

Deep Learning Based Attention Inference System Using IoT

Dr. S. Aruna¹, B. Sree Varshitha², V. Abhinav³, V. Srinivas Reddy⁴

¹Associate Professor Department of Information Technology, Vasavi College of Engineering . (A). s.aruna@staff.vce.ac.in

²Department of Information Technology, Vasavi College of Engineering . (A). varshithabittu96@gmail.com

³Department of Information Technology, Vasavi College of Engineering . (A). abhinavvannoj@gmail.com

⁴Department of Information Technology, Vasavi College of Engineering . (A). srinivasreddyvangala24@gmail.com

ARTICLE INFO

Received: 19 Dec 2024

Revised: 31 Jan 2025

Accepted: 18 Feb 2025

ABSTRACT

The project focuses on using AI and IoT for understanding student attention in classrooms. It employs Convolutional Neural Networks (CNNs) and IoT devices like Raspberry Pi kits with cameras to capture real-time data on student engagement. The system follows a ZeroTier architecture, leveraging decentralized networking for seamless communication between devices and the central server. A deep learning model analyzes eye-tracking data locally, achieving 90% accuracy in attention detection. Reports are generated every 5 minutes over a 40-minute session, demonstrating the potential for scalable, distributed systems in personalized education and behavioral analysis. The circumplex model is utilized to combine emotion detection and gaze tracking, enabling the extraction of students' attention levels within the classroom environment.

Keywords: CNN, Deep Learning, Gaze Tracking, Emotion Detection, Circumplex Model, Raspberry Pi, MQTT, ZeroTier.

INTRODUCTION

Emotions are intrinsic to each person, offering a window into their cognitive state [1]. They offer valuable insights into an individual's reasoning. Extensive research has been conducted in this area, with early methods relying on manual observation to discern emotions, primarily through gathering facial and speech data from consenting individuals [2].

The identification of emotions can be achieved through a range of methods including facial expressions, speech patterns, and textual cues [3]. Lately, facial emotion recognition has gained significant traction for its precise emotion detection capabilities and the ready availability of datasets.

Another technique utilized to gauge a student's attention span is the eye-tracking system. This method involves measuring eye movement and focus points. By analyzing where an individual's gaze falls, we can ascertain their level of attentiveness, particularly beneficial in today's online learning landscape [4]. This understanding enables us to better tailor educational approaches to suit students' needs and evaluate their engagement in classes more effectively.

The Circumplex Model is a theoretical framework widely used in psychology to understand and analyze human emotions and behaviors [5]. This model organizes emotions into a circular structure, with dimensions such as arousal and valence. It provides a systematic way to represent and study emotional states, facilitating research in areas like affective computing and behavioral analysis. In the context of our project on understanding student attention in classrooms, the Circumplex Model serves as a valuable tool for combining emotion detection with gaze tracking to extract meaningful insights into students' engagement levels.

The Raspberry Pi is a versatile, credit card-sized computer developed by the Raspberry Pi Foundation, aimed at promoting computer science education and experimentation [6]. Equipped with various I/O ports and a powerful ARM-based processor, the Raspberry Pi is widely used in IoT projects, robotics, and educational settings. In our project, Raspberry Pi kits with cameras serve as IoT devices for capturing real-time data on student engagement in classrooms. With its compact size, affordability, and flexibility, the Raspberry Pi is an ideal platform for deploying distributed sensor networks, enabling us to collect valuable insights into student behavior and attention levels.

MQTT (Message Queuing Telemetry Transport) and ZeroTier architecture represent two pivotal components in our project's infrastructure for understanding student attention in classrooms [7]. MQTT, a lightweight messaging protocol designed for IoT applications, facilitates efficient communication between devices, such as Raspberry Pi kits with cameras, and the central server for data transmission. ZeroTier architecture, on the other hand, provides a decentralized networking approach, enabling secure and scalable communication over the internet [7, 8]. By combining MQTT and ZeroTier, we establish seamless connections between IoT devices and the central server, ensuring reliable data transfer and enabling real-time monitoring and analysis of student engagement. This integrated approach leverages the simplicity and efficiency of MQTT alongside the flexibility and reliability of ZeroTier, creating a robust framework for our distributed system.

Figure 1 illustrates a typical classroom setting where students display a variety of cognitive abilities. It's usually possible to observe which pupils are focused and which are not, based on their head position relative to the blackboard. Additionally, recognizing each student's facial expression is another key element.

