

Statistical Multiframe Methodology with Agnostic Thresholding for Attendance Marking System

KUAN HENG LEE, SANJAY V. ADDICAM, ILYA KRYLOV, SERGEI NOSOV,
MEE SIM LAI, ZHAN QIANG LEE AND CHUNG SHIEN CHAI

*Intel Microelectronics Sdn. Bhd.
IOTG Retail Banking Hospitality Education
Bayan Lepas, Penang, 11900 Malaysia*

*E-mail: {kuan.heng.lee; addicam.v.sanjay; ilya.krylov; sergei.nosov; mee.sim.lai;
zhan.qiang.lee; chung.shien.chai}@intel.com*

Attendance marking is a burdensome and time-consuming task for every school teaching staff to perform manually in the classroom. It becomes very attractive if this attendance marking process can be automated through a facial recognition system. Although facial recognition works well under constrained environment, identifying each student in a dynamic classroom environment remains a challenge especially the students are in uncooperative manner. Conventional frame-based accuracy metric cannot reflect the true outcome of the attendance as it varies drastically over frames, due to the large variations of scales, poses and occlusions in the classroom environment. In this paper, a statistical methodology based on multiframe was proposed to improve the attendance marking accuracy after a convergence time. This methodology was combined with the mean thresholding scheme to achieve the same accuracy as full inference rate (*i.e.* 30 FPS) with a lower inference rate (*i.e.* 3 FPS). This drives away the need to invest an expensive hardware to maintain the same accuracy with a higher inference rate.

Keywords: artificial intelligence (AI), face detection (FD), face recognition (FR), interpupillary distance (IPD), sliding window filtering, frame per second (FPS), false positive (FP), false negative (FN), thresholding

1. INTRODUCTION

Facial Recognition (FR) has been widely deployed in many applications, including airport security, device access verification, library access and *etc.* It has been known to work well under constrained environment [1] (*e.g.* sufficient lighting, cooperative faces and close distance with test subject) with good accuracy (*i.e.* > 95%). FR can be a value added for attendance marking application to remove the need of manual inspection and attendance sheet.

Many of the FR based attendance marking solutions were proposed, such as the monitoring the students on the fixed seating position, capturing an image when all employees are gathered in a cooperative manner before matching [2-4]. However, none of them addresses the issue of attendance marking in an unconstrained environment, due to the variations of scales, poses and occlusions in an uncooperative manner. In such scenarios, FR may suffer from poor accuracy. Even we can obtain discriminative facial features through the state-of-the-art FR model, how to decide the best threshold for practical use in a dynamic environment remains a challenge [5].

Received August 24, 2020; revised October 6, 2020; accepted November 1, 2020.
Communicated by Chih-Hung Wu.

When there is an issue on AI accuracy, people often relate this to the problem of Deep Learning (DL) accuracy. Although this is partially true, there are also other factors that may influence the overall accuracy, such as the extrinsic factors (*A*), camera factors (*B*), pre-processing modules (*C*), thresholding (*D*) and application (*E*) as depicted in Fig. 1.

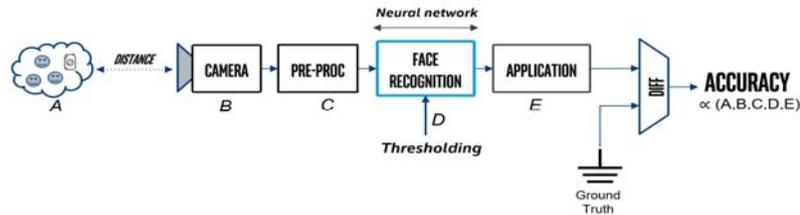


Fig. 1. Accuracy limiting factors.

Extrinsic factors (*A*) refer to the ambient conditions (*e.g.* lighting) and the conditions of the test subjects. In a real-life scenario, the test subjects are normally uncooperative. For example, the person-of-interest may not always look frontally to the camera. Camera factors (*B*) refer to the camera design, such as the Field-of-View (FOV), video settings (*e.g.* resolution, bitrate) and camera position. For example, the system integrators need to avoid backlight and low-light conditions while installing their cameras [6].

