# Classroom Attendance System Using Deep Learning

Gaurav Khanal

2025-03-31

## Table of contents

# 1. Introduction

Accurate attendance management is a recurring challenge in educational environments, where manual roll calls can be time-consuming and prone to human error. In this paper, we propose a Deep Learning framework for face-based attendance tracking, integrating modern Convolutional Neural Network (CNN) architectures. The system detects faces within classroom images, classifies each face's gender, and matches the identity to a known database of students, thereby automating attendance logs.

Deep learning has proven to be highly effective in computer vision tasks due to its capacity to learn complex patterns directly from data [1]. Specifically, face detection, gender classification, and face recognition are well-researched fields within the deep learning community. In our approach, we harness YOLOv8 for robust and real-time face detection [2], ResNet for gender classification [3], and either FaceNet [4] or a ResNet-based classifier [5] for identity recognition. The fusion of these models results in a pipeline capable of accurately and efficiently marking attendance in a classroom setting.

Objectives of the Work

1. Face Detection (Segmentation): Accurately locate faces in an image, even under challenging real-world conditions.
2. Gender Classification: Distinguish between male and female faces for demographic insights (or future analyses).
3. Face Identification: Match each detected face to a registered student identity for attendance logging.

The remainder of this paper is organized as follows: Section 2 discusses relevant works in deep learning for face-related tasks, Section 3 presents our data gathering and annotation strategies, Section 4 outlines the methodology, Section 5 details our experimental results, Section 6 describes ethical concerns, Section 7 concludes the paper, and Section 8 provides references.

## 2. Related Works

Early work on Face Detection relied on classical feature-based techniques such as Haar Cascades introduced by Viola and Jones [6]. Although these methods are computationally efficient, they struggle with large pose variations and complex illumination. Advances in CNNs led to methods like the Multi-Task Cascaded Convolutional Network (MTCNN) [7], which improved detection accuracy by learning bounding box regression and facial landmarks jointly. Alternatively, generic object detectors like YOLO [2] can be adapted to detect faces, providing a strong balance between speed and precision.

Gender Classification initially depended on handcrafted descriptors (e.g., HOG, SVM). However, CNN architectures such as ResNet [3] and MobileNet have supplanted traditional methods by directly learning robust features from images. These models remain robust under lighting changes, partial occlusions, and natural facial variations, crucial in unconstrained classroom images.

In Face Identification, FaceNet [4] introduced an embedding-based approach, mapping faces into a compact vector space where Euclidean or cosine distances indicate identity similarity. For closed-set recognition (i.e., a fixed set of known students), a simpler approach is to fine-tune the final classification layer of ResNet [3] or VGG [5]. While embedding-based systems facilitate new class additions, classification-based systems are typically more straightforward to train when the roster remains relatively constant.

Thus, our pipeline draws upon these state-of-the-art methods to generate a robust attendance system: YOLOv8 for face localization, ResNet for gender prediction, and either FaceNet or ResNet-based classification for identity recognition.

## 3. Data

Deep learning models, particularly those targeted at face detection and recognition tasks, depend critically on representative and well-annotated datasets. In this project, we combine subsets of the Labeled Faces in the Wild (LFW) dataset [8] with synthetic classroom images to address both privacy considerations and the need for realistic variation in lighting, pose, and background. Below, we describe the data pipeline in sequential steps, highlighting how each subset was created, processed, and employed in training and testing.

### 3.1 Sampling from LFW

Our work begins by focusing on approximately 20 individuals from the LFW dataset [8]. LFW contains thousands of celebrity portrait images with associated identity labels, offering a diverse range of facial expressions and poses. We specifically chose individuals who each had sufficient sample images (at least 10–15) to allow meaningful training of recognition models and

to maintain variation across genders. Although LFW supplies valuable variety, we note that it often favors posed, front-facing portraits, which may differ from typical classroom scenes.

## 3.2 Obtaining Face Bounding Boxes with MTCNN

Once the target individuals were identified, we employed the MTCNN algorithm [7] to detect faces and produce bounding boxes for the chosen LFW images. MTCNN operates through a multi-stage convolutional process, refining candidate boxes while simultaneously identifying key facial landmarks. This step ensured that each portrait in the LFW selection would be accurately cropped around the face, providing consistent inputs for subsequent tasks such as synthetic data generation and classification.

