using the quality deterval $\leq qj \forall i < j$ without compromising generality. If we define a lower constraint for the measurement of quality measure $q0 = q1 - \varepsilon$, where ε represents an arbitrarily small positive constant, we may build the interval bounds as follows:

$$l_{i} = \frac{1}{2} \left(q \left[\frac{(i-1)N_{T}}{K} \right] + q \left[\frac{(i-1)N_{T}}{K} \right] \right) \forall i \in \{1, \dots, K-1\}$$
(1)

CNN feature extraction

The network typically takes features and trace them. The network is given particular places in the frame where the features are reduced and trained to decrease spatial complexity. CNN feature extractors employ a pre-trained model predicated on the transfer learning process. This enables the network to use both a a stom-trained model for video analysis and the pre-trained model's capabilities. Once the layer that will be utilized as a feature extraction point has been identified, the remaining layers of the pour CNN are removed using this procedure. The feature extractor network that is left over is then moved to the LSTM network. To simplify memory, video data is extracted into a start they are not preserved.

These frame-level image sets extract raw abstract facial features with high-dimensional visual information as nature vectors per frame.

$$u_n^{(v)} = (u_1, u_2, \dots, u_m)$$
 (2)

The Deep CNN feature extractor, the number of frames or frame sequences is n, and the put dimension is m.

Given v as the exclusive video identifier for a total accessible No. of videos in the training set, the set of all frame feature vectors $u_n^{(v)}$ used to predict the sequence $U^{(v)}$ is represented as features. The set below shows that there are T frames.

$$U^{(v)} = U_1^{(v)}, \dots U_T^{(v)}$$
(3)

LSTM (Long short-term memory)

The resulting video-level feature vector is sent into the LSTM cell inputs. The mean pooling layer of the LSTM temporarily stores frame-level visual features in its memory unit to obtain them without storing them. The CNN feature extractor obtains the feature vectors by using face embedding.

Following the extraction of patterns from the cell, the LSTM computes future patterns in queue frames after learning the pattern for each video frame v at time t. Prior connection and patterns are retained in the hidden states, making them accessible for additional process prospective memory cell state is formed from the internal memory cell state (c t), state (g_t) , the three gates (input gate (i_t) , output gate (o_t) , forget gate (f_t) , and input gate i_t . his serves as a link between the information at hand and hidden or previously store are pati controlled by mathematical calculations,

(5)

(7)

(8)

(9)

$$i_{t} = \sigma(W^{i}x_{t} \oplus w_{t} + U^{i}h_{t-1} + b_{i})$$

$$o_{t} = \sigma(W^{o}x_{t} \oplus w_{t} + U^{o}h_{t-1} + b_{o})$$

$$f_{t} = \sigma(W^{f}x_{t} \oplus w_{t} + U^{f}h_{t-1} + b_{f})$$

$$g_{t} = tanh(W^{g}x_{t} \oplus w_{t} + U^{g}h_{t-1} + b_{f})$$

$$c_{t} = f_{t} \otimes c_{t-1} + i_{t} \otimes g$$

$$h_{t} = o_{t} \otimes tanhc_{t}$$

Here, \oplus –vector concatenation operator

ł

 \otimes -element – wise multiplication between two vectors,

W – weights of inputs to

U – weights from hidde o hida

b – biases the val e of e ch **F**eature

the LSTM pipeline: There are two

1) Th ment is the comprehensive feature, which represents every frame in a certain movie ained using the CNN model's feature vector. The mean pooled frame feature serves as and is the hodel's popul for each frame.

After Witting each video into equal frame occurrences, the feature vectors for each movie are under the target vector/label T. Here, the target label is divided into two classifications ext and authentic—and one-hot encoding vectors (y1, y2). High-dimensional sparse vectors are compressed into lower-dimensional sparse vectors via the embedding layer, which assigns weights (w1,...,wn) to each feature vector. After being replicated and concatenated with weight vectors wt at each time step, the average frame feature x is fed into the LSTM model as (x1 Nw1,...,xT NwT). The output of the LSTM cells, which are responsible for directing the convergence of the memory cells to the ultimate goal state, are the intermediate hidden layers (h1,..., ht). For every video frame v, the conditional probability would be,

$$P(y_t,\ldots,y_1|x_t\otimes w_t,\ldots,x_1\otimes w_1)=\prod_{i\leq t\leq T}P(y_t|h_{t-1})$$

4. Result and discussion

CNNs- LSTMs are combined in the suggested hybrid emotion prediction model to enhance emotion identification from video data. During preprocessing, normalization is applied to ensure data consistency and optimize model performance. The CNN extracts spatial features from viace frames, capturing essential patterns, while the LSTM analyzes temporal dependencies to improve classification accuracy. The model performance is evaluated using key evaluation metres, including precision, recall, accuracy, and F1 score, ensuring a comprehensive analysis of its effectiveness. The implementation uses Python (version 3.11) in the Spyter evelopment environment. The system runs on a 64-bit processor with 8 GB of PtM, providing encient execution and processing of video data for accurate emotion prediction.