The inter-pupillary distance (IPD) value is commonly used to indicate the FR working range. A typical Face Recognition (FR) algorithm recommends a minimum IPD of 32 pixels [7] while the Face Detection (FD) can achieve lower IPD than FR. Fig. 2 shows that the relationship of a person's IPD with respect to the camera's HFOV, video resolution and camera distance. Based on Eq. (1), an IPD pixel of a person face can be estimated:

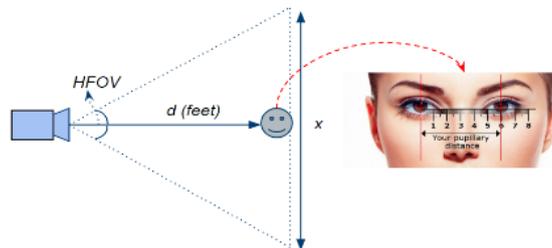


Fig. 2. The FR working range with the impact of FOV and video resolution.

$$x = 2d \tan \frac{\theta_{HFOV}}{2}, \quad IPD(\text{pixel}) = \frac{\text{vide width}(\text{pixel}) * IPD(\text{mm})}{x} \quad (1)$$

where the actual IPD range of a female is in between 51.0mm to 74.5mm, and a male is in between 53.0 and 77.0mm [8].

Figs. 3 (a) and (b) show the FD and FR working range based on a 4k resolution with different HFOV using colour code. It is noticed that the FR working range is reduced into halved (*i.e.* from 12 feet to 6 feet) by selecting a wider HFOV (*i.e.* from 90° to 120°). Although a wider camera view can improve the coverage of the captured scene, it impacts the FR working range significantly.

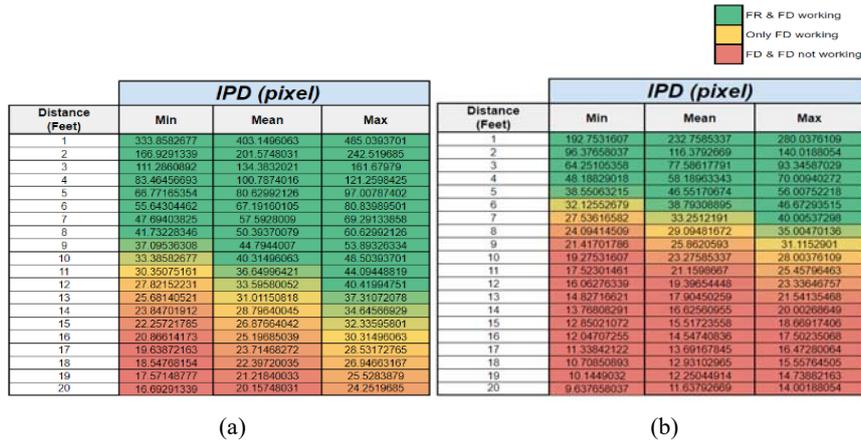


Fig. 3. Facial Recognition working distance of a 4k resolution with different HFOV; (a) HFOV = 90 degree; (b) HFOV = 120 degree.

Pre-processing (*C*), on the other hand, can be applied as an enhancement for the FR model. For example, pre-correcting a pose using 2D transformation can improve the FR accuracy [9].

For the deep learning model (*D*), selecting an optimum threshold to work well in actual environment is crucial. Figs. 4 (a) and (b) show the distribution of classes on any classification model for an ideal case and real-world case respectively. In an ideal situation, the two curves do not overlap at all where the FR model is able to distinguish between positive class (*i.e.* matched face) and negative class (non-matched face) using a fixed threshold. However, it is impossible to distinguish perfectly between positive and negative classes in many real-world applications with a fixed threshold in a real-world scenario especially in the classroom where the students are from the same age group, demographic and ethnicity. Although an adaptive threshold was proposed in [5], the evaluation was done on LFW test cases where the cases are not generalized well for classroom types of use cases.