## 3.3 Creating Synthetic Classroom Images

Since LFW images largely feature isolated portraits rather than group scenes, we addressed the lack of real classroom images by constructing synthetic classroom photographs. First, we acquired royalty-free classroom background images from the Pexels platform (under a free-use license), ensuring a variety of layouts, lighting conditions, and visual styles. Next, we pasted the MTCNN-cropped LFW faces onto these backgrounds. We applied random transformations—such as scaling, rotation, mild blur, and partial occlusion—to increase realism. Through this process, we generated between 100 and 300 synthetic classroom images, each containing multiple faces to simulate attendance scenarios.

## 3.4 Using YOLOv8 for Face Detection and Cropping

After the synthetic classroom dataset was assembled, we applied YOLOv8 [2] to detect faces within these newly created images. YOLOv8 provided bounding boxes for each face, which were verified for accuracy against the known locations of pasted faces. We then cropped the detected faces from the synthetic images and saved them to a separate folder designated as "test data," aiming to replicate how a real system would isolate faces from classroom photographs. This cropped set later proved essential in evaluating how well the trained models perform under realistic, cluttered backgrounds.

## 3.5 Gender Classification Dataset

To facilitate gender classification, we relied on the same set of ~20 individuals from LFW, for whom we had assigned gender labels. While training on LFW allowed us to capture various expressions and poses, we validated the classification model on faces cropped from the YOLOv8-detected synthetic images. This procedure tested the classifier's ability to handle

varied backgrounds and partial occlusions, even though LFW itself was predominantly portrait-oriented. To combat potential bias—given that LFW is still somewhat skewed in favor of male images—we employed either oversampling or Weighted Random Sampling, ensuring that female faces were adequately represented during model training.

## 3.6 Face Identification Dataset

For face identification, we exclusively trained our model on the original LFW data, using the same ~20 individuals. Once these networks (either embedding-based FaceNet [4] or a ResNet [3] classifier) were trained, we evaluated their performance on the cropped synthetic faces from the YOLOv8 pipeline. This approach allowed us to observe how well the models transferred from clean, mostly frontal LFW images to more cluttered scenes representative of classroom conditions, including differences in lighting, angle, and potential occlusions.

## 3.7 Final Usage

Ultimately, each facet of the pipeline—face detection, gender classification, and face identification—benefited from a slightly different subset or derivation of the original LFW images:

- Face Detection: Trained and validated on synthetic classroom images composed of LFW faces placed in open-license classroom photos.
- Gender Classification: Leveraged LFW face crops and validated on the synthetic YOLO-cropped images, with additional attention to balancing male and female examples.
- Face Identification: Trained on the original LFW portraits for each chosen identity, then evaluated on the same synthetic YOLO-cropped samples.

By carefully separating training and testing data, we aimed to provide a credible measure of real-world performance, despite relying on synthetic classrooms rather than genuine student photographs.

## 3.8 Addressing Class Imbalance

One notable challenge in assembling these datasets was the uneven distribution of male and female faces in LFW. Since male subjects outnumber female subjects, straightforward sampling could leave the network biased toward predicting male faces more frequently. To mitigate this problem, we adopted multiple strategies across different stages of the pipeline.

First, during gender classification training, we implemented oversampling of the minority (female) class so that each batch contained a proportion of female faces closer to an even split. In cases where oversampling was not ideal, we employed a Weighted Random Sampler, which assigns higher sampling probabilities to underrepresented examples. This ensures that the

training process sees each female instance more often, preventing the model from leaning excessively toward male predictions. Additionally, when preparing synthetic classroom images, we intentionally pasted female faces at slightly higher rates in some batches to approximate a more balanced scenario.

Finally, we monitored class-specific performance metrics (such as per-class F1-scores) to validate that these balancing techniques effectively reduced the discrepancy in accuracy or recall between male and female faces. These measures collectively improved the model's robustness, helping it generalize better to real-world classrooms with a diverse student body.

# 4. Methodology

Our system comprises three primary components—face detection, gender classification, and face identification—culminating in an automated attendance logging mechanism. By separating these tasks, we can individually optimize each stage for improved performance.

## 4.1 Face Detection with YOLOv8

The first step in our pipeline involves detecting faces within classroom images. We chose YOLOv8 [2] because it balances real-time speed with high detection accuracy. During training, we provided YOLOv8 with a set of synthetic classroom images where each face was labeled in the YOLO format. The training procedure typically ran for 100 epochs, during which we tracked mean Average Precision (mAP) on a validation subset to guide early stopping. By limiting training data to scenes that resemble classrooms (albeit synthetic), we aimed to minimize domain discrepancy when transferring the model to real classroom images. Once trained, YOLOv8 outputs bounding boxes around each detected face, effectively isolating regions of interest for subsequent tasks.

