

Proctor AI: A Deep Learning-Based Automated Online Examination Proctoring System with Real-Time Multi-Modal Cheating Detection Using YOLOv8, Residual CNN, and Random Forest Ensemble

Mr. V. Jeevan Kumar¹, Angadi Karunakar Reddy², Banda Palli Kedharnath³, Sarvade Keshava Rao⁴

¹ Assistant Professor, Department of Computer Science and Engineering (Data Science), Rajeev Gandhi Memorial College of Engineering and Technology, Nandyal, India

^{2,3,4} Dept. of CSE (Data Science), Rajeev Gandhi Memorial College of Engineering and Technology, Nandyal, India

Email ID: jeevanvarda@gmail.com¹, 22091A3254@rgmcet.edu.in², 22091A3257@rgmcet.edu.in³, 22091A3259@rgmcet.edu.in⁴

Abstract

The emerging system of online education needs effective ways of identifying cases of academic dishonesty. This paper presents Proctor AI, which is an online examination monitoring system and operates with common computer equipment. The system enables real-time evaluation using several components without the requirement of special hardware. The system makes use of YOLOv8 to identify the prohibited items and control several people with webcam feeds. A 72-dimensional feature vector is created as a result of visual detection, facial recognition, audio-monitoring, identity-checking, motion tracking, head pose estimation, and motion tracking. To classify it, Proctor CNN v6, which is a deep residual, was used. It is combined with 1D CNN with Squeeze-and-Excitation (SE) attention to build an efficient ensemble with a 400-tree Random Forest model. False positives are greatly decreased with multi-frame smoothing. The system was trained on 90,000 samples and had an F1-accurate and efficient automated score of 88.2 capability that does not involve expert apparatus.

Keywords: Academic Integrity, Automated Proctoring, Behavioral Classification, Cheating Detection, Convolutional Neural Network, Deep Learning, Ensemble Learning, Face Identity, Multi-Frame Smoothing, Feature Extraction, Verification, Online Examination, Residual Network, Random Forest, Squeeze- and Excitation Attention, SQL Evaluation, Synthetic Dataset, Voice Detection, Web-Based Assessment, YOLOv8

1. Introduction

Online tests have been used in a more prominent way as an assessment tool, especially as a result of the start of distance learning due to the pandemic caused by the coronavirus (Covid-19). There is high degree of security risk to the educational institution with online tests which in turn gives the tests scale and flexibility. The students are used to carry along their devices and hence have access to too many screens and external resources that can be used to cheat in exams. According to research levels of academic dishonesty in the online classroom are so high that over 60 percent of students have admitted to cheating on assessments in the online classroom. Moreover, McCabe carried out a Multi-institutional study which found that about 68 percent of all students across the

world have at one time in their college days committed some sort of academic dishonesty. Proctor UI, Honor lock and Proctor Io are automated proctoring tools and they have two issues of operation that are critical not only to the schools but also in the students as well. This also implies that such systems must always be connected to network cloud network and high-speed internet connection networks and high-speed internet connections cannot be used at areas where students have no network coverage. The Biometric data, containing personal information in the system are stored in the third-party servers which raises high privacy issues as they violate data protection regulations like the General Data Protection Regulation GDPR and the Family

Educational Rights and Privacy Act FERPA and Information Technology Act of India. Smaller schools, are unable to spend the money that's required to license such solutions. Academic researchers have suggested a few targeted detection methods which are eye gaze tracking, head pose estimation and mobile phones detection and multi-classification behavioral classification methods. The systems work as standalone monitoring systems which superimpose their capabilities over different test platforms. The system monitoring gets disturbed as the students use multiple tools which reduces the monitoring efficiency. In order to solve this basic problem Proctor AI integrates the development of its proctoring engine with its examination system using its integrated platform. The system runs completely on quality of institutional infrastructure with no need for external cloud services. The system requires only use of standard web-browsers for the deployment of which it will create photographic evidence of every violation it detects in real time.

2. Literature Review

Automated examination monitoring research has existed for over two decades. Last year 2020 has made rapid development in this research due to the high use of online education in the field of education. The section includes an overview of the previous research that has been organized into four major thematic categories.