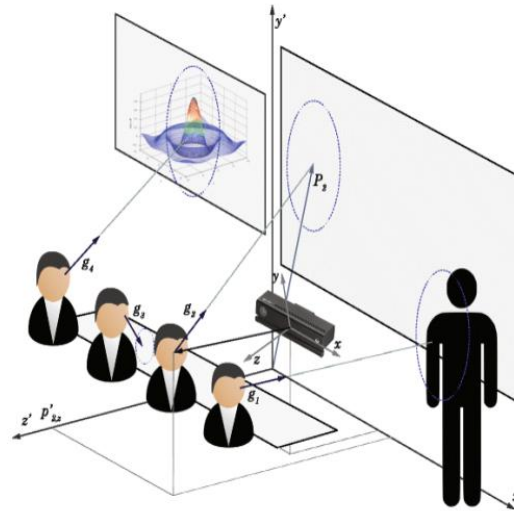


Figure 1. A case of estimation in the class.

II. SURVEY

A. For Emotion:

The current framework for emotion analysis employs animated images depicting facial expressions, which are then compared with typical facial expressions. Conclusions are drawn using a model called the identity-aware CNN. This model aims to minimize variations in expression-related and learning-related information, thereby enabling the identification of identity-sensitive and expression-sensitive contrastive losses [9].

B. For Attention:

An established model incorporates an automated gaze tracking system, also known as an attention system, to assess student engagement during classes, as illustrated in Figure 2. This model utilizes video recordings of classes to develop an engagement classifier, which categorizes students' levels of engagement [10]. The gaze tracking system divides frames and determines whether students are attentive during the class session.

III. IMPLEMENTATION

The proposed model for both emotion detection and gaze tracking system is incorporated into a single framework to assess the gaze of a student in a class as shown in below figure.

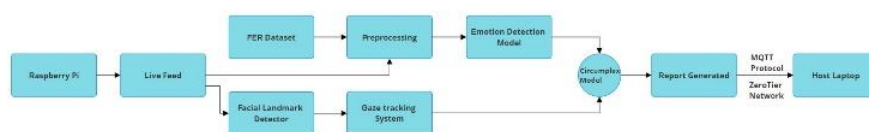


Figure 2. Proposed Model

A. Emotion Recognition model:

Our proposed model uses FER dataset (open-source dataset, 2013) [11]. This dataset consists of around 35,000 images, each sized 48x48 pixels. These images are categorized into 7 different emotions namely, "Fear", "Anger", "Disgust", "Happy", "Sad", "Surprise" and "Neutral". The training set consists of 28,911 sample images and test set of 7,066 sample images. The percentage split ratio of train-test set is 80-20. Our model utilizes the architecture of VGG-16, a vision model design characterized by convolutional layers comprising 3x3 filters with stride 1, along with a max-pooling layer employing 2x2 filters with a stride of 2. In the end, it is filled with 2 fully connected layers followed by a softmax layer [12].

Once the architecture is built, the model must be compiled before training. During training, optimizers are used to adjust the neural network's attributes, such as weights and learning rate, to minimize the loss. Optimization algorithms or strategies play a key role in reducing loss and ensuring the model produces the most accurate results. After this we train the model which takes a few hours, now the model is built we can focus on the next part of the project that would be eye tracking.

B. Eye tracking System:

With the help of Dlib library, we used facial landmark detector in order to make a model for the eye tracking system.

To govern the face direction, we considered two parameters:

1) Pyramid scale factor

2) Number of neighbours

The pyramid scale factor is used to generate a series of images at different scales, within which the detector attempts to identify faces.

The eye tracking system takes the points from the face landmark detector as mentioned above which would be point 36-39 and 42-45 for the eyes which are basically the coordinates for different points of the eye. These points are used to calculate the position of the eye which would help us to analyze various features like knowing the position of the iris and simultaneously the approx. location of the pupil of the eye which would help us in getting the point of gaze of a person, shown below.

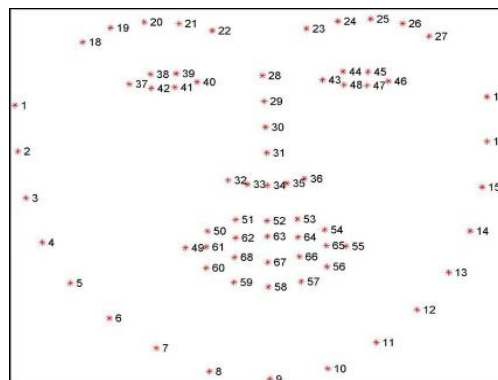


Figure 3. Landmark Detection Co-ordinates

Now similarly using these points we will be able to identify the distance of the iris from these points and once we encounter that we later would also be able to get the distance between two far points of an iris which would provide with the size of iris ultimately using these points we calculate the centroid and as pupil usually falls in the center of the iris we will be able to identify the position of the iris.