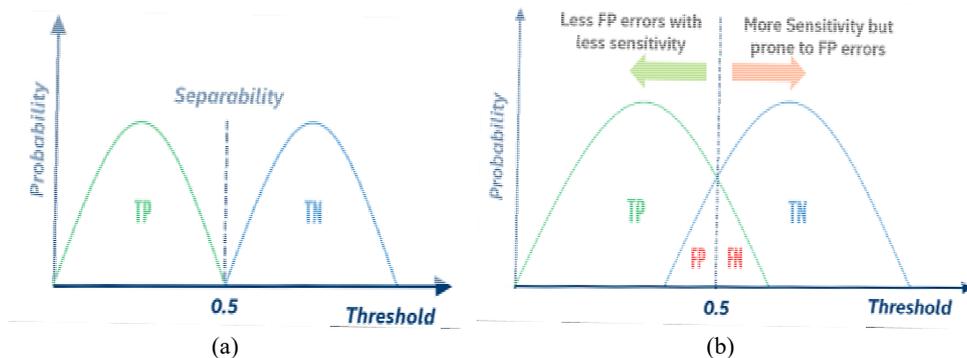


Fig. 4. Distribution of the classes on any classification model; (a) Ideal case; (b) Real-world case.

Finally, the application level (*E*) is the post-processing module to digest the inference results based on usage scenario. For example, the teachers are not concerned about the

precision and recall of FR model, but they are more interested in knowing the actual attendance in a classroom. Practically, it is impossible to conclude the student attendance based on the inference result from a single frame where the students are under uncooperative manner. Therefore, the conventional frame-based FR accuracy metrics cannot address the real-world problem for attendance marking [10]. In this paper, our motivation is to address the issue of the existing frame-based attendance marking and the manual fine-tuning of the thresholding during deployment. We compare the conventional frame-based result with our proposed methodologies to show that the superiority of the statistical multiframe and mean thresholding approaches under the unconstrained classroom environment.

2. CONVENTIONAL FACIAL RECOGNITION ACCURACY METRICS

False Positive (FP) and False Negative (FN) are the two types of prediction errors in FR system. In contrast, FP is more damaging than FN as recognizing an absentee as presence is more undesirable. Precision and recall [11] are the two fundamental metrics to define the image or frame-level accuracy. These metrics can be calculated based on Eq. (2). The mean average precision (mAP) [12] is a measure that combines both precision and recall for evaluating the accuracy of the FR model.

$$Precision(frame) = \frac{TP}{TP+FP}$$

$$Recall(frame) = \frac{TP}{TP+FN}$$

$$Accuracy(frame) = \frac{TP}{TP+FP+FN} \quad (2)$$

where, TP is the total number of matched faces in an image. FP is the total number of false accepted faces in an image. FN is the total number of false negative faces in an image.

While the frame-based accuracy is useful for training the FR model, but it fails to reflect the true attendance in a classroom as the result may be fluctuated over time due to large variations of poses and occlusions. Therefore, a statistical multiframe accuracy methodology for attendance marking is desirable.

3. PROPOSED METHODOLOGY FOR STUDENT MARKING

Attendance marking is more than just a facial recognition. A correct definition for attendance marking should be detecting and tracking of an individual and deciding the presence of a known person over an observation time. Fig. 5 illustrates this concept by identifying the presence of a known person based on the multiframe inference results. Here is the boy who played skateboard in the video. The FR inference result varies as the person was moving over the frames. In general, a higher precision FR model would produce a relatively higher TP than FP in terms of frames. In this example, majority of the frames are predicted correctly with a fewer prediction error. The attendance marking application is designed to decide if this boy is either Alan or Tom that is available in the facial database.

Since Alan has received a higher voting than Tom, therefore, the system concludes that this person is Alan and not Tom.

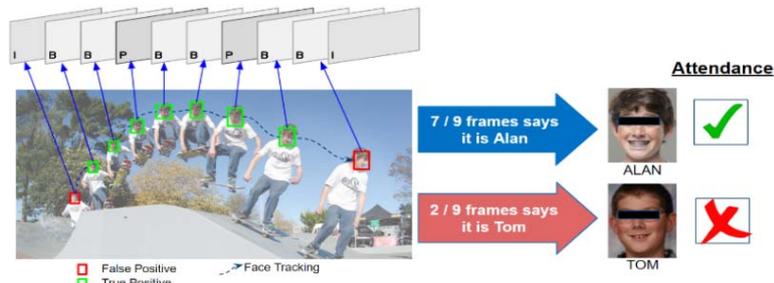


Fig. 5. Statistical multiframe concept for attendance marking.