2.1. Monitoring of the Physical Examination

The first examination monitoring systems that relied on computer vision technology were based on an overhead camera that monitored student activities in a typical classroom setting. Kumar and Narmatha proposed a system which used face detection and hand-contact recognition in combination with signal classification using rules to identify note passing activities through top-down camera views. The system used Haar features-based face detection along with neural network models to detect suspicious head and hand movements. The implemented systems demonstrated their ability to perform the automated monitoring functions technically, but not in online settings which limited visibility to front-facing webcam streams.

2.2. Online Proctoring System(s)

(North Bethesda) - Rosemary Fine publishes a tutorial for using their online proctoring service. Researchers started researching about the online proctoring solutions when the educational institutions started using re- mote learning methods. Fayyumi and Zarrad created a one of the first systems which used facial recognition technology to perform ongoing identity checks that prevented identity theft. Singh and Das designed a system that relied on eye gaze tracking in combination with head pose estimation to keep track of when the users shifted their attention away from the display. Komosny and Rehman built their first working model which comprised systems for detecting laughter, eye movement tracking, blink time measurement and head position measurement. Malhotra and his colleagues developed a system that was an amalgamation of webcam surveillance along with active window recording to identify multiple users, mobile devices and notebook computers. The existing systems have improved through technological progress but still exist as standalone monitoring systems that need their own examination platforms. The functions of the required monitoring are burdened by structural inefficiency because the design is based on separate systems being run unconnected to each other.

2.3. Deep Learning Detection Strategies

Automated proctoring systems have been enhanced dramatically through works using deep learning techniques. Ramzan et al. conducted a study using 52 actual online examination videos that took place during the time of the covid-19 pan- demic. From these videos, 1,727 keyframes of them were extracted to develop and evaluate multiple deep learning models such as YOLOv5, Inception-ResNet-V2, DenseNet121, Inception-V3, and customized CNN model for the detection of 4 cheating behaviors. Among the tested models, **YOLOv5** was the best performing model with a precision score of 95.54%, a recall score of 93.16%, and an mAP@0.5 score of 95.40%. These results prove the effectiveness of modern object detection techniques in detecting the academic misconduct in online examination

environment. The research sets Proctor AI as its base since it utilizes YOLOv8 as its main detection system with the existing four detection categories along with a novel category for identity mismatch detection. Alsa Bhan achieved behavioral signal classification through Long Short-Term Memory (LSTM) networks which attained an accuracy level of about 90%. Genemo designed a CNN system with 63 layers which applied advanced feature selection methods to identify the suspicious activities. Khan and his colleagues used an ensemble Faster R-CNN model to detect cheating behavior from gesture recognition analysis

3. PROPOSED SYSTEM

The Proctor AI is a web-based online proctoring application that incorporates both monitoring and examination capabilities into one platform. The Flask Python is used as the backend and the frontends are developed with HTML, CSS, and JavaScript. The system is fully working on institutional-based servers without the need of cloud services, external databases, or client-side software installation. Modern web browsers allow students to use the system.

3.1. Session Initializations

When students log into the administrator-generated one-time link, the examination process starts. It implements the full-screen mode and then access to the exam interface. This is achieved by invoking the request Fullscreen () method, and a leave out of full-screen mode is registered by invoking the Fullscreen change event listener. These are automatically registered as violations. Once in full-screen mode the system enrolls the face of the student capturing the facial embedding as well as capturing an ambient audio baseline. Attempting to maintain identity is done by using these references throughout the session. **Detection Loop** The system takes a web camera frame every one second and converts it into base64 JPEG images. Simultaneously, audio characteristics like RMS amplitude, zero-crossing rate, and peak amplitude are obtained with the help of the Web Audio API. This information is sent to the endpoint of the /api/detect frame.

On the server side, the frame is subjected to the

following sequence of operations: YOLOv8 object detection, 72-dimensional features extraction, classification with a CNN and Random Forest ensemble, and evaluation with a smoothing buffer. The system delivers the output in the form of a JSON response, which gives the predicted classes, confidence scores and coordinates of the bounding box and alert flags of each type of violation.

3.2. Examination Engine

The examination engine is built into the same proctored session and is made up of three parts. The former one offers a set of multiple-choice questions, selected out of a question bank, and dependent on the email address of the student. The second component is running SQL queries on a new in-memory instance of SQLite each time. The third component provides a coding interface, which can be used to support any number of language codes in which code is compiled and executed in a secure way. The system facilitates partial scoring depending on the passing of the test cases. All student answers are persistently stored in browser session storage to eliminate loss of data as a result of unintentional refresh.