Figure 2. Loss and accuracy of emotion prediction

Figure 2 shows the input video frame accuracy and the training and validation losses as a function of benumber of epochs. These measures shed light on the model's training procedure and long term generalizability. As the number of epochs increases, the trends in training and validation accuracy indicate the model's convergence and performance, while the loss values reflect the reduction in errors during the training process. This analysis is essential to understanding the effectiveness of the CNN-LSTM framework in emotion prediction.



Figure 3 shows the confusion matrix of the hybrid continuous prediction model, which produces labels such as sadness-fear and amus ment-corprise. This matrix compares the actual and predicted labels, clearly visualizing the nonel's performance in distinguishing between these emotional categories. It highlights correct classifications and misclassifications, giving valuable insights into areas where the model may need improvement. The confusion matrix is a key tool to assess the model's accuracy and ibilities accurately predict different emotional states.



Figure 4. Classification report

The hybrid emotion prediction model's classification report offers comprehensive findings regarding accuracy, recall, and F1 score, particularly emphasizing the emotion pairings amusement-surprise and sadness-fear. These metrics, which are shown in Figure 4, offer thorough evaluation of the model's effectiveness for every emotion class. Recall gauges the capacity to find all pertinent examples, accuracy shows the percentage of accurate positive predictions, and the F1 score strikes a balance between the two. This classification report is calcial to evaluate the model's overall efficacy in differentiating the designated emotion categories.



Figure 5. erformance metrics of emotion prediction

The emotion prediction odel's r rmance analysis demonstrates exceptional outcomes with precision, recall, F1 accuracy, all hitting 100%. These results, shown in Figure 5, à illustrate the model classification capability across all emotion categories, indicating its perfec high level of effect emotion recognition. The flawless performance suggests that the ness model l carned to distinguish and predict emotions with maximum accuracy, ıllycces **n** its reliability for real-world applications. idence prov





Figure bapp results

The input videos, extracted to express the elections of sadness and fear, were used to train the model with the CNN-LSTM architecture. The process shown in Figure 6 illustrates how the model was trained on this emotion-labeled use data. The model can precisely recognize and categorize the emotions of fear and saddess because the CNN component collects spatial data from the video frames, and the LSTM network analyzes the continuous biases. The figure demonstrates the training flow and highlight the bain stages of feature extraction and emotion classification.

2/47 26s 600ms/step Accuracy: 1.0 .>ecision: 1.0 Recall: 1.0 F1-score: 1.0	
<pre>1/1 1s 1s/step Predicted Emotion: amusement_surprise</pre>	
Total samples: 6616 Emotion labels: ['amusement_surprise' 'sadness_fear']	

Figure 7. Prediction result.

The CNN-LSTM model for hybrid emotion prediction trained on a dataset of emotion models successfully identifies and classifies combined emotions such as amusement_surprise. The prediction results are illustrated in Figure 7.

5. Conclusion

This study presents a hybrid approach for emotion prediction by integrating normalization as a preprocessing step with a CNN-LSTM-based classification model. By increasing data consists normalization improves model performance and feature representation. CNN and LSTM of erate together to capture temporal and spatial connections, resulting in a more reliable and ac framework for classifying emotions. Our strategy successfully manages complicated d em onal patterns, as experimental data shows it outperforms existing methods. The model formar le is assessed using accuracy, precision, recall, and F1-score; all four metrics . This fect a outstanding result shows how reliable and effective the propose ch \ in accurately appr predicting emotions. The study concludes that the CNN-LSTM del c hbined with Q-based score normalization is a dependable and incredibly successful technique for predicting emotions. Applications in affective computing, human-computer interaction, and it tal health monitoring might greatly benefit from this paradigm. Future work may optimized yperparameters, incorporate enhance system efficiency and multimodal data, and explore real-time implementati adaptability.

Conflict of interest: The authors declare network licts if interest(s).