## 4.2 Gender Classification with ResNet

Once YOLOv8 localizes a face, we crop that region and pass it to our ResNet18 classifier [3], which predicts whether the face is male or female based on the gender labels. Before training, each image is resized to a uniform resolution of 128×128 pixels. We also apply various data augmentations (e.g., random horizontal flips and slight rotations) to increase model robustness. The ResNet architecture, pretrained on ImageNet, is modified by replacing its final classification layer with a new head outputting two logits (male, female). We employ a conventional CrossEntropyLoss, often enhanced with class weights to counter the imbalance in male–female samples. The network is trained for roughly 30–50 epochs. This procedure consistently converges to high accuracy on our validation set.

### 4.3 Face Identification

Here, a ResNet (e.g., ResNet18) is adapted so that its final fully connected layer matches the number of enrolled students [5]. The ResNet model was fine tuned just like the YOLOv8 model using the subsample LFW data.During training, we collect all cropped images of each student and learn a direct mapping from face to identity. Although closed-set classification is more straightforward, it is less flexible in handling new students or labeling unseen faces as "unknown." Hence, the choice between embedding-based or classification-based methods hinges on the deployment scenario and desired scalability.

### 4.4 Attendance Logging

Upon receiving a classroom image, the system sequentially executes the above three modules. YOLOv8 locates the faces, ResNet predicts each individual's gender, and the identification model (FaceNet or ResNet) infers their identity. A final script records each recognized individual's name, gender, and detection confidence into a CSV file or database, effectively automating attendance records. By centralizing these logs, instructors or administrators can quickly confirm who was present or absent without manual roll calls.

## 5. Results

In this section, we provide detailed, full-sentence results for each stage of our deep learning pipeline. The performance metrics and inference speeds are drawn from thorough experimentation on our synthetic classroom dataset and our carefully curated subset of the LFW dataset.

### 5.1 Face Detection Performance

Our YOLOv8 model demonstrates a mean Average Precision (mAP) of approximately 98–99% on the validation portion of our synthetic classroom images. This high mAP value indicates that the trained model successfully localizes faces under a wide array of simulated classroom conditions, including variations in lighting, pose, and partial occlusions. Further, the precision and recall metrics both exceed 95% in most of the validation samples, suggesting that the model rarely misses a face (high recall) and infrequently mislabels non-face regions (high precision). These numbers validate YOLOv8's capability to serve as a reliable first step for detecting faces in real classroom scenarios.

In terms of inference speed, YOLOv8 processes each image in 10–15 milliseconds on a standard GPU, confirming its suitability for near real-time applications. This speed facilitates instant

detection in typical classroom environments, enabling educators to capture attendance data with minimal delay.

## 5.2 Gender Classification Results

We fine-tune ResNet18 to classify each detected face as male or female. Our experiments show that the model achieves an overall validation accuracy of about 94–95%, reflecting the network's strong ability to learn facial features relevant to gender recognition. By leveraging data augmentation (random flips, rotations, color jitter), we mitigate overfitting and improve robustness to minor pose and lighting differences.

In more detail, the F1-score for the female class hovers around 92%, underscoring balanced performance between precision and recall. These outcomes are especially significant given the initial class imbalance in the dataset, which was countered through weighted sampling or oversampling. Moreover, the confusion matrix reveals that the classifier seldom mislabels male faces as female or vice versa, highlighting the effectiveness of the training strategy in distinguishing the two classes.

## 5.3 Face Identification Success

When converting the ResNet model into a multi-class identifier (one class per student), we observe a Top-1 accuracy range of 93–96% on the validation set. Most misclassifications occur under conditions such as partial occlusion or extreme angles, where faces deviate substantially from the training distribution. Despite these challenges, the results confirm that a straightforward classification approach can be highly accurate in closed-set environments, provided the roster of known students remains stable.

## 5.4 End-to-End Pipeline Evaluation

Combining the three stages—face detection (YOLOv8), gender classification (ResNet), and face identification (FaceNet or ResNet-based)—results in an automated attendance pipeline that is both accurate and efficient. In test scenarios with synthetic classroom images containing up to 10 faces, our system:

Correctly detects the majority of faces with a recall of roughly 95% or above, ensuring minimal omissions in attendance. Accurately classifies gender labels for each face at a high reliability. Identifies recognized students with an accuracy often surpassing 90%, while labeling any unregistered faces as "unknown" if using the FaceNet thresholding mechanism. Logs each recognized identity in a CSV or database within a matter of seconds, showing that computational demands are modest and deployment to real-world classrooms is feasible.