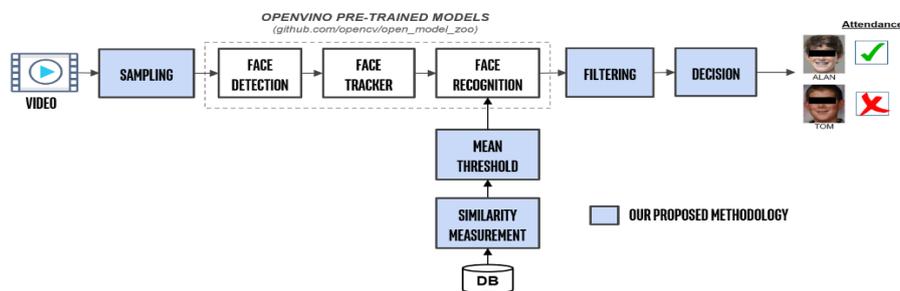


Fig. 6. Proposed attendance marking methodology.

The proposed statistical multiframe accuracy methodology with mean threshold for attendance marking is shown in Fig. 6. An FR system should contain three basic elements, namely face detection, face tracking and face recognition. To convert FR system into a student marking application, three additional blocks were implemented, namely sampling, filtering and decision. The inference results from the FR model is known as samples. The filtering block removes the excessive noise (*i.e.* FP and FN) from the samples. Finally, the decision block determines the presence of each tracked person based on the statistical analysis. Additionally, we introduce a thresholding mechanism based on mean to provide an optimum balancing between FP and FN based on the given facial database.

3.1 Sampling

The sampling frequency or inferencing rate is dependent on the network structure of the AI model, the types of hardware and the types of workload that runs in the hardware. It is ideal to infer as high as the video frame rate (*e.g.* 30 FPS) to provide better statistical analysis for accuracy. However, it comes with a cost of expensive HW, such as the neural network acceleration card. With the hardware constraint, a careful selection of the thresholding can help to reduce the inference rate as low as 3 FPS while achieving the same accuracy as full rate at 30 FPS. The explanation of the thresholding scheme will be covered in Section 3.4.

3.2 Filtering

It is unavoidable to have both FP and FN prediction errors under an unconstrained environment [13]. Sliding window (aka median filter) is one of the effective methods to filter out the random noise [14] for image processing. Here, we proposed the use of sliding window filtering to remove excessive of FP and FN from the inference samples. Fig. 7 illustrates the concept of sliding window filtering to improve TP numbers in the inference samples. The sliding window length can be adjusted for optimum result depending on the usage scenarios. A larger window size can be chosen to provide a better smoothing effect, providing the student movement in a classroom is minimal. With the use of sliding window, the system can reduce the FP and FN samples, and subsequently improving the statistical calculation for attendance decision making.



Fig. 7. Sliding window filtering.

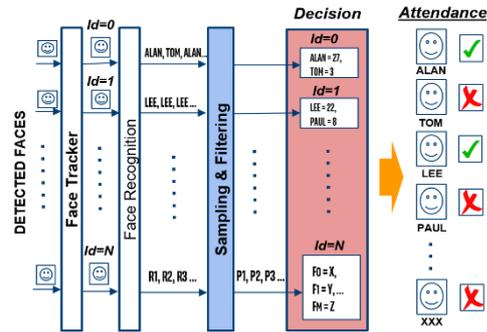


Fig. 8. Decision block based on majority voting.

3.3 Decision

The FR system may recognize more than one candidate per tracked face in an unconstrained classroom environment. The majority voting is proposed to decide the best candidate as the true attendance for each tracked face. The decision block updates the voting number per tracked face after each inference as illustrated in Fig. 8. The one with the highest voting will be selected as the true candidate who presents in the class for each tracked face based on Eq. (3).

$$Decision(per\ tracked\ face) = Majority(F_0, F_1, \dots, F_n) \quad (3)$$

where F is the number of occurrences per candidate and n is the total number of candidates per tracked face.

3.4 Mean Thresholding

A thresholding plays an important to maximize the recognition rate of a group of students by balancing between the false positive rate and false negative rate. Since the attendance marking system has a prior knowledge of the facial database of a specific classroom, we can leverage it to understand the similarity distribution within a class. Therefore, a mean value can be estimated from the similarity measurement to strive a good balance between FP and FN.