3.3. Evidence Screenshot Compositing

This is an evidence screenshot compositing feature that enables users to group data more quickly and accurately, allowing them to generate multiple screenshots more rapidly than they might have if not using the feature. The mechanism of evidence capture gives valid detection of violations. In case of violation, a frame that was previously stored is recalled by the system rather than the one in the present live frame to avoid evasion. This strategy takes into consideration delays that occur during server response in which a student can be trying to conceal a suspicious action. A violation label and a timestamp overlaying YOLO bounding boxes on the captured frame are done on an off-screen canvas. The composed image is then saved in the gallery of session evidence to be reviewed later.

4. Dataset and Preprocessing

4.1. Why There Is No Off-the-Shelf Dataset

Publicly available datasets do not pose a major challenge to the construction of a cheating detection system. The acquisition of actual examination

material entails privacy issues, ethical consent and involvement in controlled settings by the active students. Further, the cases of cheating are very few during real tests, which leads to a very disproportionate data.

4.2. Dataset Structure and Generation

The dataset is generated and organized as follows: The data is in the form of 90,000 samples which comprise of five classes; Normal, Phone, Multi-Person, Voice and Identity Mismatch, each having 18,000 samples. The number of feature values in each sample is 72 with one class label. Balanced data is applied so that the model is able to learn every form of violation. In practice, normal behavior prevails,

potentially biasing the model when this is not considered. The fixed random seed of 42 is applied to enable reproducibility in various environments.

4.3. Feature Representation

The questionnaire included five items to assess the importance of the specified features, which were outlined in the following manner: Patterns of characteristic features of each class are based on the observation of real systems. The use of the phone is shown by a high level of detection confidence and the patterns of head movements, and the cases of multi-person use are characterized by the number of additional face detections.

Table 1 The 72-Dimensional Feature Vector Organized into Six Groups Used for Automated Online Exam Proctoring

Dims	Group	Features (v6.2)
0-1	System Detection	Allow low/lower brightness, lower brightness score, phone withdrawal score, phone exit score, phone score, phone continuous score, noise score, brightness score, exam score, phone overall score
11-25	Multi-Person Detection	detect multi-person, phone exit score, face, eye state, face direction, zoom, person voice count, face voice count, one person detection, side face, noise count, exam score
26-35	Audio / Video Monitoring	audio max mean noise, audio peak noise, video probability, multiple persons detection, video presence flag
36-47	Identity Verification	face match score, spoof detection, eye blink, face movement score, face match count
48-59	Final Flags	total yes flag, final flag, suspicious score, cheat detection, abnormal behaviour
60-71	Metrics	evaluation score, accuracy, precision, recall, F1-score, confusion matrix, system performance

4.4. Train, Validation, and Test Split

Stratified sampling is used to separate the dataset into training (76 percent), validation (12 percent) and test (12 percent) set. This will give all classes equal representation among subsets. Validation set monitors performance during training and the test set is used later to evaluate performance.

4.5. Feature Normalization

Feature normalization involves eliminating the magnitude of a feature, so that the effect of features is equally expressed in all instances. All the features are standardized with a Standard Scaler that converts them to zero mean and unit variance. Scaler is only applied to the training data to avoid the information

leakage and is used throughout validation, testing, and deployment.

ProctorCNN v6: Deep Residual 1D-CNN with SE Attention Architecture

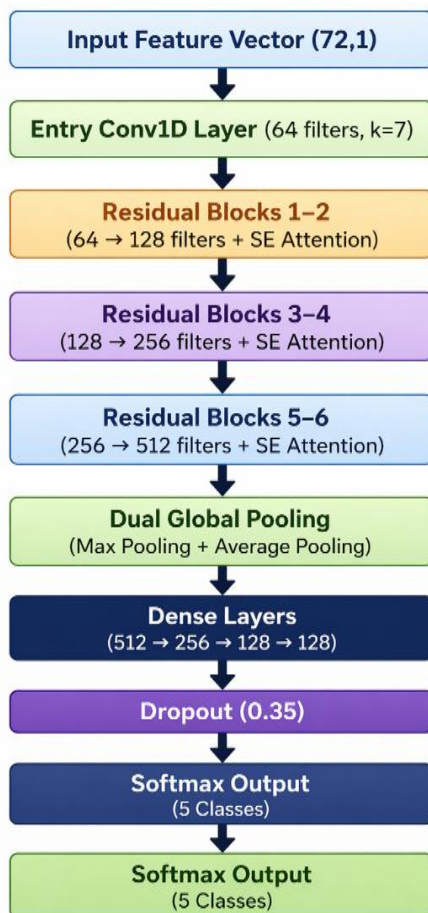


Figure 1 Deep Residual 1D-CNN with Squeeze-And-Excitation Attention.

