

# Revolutionising Attendance Face Recognition Systems

ASHISH KUMAR<sup>1</sup>, AKASH DWIVEDI<sup>2</sup>, SURESH KUMAR TIWARI<sup>3</sup>, DR. SANJAY PACHAURI<sup>4</sup>

<sup>1,2</sup>B. Tech 4th Year (Computer Science & Design) GNIOT, Greater Noida, India  
Department of Data Science & Design GNIOT, Greater Noida, India

**Abstract** — *Reliable documentation of attendance remains a persistent operational burden in educational institutions and corporate settings alike. Conventional solutions — spanning handwritten roll calls to token-based card readers — are routinely compromised by proxy attendance, labour-intensive record keeping, and limited scalability. Growing dissatisfaction with these limitations motivates the development of intelligent, low-infrastructure alternatives rooted in contemporary AI techniques. This paper proposes an automated, touchless attendance capture system founded on a sequential three-layer deep learning pipeline. The detection layer employs a Multi-Task Cascaded Convolutional Network (MTCNN) to scan live video and isolate individual faces along with associated anatomical landmarks. The representation layer then processes the geometrically normalised face crops through a FaceNet model built on an Inception-ResNet-v1 backbone, producing a compact 128-dimensional feature vector per face via triplet-loss-driven metric learning. Finally, the decision layer applies a Support Vector Machine with a Radial Basis Function kernel to associate each embedding with a registered identity that was enrolled offline. Evaluation spanned two distinct data sources: the publicly accessible Labeled Faces in the Wild (LFW) benchmark alongside an institution-specific corpus drawn from 50 volunteers photographed across varied lighting environments. The pipeline achieved 98.7% accuracy on the local dataset and 99.1% on LFW, with per-frame processing completing in 120 milliseconds. The False Acceptance Rate was confined to 0.8%. A Flask-powered administrative dashboard facilitates live monitoring and automated report export.*

**Keywords** — *touchless attendance, facial identification, MTCNN, FaceNet embeddings, SVM classifier, deep metric learning, real-time recognition, edge deployment, biometric authentication.*

## I. INTRODUCTION

Despite how routine the act of taking attendance appears, it remains an unsolved operational challenge at institutions ranging from primary schools to large enterprises. Conventional paper-based rolls invite falsification, while proximity card terminals can be gamed by passing a credential to an absent colleague.

Beyond these security gaps, manual data compilation consumes valuable administrative time with no corresponding gain in institutional effectiveness. There is a compelling case for a passive, accurate, non-contact solution.

Facial biometrics offer a compelling avenue. Unlike fingerprints or retinal scans, faces can be captured at a natural stand-off distance without requiring deliberate interaction from the subject. The maturation of GPU-accelerated deep learning and freely available open-source toolkits has made real-time face processing practical on hardware that fits within the budget of a typical department rather than a specialised laboratory.

Academic interest in computational face analysis has grown substantially over the past decade. Contemporary architectures trained on tens of millions of images achieve identification benchmarks that match or exceed human-level performance in controlled settings. Translating these capabilities into a functional attendance solution requires integrating the recognition model within a broader engineering pipeline that handles live video capture, face alignment, identity management, and structured reporting.

The framework introduced here addresses every link in that chain. By connecting MTCNN-based face detection, FaceNet-based embedding generation, and SVM-based identity assignment, the system produces an end-to-end workflow that writes confirmed attendance events to a relational database and presents them through a web interface. The design prioritises accuracy, real-time throughput, and deployment viability on commodity hardware.

### *Primary Contributions:*

- A fully contactless, camera-driven attendance pipeline that eliminates any requirement for physical interaction from participants.
- A sequentially connected detection–alignment–encoding–classification chain that sustains high

identification rates under realistic lighting and pose variation.

- Quantitative validation against both a self-assembled institutional dataset and a widely adopted public benchmark, demonstrating 98.7% accuracy.
- A production-ready implementation tested on standard desktop workstations and embedded computing boards.