Data Availability Statement: The Datasets used and for analysed during the current study available from the corresponding author on a sonable request.

Funding: No funding.

Consent to Publish: All authors save permission to consent to publish.

References

- 1. M. Bakiaraj and B. Subamani, "Optimized hybrid deep learning pipelines for processing heterogeneous final expression datasets", Measurement: Sensors, Vol. 31, 2024.
- Wahab, M. N. A. Nazir, L., Zhen Ren, A. T., Mohd Noor, M. H., Akbar, M. F., & Mohamed, A. S. (1922) Environmet-Lite and Hybrid CNN-KNN implementation for facial expression recognition on Representation on Representa
- 3. Japos. Vincentius Manalu, Achmad Pratama Rifai, "Detection of human emotions through factal expressions using hybrid convolutional neural network-recurrent neural network algorium", Intelligent Systems with Applications, Vol. 21, pp. 1-18, 2024.
- Manra, R. K., Urolagin, S., Arul Jothi, J. A., & Gaur, P. (2022). Deep hybrid learning for facial expression binary classifications and predictions. Image and Vision Computing, 128(104573), 1–12.
- Morshed, G., Ujir, H., & Hipiny, I. (2021). Customer's spontaneous facial expression recognition. Indonesian Journal of Electrical Engineering and Computer Science, 22(3), 1436– 1445.

- Lozoya, S. M. G., De La Calleja, J., Pellegrin, L., Escalante, H. J., Medina, Ma. A., & Benitez-Ruiz, A. (2020). Recognition of facial expressions based on CNN features. Multimedia Tools and Applications, 79(19–20), 13987–14007.
- Fei Ma, Yang Li, Shiguang Ni, Shao-Lun Huang and Lin Zhang, "Data Augmentation for Audio–Visual Emotion Recognition with an Efficient Multimodal Conditional GAN", Applie Science, Vol. 12, pp. 1-24, 2022.
- Soha Safwat, Ayat Mahmoud, Ibrahim Eldesouky Fattoh, And Farid Al, "Hybrid Der Learning Model Based on GAN and RESNET for Detecting Fake Faces", Vol. 12, pp. 6391 – 86402, 2024.
- 9. Chenquan Gan, Yucheng Yang, Qingyi Zhu, Deepak Kumar Jain, Vitomir Struc, DHF-vet: A hierarchical feature interactive fusion network for dialogue emotion recommon", Epert Systems with Applications, Vol. 210, 2022.
- Eman M. G. Younis, Someya Mohsen, Essam H. Houssein, Osmai Ali Ludek Ibrahim, "Machine learning for human emotion recognition: a concrehencive review", Neural Computing and Applications, Vol. 36, pp. 8901–8947, 2024.
- 11. Akshi Kumar, Kathiravan Srinivasan, Wen-Huang Cheng, Albert X. Zunaya, "Hybrid context enriched deep learning model for fine-grained sentiment analysis in textual and visual semiotic modality social data", Information Processing & Management, Vol. 57, No. 1, 2020.
- 12. Rui Zhang and Vladimir Y. Mariano, "Integration of Leadins, Factors with GAN Algorithm in Stock Price Prediction Method Research", IEC Access, vol. 12, pp 77368-77378, 2024.
- 13. Fares Abdulhafidh Dael, Omer Çagrı Yavuz a'd Uger Yavuz, "Stock Market Prediction Using Generative Adversarial Networks (14Ns); Aybrid Intelligent Model", Computer Systems science and Engineering, vol.47, no.1, 2005.
- 14. Zhen Liang, Rushuang Zhou, Li Zhang, Lhuing Li, Gan Huang, Zhiguo Zhang and Shin Ishi, "EEGFuseNet: Hybrid Unsure Dised Deep Feature Characterization and Fusion for High-Dimensional EEG With an Approximation to Emotion Recognition", IEEE Transactions on Neural Systems and Rehabilitation Engineering, Vol. 29, pp. 1913 – 1925, 2021.
- 15. Kia Jahanbin And Moha, mad Ali Zare Chahooki, "A Hybrid Deep Implicit Neural Model for Sentiment Analy, s Via Transfer Learning", IEEE Access, Vol. 12, pp. 131468-131486, 2024.
- 16. Muhammad Asi, Razzac Jamil Hussain, Jaehun Bang, Cam-Hao Hua, Fahad Ahmed Satti, Ubaid Uri Luman, Heiz Syed Muhammad Bilal, Seong Tae Kim and Sungyoung Lee, "A Hybrid a ultimetral Emotion Recognition Framework for UX Evaluation Using Generalized Mustare Functions", Sensors, Vol. 23, pp. 1-25, 2023.
- 17. Muha, mad Laved, Zhaohui Zhang, Fida Hussain Dahri and Asif Ali Laghari, "Real-Time De nfake video Detection Using Eye Movement Analysis with a Hybrid Deep Learning Approach", Electronics, Vol. 13, Issues, 15, 2024.
- 18. Die Obiedat, Raneem Qaddoura, Ala M. Al-Zoubi1, Laila Al-Qaisi, Osama Harfoushi, Mo'ath Alrefai, and Hossam Faris, "Sentiment Analysis of Customers' Reviews Using a Hybrid Evolutionary SVM-Based Approach in an Imbalanced Data Distribution", IEEE Access, Vol. 10, pp. 22260-22273, 2022.
- 19. Cristina Luna-Jiménez, David Griol, Zoraida Calleja, Ricardo Kleinlein, Juan M. Montero and Fernando Fernández-Martínez, "Multimodal Emotion Recognition on RAVDESS Dataset Using Transfer Learning", Sensors, Vol. 21, pp. 1-29, 2021.