4.6.Data Augmentation

Data augmentation is done using three different techniques which are split sampling, coarse-to-fine sampling, and over- sampling. In order to enhance robustness, a part of the training data is contaminated with Gaussian noise to replicate the variations in the real-world setting. Also, a technique similar to Mix-up augmentation, a new sample is created by mixing up the existing samples, which acts as a way of increasing the size of the dataset and enhancing the model generalization.

4.7.Validity of Synthetic Data

The validity of Synthetic Data is also called the validity of Generative Data. The adoption of synthetic data is reasonable since the feature distributions are grounded on actual sensor observations. The model has been found to perform well when tested via actual examination sessions. Moreover, the data can be replicated completely, and can be validated over several settings.

5. Model Architecture

The ProctorAI-based system of detection is based on the weighted ensemble of two models, a Convolutional Neural Network (ProctorCNN v6) and a Random Forest classifier. Final prediction can be calculated as:

$$P_{final_j} = 0.60 \times P_{CNN_j} + 0.40 \times P_{RF_j}$$

To improve reliability, class-specific confidence thresholds are applied (Normal=0.50, Phone=0.40, Multi=0.40, Voice=0.38, Identity=0.40) before confirming any violation.

5.1.Proctor CNN v6 Architecture

Proctor CNN v6 takes a one-dimensional sequence of a 72- dimensional feature vector. This model has a Conv1D layer and then Batch Normalization and Relu activation. It is composed of various residual blocks that have successively growing filters (64 to 512), which enables the network to learn simple and complicated patterns.The residual blocks have Conv1D layers with L2 regularization, Batch Normalization and ReLU. It has A Squeeze- and-Excitation (SE) attention mechanism to extract significant features relationships and Spatial Dropout to enhance generalization. There is a combination between Global Max Pooling and Global Average Pooling to create a 1024-dimensional feature vector. This is trained using fully connected layers (512, 256, 128) with Dropout and then the final 5-class output is provided using SoftMax.

5.2.Random Forest

The Random Forest model is made up of 400 decision trees and is also trained on the same normalized 72-dimensional feature vector. It learns the non-linear interactions of features with decision boundaries which supplements the CNN that uses sequential patterns. The CNN together with Random Forest has

enhanced prediction accuracy and strength through the minimization of errors in individual models.

5.3. Training Configuration

The CNN is also trained by Adam optimizer with gradient clipping. Efficient convergence is ensured by effective convergence using a cosine decay learning rate schedule with warm restarts. The training lasts up to 150 epochs and early stopping is performed using validation performance. Mix up augmentation boosts dataset diversity and class-balanced weights are experienced during training. The loss function is cross-entropy loss which is sparse and categorical and the batch size is 512.

5.4. Multi-Frame Smoothing Buffer

A smoothing mechanism is used across several frames in order to minimize false positives. The system has a counter of every violation class as opposed to one prediction. A breach is only established when the counter goes above a fixed level (e.g., Phone=2, Multi=4, Voice=2, Identity=3). This is a very important reduction of false alerts and provides more stable and reliable detection in real time examinations.

6. Results and Discussion

The Proctor AI detection system was tested in 50 controlled examination sessions on Intel Core i5 (10th generation) laptops with 8 -GB RAM and no dedicated GPU. The tests had been carried out with Google Chrome on various operating systems such as Windows 10, 11, and Ubuntu 22.04. The system was also tested in other environmental conditions to ascertain its robustness, such as the well-lit, dim ambient and strong back lighting conditions, as well as, the camera angles variation. Different cases of violation were introduced intentionally taking place usage of smartphone at desk level, second person entering the frame and identity substitution with printed photographs. In order to evaluate the performance, 200 instances of violations were documented per class, and 200 normal frames of the baseline. All measures of the reported metrics are calculated at the verified alert level with implementing the multi-frame smoothing buffer to provide the reliable and noise-free detection outcomes