## II. LITERATURE REVIEW

Efforts to automate attendance tracking have evolved through successive generations of technology. Early digital systems traded paper registers for magnetic stripe cards and, subsequently, RFID proximity tokens. Although these approaches reduced transcription errors, they failed to bind the credential to a specific individual — any possessor of the card could register attendance on behalf of the enrolled holder [1]. Fingerprint readers were later introduced as a biometric countermeasure, tying presence to a physiological attribute. However, repeated contact with sensor surfaces raises sanitary concerns in shared environments, and minor skin irregularities such as dryness or superficial wounds frequently produce erroneous rejections [2].

Computational face recognition entered the research mainstream during the 1990s. Turk and Pentland's Eigenface formulation projected face imagery into a low-dimensional principal component subspace, enabling recognition on hardware of the era [3]. Belhumeur and colleagues subsequently proposed a class-discriminant variant — Fisherfaces — that more effectively distinguished individuals who shared similar illumination conditions [4]. Local Binary Pattern Histogram descriptors advanced lighting robustness further by encoding neighbourhood contrast relationships rather than raw pixel values [5]. These handcrafted feature families, however, possessed an inherent representational ceiling that prevented them from generalising to unconstrained, real-world acquisition scenarios.

Deep convolutional architectures fundamentally altered the field's trajectory. Taigman et al. demonstrated that a nine-layer network trained with carefully aligned face images could approach human verification accuracy on the LFW benchmark [7]. Schroff and colleagues subsequently introduced triplet loss as a training objective that directly shapes the geometry of a learned embedding space, causing

same-person images to cluster while pushing apart images of distinct individuals [8]. Their FaceNet system set a new performance ceiling that remained influential for subsequent years. Deng et al. refined this further by incorporating an angular margin penalty — ArcFace — that tightens within-class distributions and widens inter-class separation at the hypersphere surface [9].

Concurrent work in face detection led to the MTCNN framework of Zhang et al. [10], which organises three convolutional stages in a coarse-to-fine cascade: early stages rapidly reject background regions, while the final stage refines bounding boxes and localises five facial landmarks with sufficient precision to guide affine alignment. Its efficiency and reliability have made it a preferred detection front-end for many downstream recognition pipelines.

In the attendance domain specifically, prior studies have demonstrated functional but limited systems. A convolutional classifier achieved identification rates marginally above 96% under controlled laboratory conditions with a fixed overhead camera [11]. A subsequent investigation replaced the classification back-end with a k-nearest-neighbour decision rule and observed improved throughput on small groups, though accuracy degraded under varying outdoor illumination [12]. A comprehensive review of face recognition methodologies concluded that deep learning pipelines consistently outperform classical approaches on every major evaluation criterion [13]. Nevertheless, a recurring finding across this literature is a substantial gap between laboratory demonstration and real-world deployability: most prototypes assume uniform lighting, cooperative participants, and modest population scales. The present work addresses those overlooked deployment conditions.

## III. DESIGN METHODOLOGY

The proposed system is organised as a sequential processing chain in which each component receives the output of its predecessor and passes a refined result to the next. Four principal modules — video ingestion, face region extraction, identity encoding, and attendance recording — are described in the subsections below.

### *A. Video Ingestion and Frame Conditioning*

A USB or CSI camera delivers a continuous video stream to the host at thirty frames per second. Prior

to any face-specific computation, each frame undergoes three preparatory operations: uniform spatial rescaling, pixel intensity normalisation to the numeric range expected by downstream neural networks, and light Gaussian blurring to suppress sensor-level noise that would otherwise introduce misleading gradient information into the detection stage.

#### *B. Facial Region Extraction via MTCNN*

The conditioned frame is submitted to MTCNN [10], a three-network cascade designed to balance detection speed against localisation precision. The first sub-network (P-Net) scans the image at multiple resolution scales using a lightweight convolutional filter, yielding a large pool of candidate face proposals. The second sub-network (R-Net) evaluates each candidate and discards the vast majority as background, while refining bounding box coordinates for the survivors. The third sub-network (O-Net) performs a final classification pass and simultaneously regresses the pixel coordinates of five anatomical landmarks — both eye centres, the nasal tip, and the two mouth corners. These landmark positions drive an affine warp that corrects in-plane rotation and crops each face to a standardised  $160 \times 160$  pixel representation.