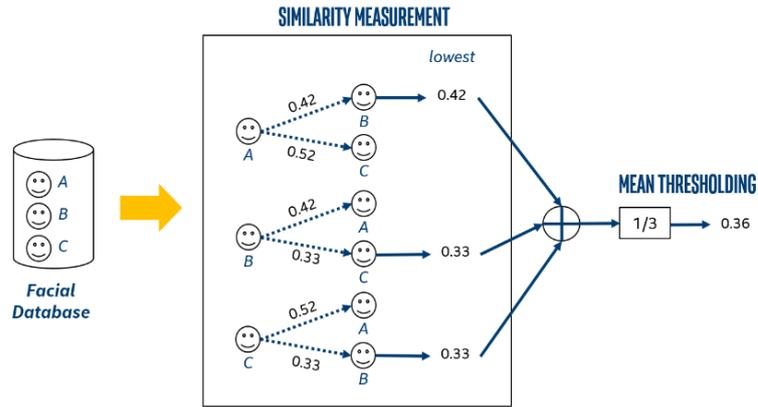


Fig. 9. Mean thresholding methodology.

As shown in Fig. 9, each registered face is compared with the rest of the facial embeddings in the database for cosine distance and the one with the lowest distance will be selected. The cosine distance calculation is given in Eq. (4). A mean thresholding is estimated by averaging the selected distance over all the faces in the database based on Eq. (5).

$$\text{cosine distance}(two\ vectors) = 1 - \cos \theta = 1 - \frac{AB}{\|A\| \|B\|} \tag{4}$$

where A and B represent two set of face embeddings or vectors and θ is the angle between the two vectors.

$$\text{Mean Threshold(per database)} = (T_0, T_1, \dots, T_n)/n \tag{5}$$

where, T is minimum cosine distance per registered face by comparing all the embeddings in the database. n is the total number of registered faces in the database.

4. PERFORMANCE EVALUATION

An accuracy tool for attendance marking was developed and evaluated using OpenVINO. Some internal test video sequences were used for the evaluation as shown in Fig. 10.



Fig. 10. Classroom test videos with different classroom configuration. (To protect privacy of students, we hide their eyes in the paper)

Each of the video length is 2 minutes with a frame rate of 30 FPS. It contains 48 students in the classroom. The students were uncooperative in the video, where it contains many student activities, such as hand-raising, standing, sitting, reading books, sleeping, and chatting. The “face-detection-adas-0001” [15] and the “face-reidentification-retail-0095” [16] pre-trained models were used and the sample codes are available in [17]. These pre-trained models have been optimized for Intel processor. We evaluated both the accuracy of the frame-based and the proposed statistical multiframe approach for attendance marking.

Figs. 11 (a) and (b) show the frame-based student attendance results over 120 seconds or 3600 frames, with the cases of without and with sliding window. Blue curves indicate the true attendance (aka TP), and red curves indicate the false attendance (aka FP). The student attendance was significantly improved after applying the sliding window filtering. However, it is still impossible to get a full student attendance in any of the frame even with the support of sliding window filtering. Additionally, the results also consistently show one to two false attendees throughout the whole video sequence, which are unable to be removed by the sliding window filtering.

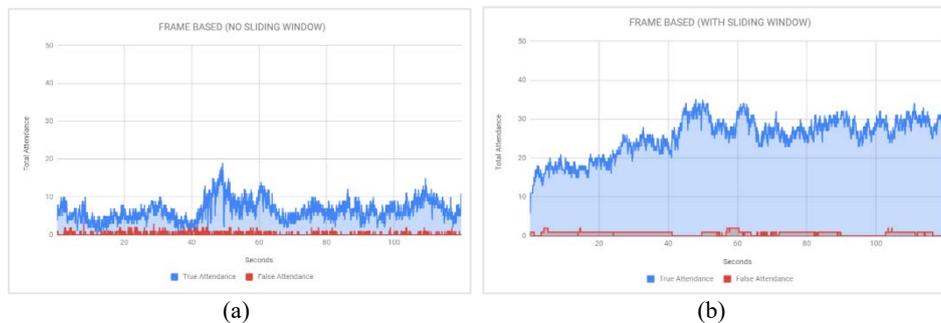


Fig. 11. Frame-based student analysis (a) no sliding window; (b) with sliding window.