Overall, these results confirm the viability of our deep learning approach for automated classroom attendance. The high detection accuracy, efficient runtime, and strong classification/recognition performance collectively demonstrate that the proposed pipeline can significantly reduce the administrative burden on instructors while maintaining comprehensive attendance records.

# 6. Ethical Considerations

While automated attendance offers clear advantages in efficiency and accuracy, the use of facial recognition in educational settings raises important questions about privacy, consent, and fairness. Biometrics-based technologies can capture and store sensitive information about individuals, potentially subjecting them to unintended risks if systems are not designed and governed responsibly. This section highlights the core ethical dimensions institutions should consider—covering the handling of photographic data, broader privacy obligations, awareness of non-binary gender identities, the potential for increased surveillance, and the need to comply with relevant data protection laws.

## 6.1. Ethical Use of Photographs

The dataset described herein relies on face images collected in a manner potentially beyond the original subjects' awareness or consent [10]. Institutions should secure informed consent where possible, clarify data retention policies, and consider anonymizing or discarding raw images in favor of embeddings.

## 6.2. Privacy

Facial embeddings constitute biometric data, considered sensitive under regulations like GDPR [11]. Storing such data long-term can expose individuals to risks like identity theft or secondary surveillance. Implementing encryption, access control, and limited data retention is advisable, alongside providing an opt-out if feasible.

## 6.3. Non-Binary Gender Recognition

Systems that classify only "male" and "female" risk misrepresenting non-binary, gender-fluid, or transgender identities [10]. Institutions should question whether gender is necessary for attendance or, if included, consider a more inclusive labeling scheme.

## 6.4. Face Recognition and Surveillance

Face recognition can expand beyond attendance to track individual movements, posing ethical dilemmas around mass surveillance [4]. Clear policy statements must limit usage, ensuring that data collected for attendance does not creep into broader surveillance applications.

## 6.5. Regulatory Compliance

Laws like the General Data Protection Regulation (GDPR) [7] demand a legitimate basis for processing biometrics, data minimization, and a path for individuals to request erasure. Institutions in other jurisdictions (e.g., California with CCPA) should similarly review local data protection standards, possibly requiring a Data Protection Impact Assessment (DPIA).

# 7. Conclusion

We have presented a Deep Learning-based solution that automates the attendance process in classroom environments. The system leverages YOLOv8 for robust and efficient face detection, ResNet for gender classification, and either FaceNet or ResNet-based classifiers to identify individual students. Our synthetic classroom experiments demonstrate that these models can achieve high accuracy in face detection (mAP~98–99%), gender classification (~94–95% accuracy), and identity recognition (above 90%). These performance metrics suggest that the proposed pipeline can significantly streamline administrative tasks associated with attendance while maintaining reliability.

Moving forward, future improvements may involve validating the models on real-world classroom images, incorporating more diverse and comprehensive training sets, and addressing complex identity issues such as new or unregistered students. Our results and analyses emphasize the importance of anticipating ethical challenges, particularly with respect to privacy, data bias, and informed consent. By combining technical innovation with responsible deployment strategies, we can help schools realize the benefits of automated attendance without compromising fundamental rights or inclusivity.

# 8. References

[1] I. Goodfellow, Y. Bengio, and A. Courville, Deep Learning, MIT Press, 2016.

[2] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You Only Look Once: Unified, real-time object detection," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016.

[3] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016.

[4] F. Schroff, D. Kalenichenko, and J. Philbin, "FaceNet: A unified embedding for face recognition and clustering," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2015.

[5] O. M. Parkhi, A. Vedaldi, and A. Zisserman, "Deep Face Recognition," in British Machine Vision Conference (BMVC), 2015.

[6] P. Viola and M. Jones, "Rapid object detection using a boosted cascade of simple features," in Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2001.

[7] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao, "Joint face detection and alignment using multitask cascaded convolutional networks," in IEEE Signal Processing Letters, vol. 23, no. 10, 2016.

[8] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller, "Labeled Faces in the Wild: A database for studying face recognition in unconstrained environments," UMass Amherst Technical Report, 2007.

[9] J. Buolamwini and T. Gebru, "Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification," in Proceedings of the Conference on Fairness, Accountability and Transparency (FAT), 2018.

[10] C. Garvie, A. Bedoya, and J. Frankle, The Perpetual Line-Up: Unregulated Police Face Recognition in America, Georgetown Law, Center on Privacy & Technology, 2016.

[11] P. Voigt and A. Von dem Bussche, The EU General Data Protection Regulation (GDPR): A Practical Guide, Springer International Publishing, 2017.