#### *C. Identity Encoding with FaceNet*

Each normalised face patch is forwarded to a FaceNet model whose backbone follows the Inception-ResNet-v1 topology [8]. The network was pre-trained on the VGGFace2 corpus — exceeding three million images across thousands of subjects — using triplet loss as the training criterion. Triplet loss simultaneously reduces the embedding distance between two images of the same person (an anchor-positive pair) and enlarges the distance to an image of a different person (a negative). The resulting 128-dimensional vector encodes facial identity compactly, such that Euclidean distance between vectors is a reliable proxy for facial similarity. During an offline enrolment phase, representative embedding vectors for each registered participant are computed and stored alongside identity metadata.

#### *D. Classification and Identity Decision*

At inference time, the query embedding is presented to a Support Vector Machine trained with a Radial Basis Function kernel on stored enrolment embeddings. The SVM assigns the query to the most probable registered identity. Each prediction is

accompanied by a confidence score; only predictions meeting or exceeding a 95% confidence threshold are accepted as positive identifications. Queries falling short of this threshold are directed to a separate anomaly log for subsequent manual review, ensuring that unverified individuals are neither silently admitted nor incorrectly rejected without supervisory awareness.

#### *E. Attendance Recording and Reporting*

A confirmed identification triggers an atomic write to a MySQL database table containing the participant identifier, full name, session identifier, and a server-generated UTC timestamp. A session-scoping guard prevents duplicate entries: once an individual has been logged within a configurable time window (typically the duration of a class or shift), subsequent detections of the same face are acknowledged internally but not re-written to the database. At session close, formatted reports are exported in both CSV and PDF formats. A Flask-based dashboard running on the same host delivers live status feeds and allows supervisors to query historical records without direct database interaction.

#### *F. Processing Pipeline Summary*

Step 1 — Camera delivers a raw frame to the ingestion module.

Step 2 — Frame conditioning standardises resolution, brightness range, and noise level.

Step 3 — MTCNN detects and landmark-aligns face crops within the conditioned frame.

Step 4 — FaceNet converts each aligned crop into a 128-dimensional embedding vector.

Step 5 — SVM assigns an identity label with an associated confidence score.

Step 6 — Detections surpassing the 95% threshold are committed to the attendance database.

Step 7 — Sub-threshold detections are routed to the anomaly review log.

Step 8 — Dashboard and report generator reflect all new entries in real time.

## IV. BLOCK DIAGRAM OF THE PROPOSED DESIGN AND WORKING

Fig. 1 illustrates the architectural block diagram of the proposed system. Processing flows from raw video input at the top toward final database storage at the bottom, with a lateral branch channelling unrecognised frames into a separate anomaly handling path. Each block operates as a self-

contained unit communicating with adjacent blocks through well-defined data interfaces, so that individual components can be upgraded or replaced without restructuring the surrounding pipeline.

[Fig. 1. Architectural block diagram of the proposed face recognition attendance system.]

The first decision node following the MTCNN stage governs downstream routing: frames containing no detectable face are discarded and recorded as empty frames, while frames yielding at least one detection proceed to landmark alignment and embedding generation. After the SVM produces its verdict, a second decision node applies the confidence threshold, separating confirmed identifications from uncertain queries before any database write is authorised.

## V. RESULTS AND DISCUSSION

### A. Implementation Environment

All experiments ran on a workstation configured with Ubuntu 20.04, an Intel Core i5-10400 processor, 8 GB DDR4 RAM, and an NVIDIA GeForce GTX 1650 GPU carrying 4 GB of dedicated VRAM. The software environment consisted of Python 3.8, TensorFlow 2.9, OpenCV 4.6, and scikit-learn 1.1. Two datasets were employed: the publicly accessible LFW benchmark and a locally assembled collection of 50 volunteers who each contributed 30 photographs captured under a range of controlled and uncontrolled lighting scenarios.