Fig. 12 shows the result of the proposed statistical multiframe analysis with sliding window filtering and majority voting. The true attendance rate was able to converge from 6 to 40 students after 50 seconds of convergence time. For false attendance rate, it takes slightly longer (*i.e.* 110 seconds) to minimize it to zero.

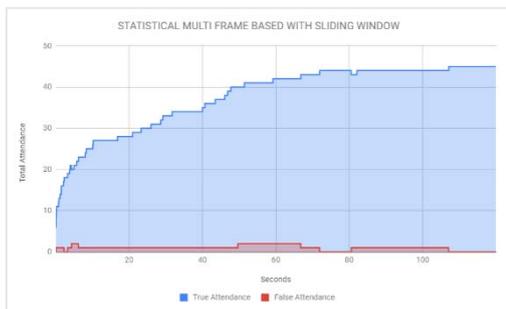


Fig. 12. Statistical multiframe analysis for attendance marking.

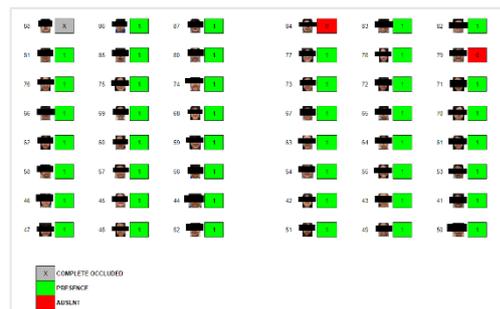


Fig. 13. The classroom attendance after 110 seconds of convergence time.

Fig. 13 shows the status of the attendance marking tool according to their seating position after 110 seconds of convergence time. With the proposed statistical multiframe accuracy methodology, the system is able recognize 45 out of 48 students correctly within 2 minutes (*i.e.* the accuracy of 93.75%) by inferencing the same rate as the video rate. The other 3 unrecognized students (Id = 88, 84 & 79) suffer from low-resolution and severe occlusion, therefore, the statistical multiframe approach cannot improve in such conditions.

Furthermore, we evaluated the three thresholding methods, namely fixed, adaptive and mean threshold. All three thresholding methods were combined with the statistical multiframe methodology for the attendance marking accuracy. A fixed threshold of 0.43 was used, as this value was estimated over a set of test videos to achieve best balance between the positive and negative cases using the conventional frame-based accuracy metric. The adaptive threshold was calculated based on the proposed method in [5], and the mean threshold was calculated based on the proposed method in this paper.

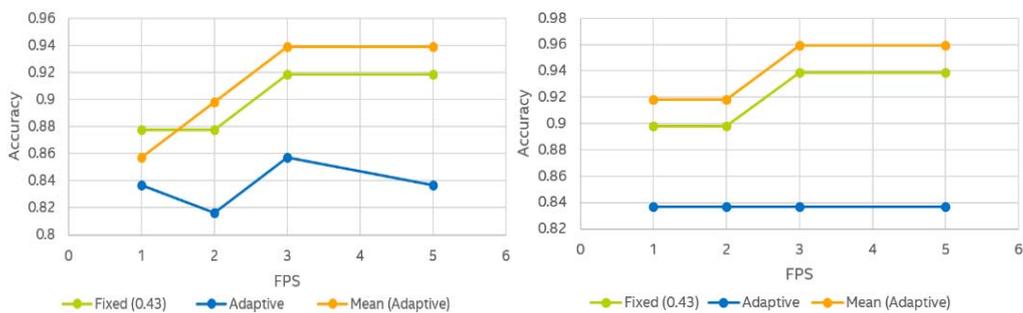


Fig. 14. Thresholding comparison using with different classroom videos at different inference rate FPS.

Fig. 14 shows the comparison of thresholding with different classroom videos. Unfortunately, the adaptive threshold was the worst among all the cases. Although adaptive threshold outperforms fixed threshold on LFW test cases [5], it underperforms under an unconstrained classroom environment. We found that the faces with the stringent threshold have the hard time for recognizing them especially for those who sit in the last few rows of the classroom while the faces with the lenient threshold introduce more FP than TP. Thus, adaptive threshold is not suitable for classroom types of use case.