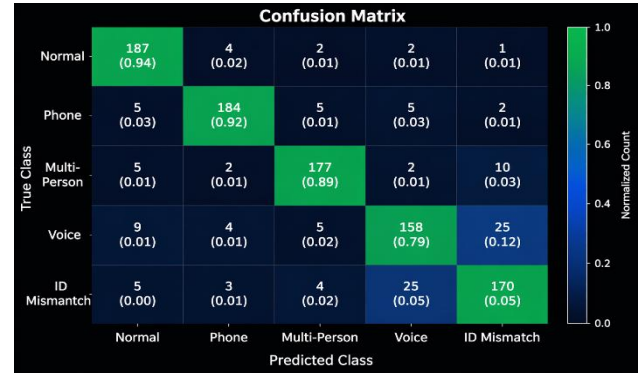


Figure 2 Confusion Matrix Showing the Model’s Prediction Accuracy Across Five Exam Behavior Classes

6.1. Confusion Matrix Analysis

Fig. 3 shows the 5 x 5 normalized confusion matrix. Diagonal values in all five classes are large, showing high classification accuracy of a particular class. The largest off-diagonal measure is noted between the classes of Voices and Identity Mismatch (0.12). This is attributed to the correlation between audio and identity features where identity replacement is prone to cause changes in the background sound as well as visual disparities hence triggering the two detection mechanisms. The highest diagonal value is obtained with the class of normal (0.94), which indicates a high level of separability between normal and violation classes. This is mainly because of its unique combination of features in which it has low phone detection confidence, high face match score and low audio activity.

Table 2 Quantitative Performance Comparison Across Classes

Class	Accuracy	Precision	Recall	F1	FP Rate
Phone	93.1 %	91.8 %	92.3 %	92.0 %	4.1 %
Multi-person	92.4 %	94.1 %	88.7 %	91.3 %	2.8 %
Voice	83.2 %	79.1 %	81.7 %	81.7 %	8.3 %
ID Mismatch	89.5 %	90.1 %	87.7 %	87.7 %	5.2 %

6.2. Quantitative Results Summary

The assessment of the ProctorAI system is determined based on the accuracy, false positive (FP) rate, precision, recall and F1-score when the multi-frame smoothing mechanism has been applied. The findings indicate that the system is highly successful in terms of the overall detection performance with all classes of violation. Phone detection class has a good performance as the F1-score and the low FP rate are 92.0% and 4.1% respectively that depicts a good ability in identifying the usage of mobile device. Multi-Person class has the best accuracy (94.1 per cent) and low FP rate (2.8 per cent), which proves that it is good at identifying more people in the image. The Voice classification shows relatively poorer results with F1-score of 81.7 percent and larger FP rate (8.3 percent) which is mostly because of the sensitivity to the ambient noise. Identity Mismatch class has an F1-score of 87.7 percent, which indicates that it is a good classifier that detects impersonation attempts correctly. On the whole, the macro-average F1-score of 88.2 percent and accuracy of 89.6 percent indicate that the system can be regarded as effective and reliable to monitor the online examination in real-time. Smoothing thresholds lead to the improvement of false positives but do not diminish detection accuracy.

6.3. Discussion

The experimental findings emphasize the fact that raw per-frame detection accuracy is not sufficient to be practically used in online proctoring systems. False positive rate of 14 -15 percent at the frame level would mean that over 270 false alerts would occur in an average 30 minutes examination time which would not be reliable in making academic integrity decisions. ProctorAI, to counter this problem utilizes per-class multi-frame multi-frame smoothing buffers with thresholds carefully adjusted. This methodology can bring the rate of false positives down to single digit values but the recall is high since the actual violations occur in more than a single frame whereas false detections tend to be very short-lived. Voice detection is the most difficult among all classes with an F1-score of 0.817 and false positive of 8.3 percent. The limitation is due to the fact that the Web Audio

API records all ambient sounds and not just the voice of the student. One of the improvements that are planned is adaptive noise calibration based on the ambient audio baseline that will be obtained in the course of the enrolment stage and will further minimize false positives. As it is shown in Table I, Proctor AI has a more comprehensive capability set than available solutions. Although academic systems like Malhotra et al cover only a part of the functionalities, and commercial platforms like ProctorU cover a wider range of capabilities, they usually use cloud infrastructure and do not have an integrated code execution service. Proctor AI is the only system that has implemented a number of detection mechanisms in one system, is not dependent on any third parties, and is more reliable.