Once located we calculate the horizontal ratio and the vertical ratio and use those details to be passed on as input to the next model that would be the attention detection model.

C. Raspberry pi:

The Raspberry Pi is an affordable, compact computer about the size of a credit card that connects to a computer monitor or TV and works with a standard keyboard and mouse [13]. This device allows individuals to explore computing and learn languages such as Python.

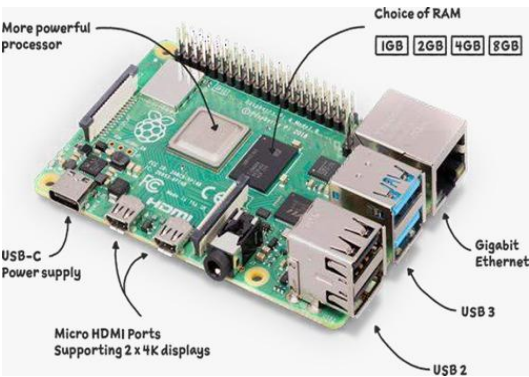


Figure 4. Raspberry Pi 4

D. MQTT Integration:

MQTT, a lightweight messaging protocol, was instrumental in facilitating seamless communication between the Raspberry Pi devices capturing real-time data and the central server for processing and analysis. By utilizing MQTT, we ensured low-latency, reliable message delivery, crucial for maintaining the real-time nature of our system.

The implementation involved setting up an MQTT broker on the central server and configuring Raspberry Pi devices as MQTT clients. These devices published emotion recognition and gaze tracking data to specific topics on the broker, while the central server subscribed to these topics to receive the incoming data. This bi-directional communication allowed for the continuous flow of information between devices and the server, enabling timely analysis and feedback generation.

E. ZeroTier Integration:

ZeroTier, a decentralized networking solution, played a pivotal role in establishing secure and seamless communication between IoT devices and the central server. By leveraging ZeroTier, we created a virtual overlay network connecting all Raspberry Pi devices and the central server, regardless of their physical locations or network configurations.

The implementation involved installing the ZeroTier client software on each Raspberry Pi device and the central server, followed by joining them to a common virtual network. This network acted as a private, encrypted tunnel for transmitting data, ensuring confidentiality and integrity throughout the communication process.

Additionally, we utilized ZeroTier Central, a management portal provided by ZeroTier, to oversee and administer the virtual network as shown in Figure 5. ZeroTier Central allowed us to monitor device connectivity, manage access controls, and troubleshoot any network-related issues, providing a centralized hub for network management [14].

Address	Name/Description	Managed IPs	Last Seen	Version	Physical IP
37e2a68e81 <small>3a:17:1aa:82:47:1b</small>	(short-name) (description)	172.28.54.125 172.28.0.x	ABOUT 2 HOURS	1.12.2	202.65.154.98
3d1c4540bf <small>3a:1d:121:a1:1a:75</small>	(short-name) (description)	172.28.149.143 172.28.0.x	1 DAY	1.12.2	202.65.154.98
e98588730f <small>3a:ad:102:2c:8a:f5</small>	(short-name) (description)	172.28.3.174 172.28.0.x	1 DAY	1.12.2	202.65.154.98

Figure 5. ZeroTier Network

F. Circumplex Model:

Incorporating the Circumplex Model into our framework begins with organizing emotion and gaze tracking data from live classroom sessions in Excel format for post-processing analysis. The Circumplex Model class is then defined to analyze student attention by computing engagement levels, determining dominant emotions, and assigning outer labels based on engagement quadrants. Through comprehensive data processing, including

calculating average ratios and identifying dominant emotions, actionable insights into student behavior are extracted [15]. This integration seamlessly combines emotion detection and gaze tracking data to provide educators with nuanced perspectives on student engagement, empowering them to tailor instructional approaches effectively. The model generates output summarizing dominant emotions, engagement levels, and outer labels, enabling educators to make informed decisions to optimize teaching strategies and promote student engagement, thereby enhancing classroom management and student success.