Overall, the mean threshold gives the best attendance marking accuracy among the three different thresholding methods. It shows its capability to achieve a good balance TP and FP during the majority calculation so that more genuine faces could be recognized. The results also show that the mean threshold can maximize the accuracy level in the classroom even with a lower inference rate as low as 3 FPS. On the other hand, fixed threshold is comparable to mean threshold, but it requires some efforts to estimate an optimum threshold point.

5. CONCLUSIONS

The evaluation results show that the conventional frame-based accuracy metric cannot address the issue for attendance marking due to the occlusion and variation of students'

poses in an uncooperative manner. The proposed statistical multiframe with sliding window filtering and majority voting show the improvement of attendance marking accuracy after a convergence time (*i.e.* 110 seconds). On the other hand, the proposed mean thresholding is more realistic and systematic approach for calculating threshold value based on the prior knowledge of a specific classroom facial database. When the proposed thresholding scheme is combined with the statistical multiframe methodology, the system is able to reduce the inference rate to as low as 3 FPS while achieving the same attendance accuracy at full inference rate (*i.e.* 30 FPS). This drives away the need to invest an expensive hardware to maintain the same accuracy with a higher inference rate.

REFERENCES

1. M. Hassaballah and S. Aly, "Face recognition: challenges, achievements and future directions," *IET Computer Vision*, Vol. 9, 2015, pp. 614-626.
2. M. Sajid, R. Hussain, and M. Usman, "A conceptual model for automated attendance marking system using facial recognition," in *Proceedings of the 9th International Conference on Digital Information Management*, 2014, pp. 7-10.
3. K. Yohei, T. Shoji, W. Lin, K. Kakusho, and M. Minoh, "Face recognition-based lecture attendance system," in *Proceedings of the 3rd AEARU Workshop on Network Education*, 2005, pp. 1-5.
4. K. MuthuKalyani and A. VeeraMuthu, "Smart application for AMS using face recognition," *Computer Science and Engineering*, Vol. 3, 2013, pp. 13-20.
5. H. Chou, J. Lee, Y. Chan, and C. Chen, "Data-specific adaptive threshold for face recognition and authentication," in *Proceedings of IEEE Conference on Multimedia Information Processing and Retrieval*, 2019, pp. 153-156.
6. F. Wheeler, X. M. Liu, and P. H. Tu, "Face recognition at a distance," *Handbook of Face Recognition*, 2nd ed., Springer, London, pp. 358-381.
7. Neurotechnology, "Verilook SDK product page – recommendation for minimal distance between eyes," <https://www.neurotechnology.com/verilook-technical-specifications.html>.
8. Wikipedia, "Pupillary distance for male and female," https://en.wikipedia.org/wiki/Pupillary_distance.
9. X. J. Chai, S. G. Shan, and W. Gao, "Pose normalization for robust face recognition based on statistical affine transformation," in *Proceedings of Joint Conference of the 4th International Conference*, Vol. 3, 2003, pp. 1-5.
10. K. H. Lee, S. V. Addicam, I. Krylov, S. Nosov, M. S. Lai, Z. Q. Lee, and C. S. Chai, "Statistical multiframe accuracy methodology for attendance marking system," in *Proceedings of International Conference on Technologies and Applications of Artificial Intelligence*, 2019, pp. 1-5.
11. M. Sundaram and A. Mani, "Face recognition: demystification of multifarious aspect in evaluation metrics," *Face Recognition Book*, Chapter 5, 2016.
12. K. Oksuz, B. Cam, E. Akbas, and S. Kalkan, "Localization recall precision (LRP): a new performance metric for object detection," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 504-519.
13. R. Verschae, J. Ruiz-del-Solar, and M. Correa, "Face recognition in unconstrained en-

- vironment: a comparative study,” *Journal on Advances in Signal Processing*, 2008, pp. 1-12.
14. L. Tan and J. Jian, “Median filter,” *Digital Signal Processing*, 2nd ed., Chapter 14, 2013, pp. 712-713.
 15. Intel OpenVINO, “Face detection pre-trained model based on mobilenet backbone,” https://docs.openvino toolkit.org/latest/_intel_models_face_detection_adas_0001_description_face_detection_adas_0001.html.
 16. Intel OpenVINO, “Face recognition pre-trained model based on mobilenet V2 backbone,” https://docs.openvino toolkit.org/latest/_intel_models_face_reidentification_retail_0095_description_face_reidentification_retail_0095.html
 17. Intel OpenVINO, “Open model zoo github repository,” https://github.com/opencv/open_model_zoo.