- 20. Mohammad Farukh Hashmi, B. Kiran Kumar Ashish, Avinash G. Keskar, Neeraj Dhanraj Bokde, Jin Hee Yoon, And Zong Woo Geem, "An Exploratory Analysis on Visual Counterfeits Using Conv-LSTM Hybrid Architecture", Vol. 8, pp. 101293-101308, 2020.
- 21. Mahsa Pourhosein Kalashami, Mir Mohsen Pedram and Hossein Sadr, "EEG Feature Extraction and Data Augmentation in Emotion Recognition", Computational Intelligence an Neuroscience, Vol. 2022, Issue.1, pp. 1-16, 2022.
- 22. Eva Lieskovská, Maroš Jakubec, Roman Jarina and Michal Chmulík, "A Review on Specific Emotion Recognition Using Deep Learning and Attention Mechanism", Electronics, V. 10, pp. 1-16, 2021.
- 23. F. Kebire Bardak, M. Nuri Seyman and Feyzullah Temurtas, "Adaptive neuroriuzzy ased hybrid classification model for emotion recognition from EEG signals", Nura Computing and Applications, Vol. 36, pp. 9189–9202, 2024.
- 24. Jehosheba Margaret Matthew, Masoodhu Banu Noordheen and Lohan and Mustafa, "Enhancement of Hybrid Deep Neural Network Using Activation Function for EEG Based Emotion Recognition", Traitement du Signal, Vol. 41, No. 4, 2024
- 25. Zhibin Su, Yiming Feng, Jinyu Liu, Jing Peng, Wei Jiang and Jingy Liu, "An Audiovisual Correlation Matching Method Based on Fine-Grained Emotion and Feature Fusion", Sensors, Vol. 24, Issue 17, 2024.
- 26. Shofiqul Islam, Nahidul Islam, Noramiza Hashira, Noraunu Rashid, Bifta Sama Bari and Fahmid Al Farid, "New Hybrid Deep Learning Appenden Using BiGRU-BiLSTM and Multilayered Dilated CNN to Detect Arrhythesia", SEE Access, Vol. 10, pp. 58081 58082, 2022.
- 27. R. Geethanjali and A. Valarmathi, "A numer hybrid deep learning IChOA-CNN-LSTM model for modality-enriched and multilingual enotion recognition in social media", Scientific Reports, Vol. 14, 2024.
- 28. Xiaohu Wang, Yongmei R h2, 200 Luo, Wei He, Jun Hong and Yinzhen Huang, "Deep learning-based EEG emption recognition: Current trends and future perspectives", Frontiers in Psychology, Vol. 14, 1 1126994, 2023.
- 29. Sara Bagherzadra, Keuran Linghooli, Ahmad Shalbaf and Arash Maghsoudi, "A Hybrid EEGbased Emotion decogni on Approach Using Wavelet Convolutional Neural Networks and Support Value Machine, Basic clinical neuro science, Vol. 14, Issue 1, 2023.
- 30. Liumeterhang, Bowen Xia, Yichuan Wang, Wei Zhang and Yu Han, "A Fine-Grained Approach or EEG-Based Emotion Recognition Using Clustering and Hybrid Deep Neural Networks", Electronics, Vol. 12, No. 23, 2023.