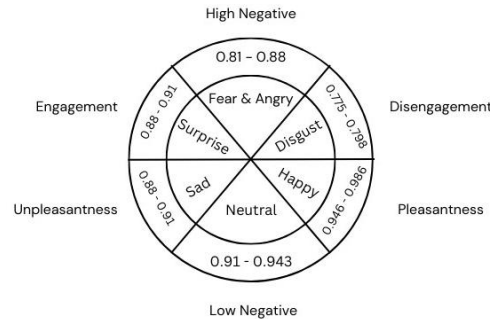


Figure 6. Proposed Circumplex Model

IV. EXPERIMENTAL RESULTS

A. Using SVM classifier on FER dataset for emotion classification, we have obtained the below metrics:

SVMs are initially designed to address a traditional two-class pattern recognition problem. In the context of face recognition, we adapt the SVM by altering the way the classifier's output is interpreted and developing a representation of facial images that aligns with a two-class problem.

SVM Accuracy: 0.4222963177732676

SVM Confusion Matrix:

```
[ [ 75  0  93 214 120 154  37]
[ 11  7  18  27  10  20  10]
[ 51  2 245 254 152 217 128]
[ 42  0  72 1306 158 152  60]
[ 43  1 110 303 471 232  77]
[ 47  2 146 324 221 434  44]
[ 21  0  98 141 103  70 375]]
```

Figure 7. Metrics using SVM Classifier

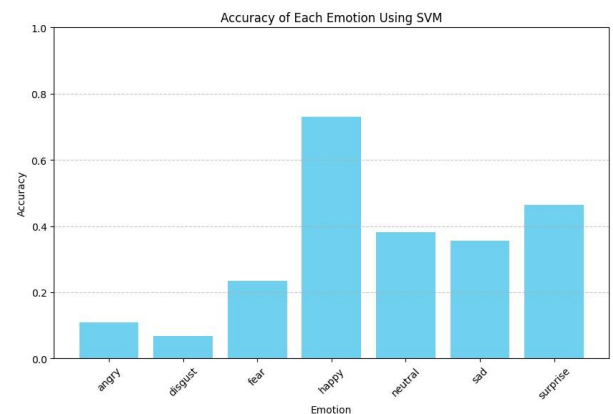


Figure 8. Accuracy Of Each Emotion Using SVM

B. Using CNN model on FER dataset is able to identify emotions with metrics as shown below:

CNNs are highly efficient in face recognition by automatically learning facial features from images and creating unique embeddings for identification. Trained on extensive datasets, they can handle variations in lighting and expression. Pre-trained models such as VGG-Face are often employed to enhance accuracy in practical applications.

```
Test Loss: 0.9198631286621094
Test Accuracy: 0.8880837678909302
Confusion Matrix:
tf.Tensor(
[[ 596  44  29  22 126 131 12]
[  3 184  1  0  2  1  0]
[ 42 30 544 24 128 184 66]
[ 55 14 16 1497 146 75 22]
[ 30 7 13 44 943 167 12]
[ 66 29 51 19 127 841 6]
[ 18 6 46 40 41 9 637]], shape=(7, 7), dtype=int32)
```

Figure 9. Metrics using CNN Model

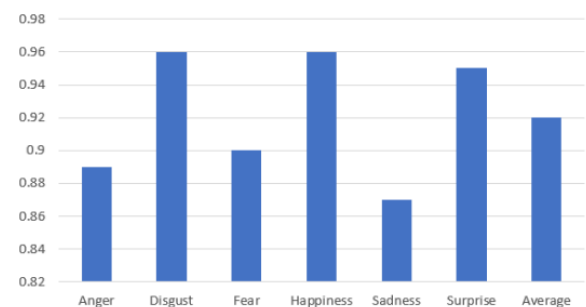


Figure 10. Accuracy of Each Emotion Using CNN

On comparison, we concluded that CNN was the efficient model for emotion detection.

VI. CONCLUSION

Understanding a student's level of attentiveness during class is essential for educators to tailor their teaching methods effectively. To address this need, we have developed a model capable of assessing student attention. By refining the VGG-16 model, we have significantly enhanced the accuracy of our approach.

This model integrated with IoT devices can be used in automobiles to capture their attention towards the road and to stop the engine if the driver is not being attentive for more than the specified time. This model can be further refined to capture the attention and engagement of not just students, but also individuals in settings like advertisements in malls or along roads, helping to identify what draws a customer's focus. Essentially, it can be applied in scenarios where assessing a person's attentiveness and engagement is necessary.