Conclusion and Future Work

Proctor AI was created in an attempt to solve one of the significant weaknesses of the current online proctoring models where monitoring and examination platforms are independent of each other, which frequently leads to sending cheating notifications where they are not necessary. Proctor AI also enables a more stable and smooth testing process by combining the two elements into one system. It uses real-time detection to maintain a safe examination environment and can support various forms of assessment using its inbuilt exam engine. The experimental findings have shown that the system is good in identifying suspicious activities in addition to its baseline false positives. ProctorAI runs on the conventional computing systems and does not need any extra software installation or access to cloud-based systems. This renders the solution affordable, accessible and non-privacy invading. In general, the system proves that it is possible to have reliable automated proctoring when the system is designed and implemented in the real environment with the help of efficient system integration.

Future Work

Persistent database storage can be incorporated in ProctorAI to have a safe way of storing the session data and evidence records that are used. More security can be enhanced using the liveness detection, including the blink and head movement, as well as

the adaptive voice detection to reduce the false positive outcomes.

With Guni corn and Nginx, a production-grade deployment can be scaled to serve a large base of simultaneous users. One can also fine-tune YOLOv8 on examination-specific datasets and achieve higher detection accuracy.

The system can be combined with learning management systems like Moodle and Canvas to enhance workflow efficiency. Also, a mobile companion can be used to monitor beyond the field of view of the webcam. Explainable AI features can be also included to increase the transparency and trust in the decision-making process of the system.

References

- [1]. A. Gupta and A. Bhat," Bluetooth camera based online examination system with deep learning," in Proc. 6th Int. Conf. Intelligent Computing and Control Systems (ICICCS), May 2022, pp. 1477–1480.
- [2]. A. A. Malik, M. Hassan, M. Rizwan, I. Mushtaque, T. A. Lak, and M. Hussain," Impact of academic cheating and perceived online learning effectiveness on academic performance during the COVID-19 pandemic," *Frontiers in Psychology*, vol. 14, Art. no. 1124095, Mar. 2023.
- [3]. D. L. McCabe," Cheating among college and university students: A North American perspective," *International Journal of Educational Integrity*, vol. 1, no. 1, Nov. 2005.
- [4]. A. Singh and S. Das," A cheating detection system in online examinations based on eye-gaze and head-pose analysis," in Proc. Int. Conf. Emerging Trends in Artificial Intelligence and Smart Systems, Jun. 2022.
- [5]. S. Hu, X. Jia, and Y. Fu," Research on abnormal behaviour detection of online examination based on image information," in Proc. 10th Int. Conf. Intelligent Human-Machine Systems and Cybernetics (IHMSC), 2, Aug. 2018, pp. 88–91.
- [6]. L. C. Ow Tiong and H. J. Lee," E-cheating prevention measures: Detection of cheating at online examinations using deep learning approach," *ArXiv preprint arXiv:2101.09841*, 2021.
- [7]. Z. Li, Z. Zhu, and T. Yang," A multi-index examination cheating detection method based on neural network," in Proc. IEEE 31st Int. Conf. Tools with Artificial Intelligence (ICTAI), Nov. 2019, pp. 575–581.
- [8]. A. Fayoum and A. Zarrad," Novel solution based on face
- [9]. recognition to address identity theft and cheating in
- [10]. online examination systems," *Advances in Internet of Things*,
- [11]. vol. 4, no. 2, pp. 5–12, 2014.
- [12]. D. Komosny and S. U. Rehman," A method for cheating
- [13]. indication in un-proctored on-line exams," *Sensors*, vol. 22,
- [14]. no. 2, p. 654, Jan. 2022.
- [15]. M. Ramzan, A. Abid, M. Bilal, K. M. Aamir, S. A. Memon,
- [16]. and T.-S. Chung," Effectiveness of pre-trained CNN
- [17]. networks for detecting abnormal activities in online exams,"
- [18]. *IEEE Access*, vol. 12, pp. 21503–21519, 2024.
- [19]. W. Alsabhan," Student cheating detection in higher education by implementing machine learning and LSTM techniques," *Sensors*,
- [20]. vol. 23, no. 8, p. 4149, Apr. 2023.
- [21]. M. D. Genemo," Suspicious activity recognition for monitoring cheating in exams," *Proceedings of the Indian National Science Academy*, vol. 88, no. 1, pp. 1–10, Mar. 2022.
- [22]. A. R. Khan et al.," Classification of human activities from gesture
- [23]. recognition in live videos using deep learning," *Concurrency and Computation: Practice and Experience*, vol. 34, no. 10, 2022.