Kuan Heng Lee (李冠興) is a platform solutions architect in Intel, who is currently working closely with the CTO office in driving AI accuracy redefinitions for RBHE (Retail, Banking, Hospitality and Education). Previously he was the Principal Software Engineer in Motorola Solutions, focusing on the Connected Police Officer applications, such as live video streaming, mobile facial recognition, gun-holster triggered recording and *etc.* Prior to Motorola, he also worked in Ericsson on the video compression algorithms improvement. More details can be found: <https://my.linkedin.com/in/kuan-heng-lee-8434153a>.



Sanjay Addicam is an Intel Sr Principal Engineer with a focus on video analytics, Deep learning and data mining algorithms. He has around 30 pending patents and numerous papers in tier 1 conferences like KDD. He is the author of the book *Building Intelligent Systems: utilizing computer vision, data mining and machine learning*. His current focus is on sensor data fusion involving RFID and computer vision and creating Edge systems which can auto generate ground truth data and deploy deep learning models without any manual intervention. He is also a Kaggle master.



Ilya Krylov is a Senior Deep Learning Engineer in Intel’s Internet of Things Group. He works in a team of deep learning engineer developing computer vision algorithms that can be run at the edge. He contributes to the OpenVINO toolkit, especially the Open Model Zoo and OpenVINO Training Extensions. Ilya earned bachelor’s and master’s degrees, both in Mathematical and Computer Science, from the State University of Nizhni Novgorod, Russia.



Sergei Nosov is a Software Engineering Manager at ICV leading one of the algorithm teams developing Intel Models for OpenVINO toolkit. Sergei's interests are in computer vision, deep learning, computer science and programming languages. Sergei has about 8 years of experience developing computer vision algorithms including panorama stitching, camera calibration and object detection. His current area of R&D lies in deep learning with algorithms such as object detection and semantic segmentation, attributes classification and image processing. Sergei holds a Master of Science degree from Nizhny Novgorod State University in Computer Science.



Lai Mee Sim (赖美璇) is the Platform Solution Architect of Intel IoTG who responsible for bringing Artificial Intelligence into Education space with the latest Intel technologies. Her focus is to create new AI innovations with Intel building blocks, advocate the values and making sure it meets the needs of education segment. She also works very closely with the education ecosystem globally to help connect customers to technology providers and help promote solutions that help uplift teaching and learning experience with goal of education: Personalization, Equality, Efficiency. Prior to this, she has 12 years of silicon, hardware and software customers support and enabling experience based on Intel architecture and IXP architecture. Michelle holds a master's degree in computer science from the University of Hertfordshire, UK.



Zhan Qiang Lee (李展强) is currently a Platform Solution Engineer with Internet of Things Group (IOTG) at Intel with focusing in building solution exercising AI at the edge. His interests include vision based deep learning, media, and security technologies. He holds a bachelor's degree in Electronic engineering from Multimedia University Malaysia (MMU) and an MBA degree from University Utara Malaysia (UUM).



Chai Chung Shien (蔡镇宇) joined Altera in 2006 as a new college graduate from University Technology of Malaysia. Since then, he's established his reputation with significant contributions to FPGA and flash software configuration. He has accumulated in depth knowledge for configuration and has the holistic view of the complete configuration chain. With that, he delivered bitstream and programmer solutions to enable > 5 FPGA product generations to configure successfully. His recent contribution is on S10 bitstream assembler tool; a key component in S10 Configuration flow to ensure successful delivery of Stratix10 early access program (EAP) devices to customers. He has 1 patent granted, 2 pending, 1 poster presentation at Altera Technical Symposium (ATS) 2010 and awarded with best presenter for ATS 2015.