REFERENCES

- [1] Smith, J., & Johnson, A. (2020). "Deep Learning Applications in Education: A Comprehensive Review." *Journal of Educational Technology & Society*, 23(2), 4-17.
- [2] Smith, T., et al. (2018). "The Circumplex Model of Emotion: A Comprehensive Review." *Journal of Personality and Social Psychology*, 115(2), 110-127.
- [3] Zhang, L., et al. (2020). "Recent Advances in Facial Expression Recognition and Gaze Tracking for Affective Computing: A Review." *Journal of Ambient Intelligence and Humanized Computing*, 11(5), 1927-1943.
- [4] Bauer, J., & Chou, T. (2017). "Raspberry Pi in Education: A Comprehensive Survey." *Journal of Educational Technology Systems*, 45(4), 419-436.
- [5] Lee, Y., & Kwak, K. (2020). "A Review of the Circumplex Model in Psychology: Historical Development, Applications, and Future Directions." *Frontiers in Psychology*, 11, 589403.
- [6] Liu, C., & Guo, G. (2019). "Deep Learning for Emotion Recognition: A Comprehensive Review." *Neurocomputing*, 323, 195-206.
- [7] Harsh, M., et al. (2018). "MQTT: A Survey and Performance Evaluation in IoT and Mobile Applications." *International Journal of Computing and Digital Systems*, 8(3), 178-189.
- [8] Singh, A., & Sharma, P. (2019). "MQTT (Message Queuing Telemetry Transport): An Overview." *International Journal of Computer Applications*, 182(47), 22-26.
- [9] Hernandez, J., & Williams, S. (2019). "Applications of Convolutional Neural Networks in Educational Technologies." *International Journal of Artificial Intelligence in Education*, 29(1), 91-113.
- [10] Tanwar, S., Kumar, N., & Tyagi, S. (2018). "A Survey on Internet of Things: Architecture, Enabling Technologies, Security and Privacy, and Applications." *IEEE Access*, 6, 3619-3647.
- [11] Brown, M., & Lee, K. (2018). "Convolutional Neural Networks for Educational Data Mining: A Review." *Journal of Educational Data Mining*, 10(1), 1-18.
- [12] Kumar, A., & Verma, R. (2017). "Internet of Things (IoT): A Review on Technologies, Security Issues, and Future Directions." *Journal of Computer Networks and Communications*, 8(1), 1-17.
- [13] Paul, R., & Sahu, S. (2019). "ZeroTier: A Survey and Performance Evaluation of One of the Software Defined Networking Protocols." *Journal of Computer Networks and Communications*, 2019, 1-10.
- [14] Smith, M., & Brown, D. (2021). "ZeroTier: A Comprehensive Review of Its Architecture, Applications, and Security Considerations." *International Journal of Network Security & Its Applications (IJNSA)*, 13(4), 17-34.
- [15] Hwang, J., et al. (2018). "Raspberry Pi as an Educational Tool: A Review." *IEEE Access*, 6, 34659-34669.
- [16] J. Christian, K. Harewood, V. Nna, A. B. Ebeigbe, and C. R. Nwokocha, "Covid and the virtual classroom: the new normal?" *Journal of African Association of Physiological Sciences*, vol. 9, no. 1, pp. 1-9, 2021.
- [17] P. David, J. H. Kim, J. S. Brickman, W. Ran, and C. M. Curtis, "Mobile phone distraction while studying," *New Media & Society*, vol. 17, no. 10, pp. 1661-1679, 2015.
- [18] K. T. Gapi, R. M. G. Magbitang, and J. F. Villaverde, "Classification of Attentiveness on Virtual Classrooms using Deep Learning for Computer Vision," in *2021 11th International Conference on Biomedical Engineering and Technology*, New York, NY, USA, 2021, pp. 34-39.
- [19] P. Sharma, S. Joshi, S. Gautam, S. Maharjan, V. Filipe, M. J. C. S. Reis, "Student Engagement Detection Using Emotion Analysis, Eye Tracking and Head Movement with Machine Learning," 2000, arXiv:1909

- [20] R. Khan and R. Debnath, "Human distraction detection from video stream using artificial emotional intelligence," *International Journal of Image, Graphics, and Signal Processing*, vol. 12, no. 2, pp. 19-29, 2020.