### B. Evaluation Criteria

Performance was quantified using six standard metrics. Identification Accuracy captured the proportion of test queries correctly attributed to the enrolled individual. Precision and Recall jointly characterised the balance between over-acceptance and under-acceptance errors. The F1-Score synthesised these into a single harmonic mean. The False Acceptance Rate (FAR) quantified how often an impostor query was incorrectly confirmed, while the False Rejection Rate (FRR) measured how often a legitimate query was incorrectly refused.

### C. Accuracy Outcomes

Against the locally assembled dataset, the pipeline delivered an identification accuracy of 98.7%, while the LFW evaluation yielded 99.1%. The modest differential between the two figures reflects the broader subject diversity and more variable acquisition conditions represented in LFW relative to the locally controlled subset. Across both datasets, Precision reached 98.4% and Recall 98.9%, yielding an F1-Score of 98.6%. These figures held even when evaluation images were acquired under fluorescent overhead lighting substantially different from the daylight conditions used during enrolment.

### D. Comparative Evaluation

Table I contextualises the proposed framework relative to four previously published methods, selected to span the spectrum from classical machine learning to modern deep learning baselines.

TABLE I  
 Benchmarking Against Prior Attendance Recognition Systems

Method	Approach	Data	Accuracy	Latency
Ahuja et al. [1]	PCA + SVM	Custom	91.2%	380 ms
Guo & Zhang [2]	CNN Classifier	Custom	96.4%	210 ms
Pandey et al. [3]	FaceNet + KNN	LFW	97.1%	175 ms
Kortli et al. [4]	DeepFace	VGGFace2	97.8%	160 ms
This Work	MTCNN+FaceNet+SVM	LFW+Custom	98.7%	120 ms

### E. Throughput and Latency

The mean end-to-end processing time per frame — spanning raw pixel ingestion through to database confirmation — measured 120 milliseconds in CPU-only mode and approximately 38 milliseconds with GPU acceleration enabled. The CPU throughput of roughly eight identifications per second is sufficient to manage a queue of students entering a lecture hall.

GPU-enabled throughput of approximately 26 frames per second supports real-time crowd monitoring without frame skipping.

### F. Security Threshold Sensitivity

At the default 95% confidence threshold, the system recorded a FAR of 0.8% and an FRR of 1.3%. Raising the threshold to 98% reduced the FAR to

0.3% at the cost of increasing the FRR to 2.9% — a trade-off appropriate for high-security contexts. Conversely, lowering the threshold to 90% nearly eliminated rejections of legitimate users but elevated the FAR to 2.1%, which may be acceptable in low-stakes scenarios such as informal session check-ins.

*G. Behaviour Under Adverse Conditions*

Supplementary tests applied controlled degradations: synthetic occlusion emulating facial covering, illumination restricted to a single low-power desk lamp, and deliberate head tilts of up to 40 degrees. The occlusion scenario produced the most pronounced accuracy decline, approximately 4.8 percentage points below baseline. Extreme lateral

head tilt reduced accuracy by 2.3 percentage points. Dim-light conditions without supplementary illumination incurred a 3.1 percentage-point drop.

*H. Administrative Efficiency*

Over a ten-day operational pilot involving 60 enrolled participants, the logging module registered zero duplicate entries across 3,400 individual attendance events. Report generation for the full cohort completed in under two seconds. Participating administrators noted a substantial reduction in time spent compiling end-of-week attendance summaries compared with the manual spreadsheet process the system replaced.

TABLE II  
*Overall System Performance Summary*

Metric	Observed Value
Identification Accuracy (in-house dataset)	98.7%
Identification Accuracy (LFW benchmark)	99.1%
Precision	98.4%
Recall	98.9%
F1-Score	98.6%
False Acceptance Rate (FAR)	0.8%
False Rejection Rate (FRR)	1.3%
Per-frame Latency (CPU)	120 ms
Per-frame Latency (GPU)	38 ms
Enrolment Images per Participant	20–30
Report Generation Time (60 users)	< 2 s

VI. CONCLUSION

This paper has introduced and validated a touchless, camera-driven attendance recording system that supersedes manual and token-dependent processes with a fully automated deep learning workflow. The sequential orchestration of MTCNN for face detection, FaceNet for identity embedding, and SVM for classification produces a practical equilibrium among recognition accuracy, inference speed, and hardware accessibility. Sustained accuracy exceeding 98% across two independent datasets — one locally compiled and one drawn from a broadly used public benchmark — establishes that the approach generalises beyond the controlled boundaries of a single experimental context.

The system's administrative subsystem automates the functions most susceptible to human error: timestamp assignment, duplicate event prevention, anomaly flagging, and report distribution. A ten-day pilot confirmed that these mechanisms function reliably at realistic cohort scales, delivering concrete reductions in administrative overhead.

Several directions remain open for future investigation. Recognition under partial facial occlusion — notably from face coverings — warrants dedicated attention, potentially through occlusion-aware training strategies or the incorporation of complementary modalities such as periocular or gait-based cues. Vision transformer backbones, which have demonstrated strong recent performance on facial tasks, represent a promising alternative

encoding architecture. Liveness detection mechanisms that distinguish genuine faces from printed photographs or video replay attacks would materially strengthen the system against deliberate circumvention. Finally, model compression methods — including quantisation, structured pruning, and knowledge distillation — could bring the full pipeline within the memory and compute envelope of low-cost microcontroller boards.

#### REFERENCES

- [1] A. K. Jain, A. Ross, and S. Prabhakar, "An introduction to biometric recognition," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 14, no. 1, pp. 4–20, Jan. 2004.
- [2] D. Kumar, C. S. Rai, and S. Kumar, "A novel approach for attendance management system using RFID and biometric," *Int. J. Comput. Appl.*, vol. 58, no. 2, pp. 18–23, Nov. 2012.
- [3] M. Turk and A. Pentland, "Eigenfaces for recognition," *J. Cognitive Neurosci.*, vol. 3, no. 1, pp. 71–86, Jan. 1991.
- [4] P. N. Belhumeur, J. P. Hespanha, and D. J. Kriegman, "Eigenfaces vs. Fisherfaces: Recognition using class specific linear projection," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 19, no. 7, pp. 711–720, Jul. 1997.
- [5] T. Ahonen, A. Hadid, and M. Pietikainen, "Face description with local binary patterns: Application to face recognition," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 28, no. 12, pp. 2037–2041, Dec. 2006.
- [6] R. Ahuja, A. Jain, and M. Sharma, "Face recognition based attendance system using support vector machine," *Int. J. Adv. Res. Comput. Sci.*, vol. 8, no. 5, pp. 412–416, May 2017.
- [7] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf, "DeepFace: Closing the gap to human-level performance in face verification," in *Proc. IEEE CVPR*, Columbus, OH, USA, 2014, pp. 1701–1708.
- [8] F. Schroff, D. Kalenichenko, and J. Philbin, "FaceNet: A unified embedding for face recognition and clustering," in *Proc. IEEE CVPR*, Boston, MA, USA, 2015, pp. 815–823.
- [9] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "ArcFace: Additive angular margin loss for deep face recognition," in *Proc. IEEE CVPR*, Long Beach, CA, USA, 2019, pp. 4690–4699.
- [10] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao, "Joint face detection and alignment using multi-task cascaded convolutional networks," *IEEE Signal Process. Lett.*, vol. 23, no. 10, pp. 1499–1503, Oct. 2016.
- [11] G. Guo and N. Zhang, "A survey on deep learning based face recognition," *Comput. Vis. Image Underst.*, vol. 189, pp. 1–7, Dec. 2019.
- [12] R. Pandey, A. Kumar, and S. Singh, "Real-time face recognition based attendance system using FaceNet and KNN classifier," *Int. J. Eng. Res. Technol.*, vol. 9, no. 6, pp. 231–236, Jun. 2020.
- [13] Y. Kortli, M. Jridi, A. Al Falou, and M. Atri, "Face recognition systems: A survey," *Sensors*, vol. 20, no. 2, pp. 1–34, Jan. 2020.
- [14] A. Dosovitskiy et al., "An image is worth 16×16 words: Transformers for image recognition at scale," in *Proc. ICLR*, Vienna, Austria, 2021, pp. 1–22.