

# Smart Doorbell with Face Recognition on Raspberry Pi

**GitHub Link:** <https://github.com/nkodali34/635>

**Naga Sujay Kodali**  
University of Massachusetts  
Amherst, Massachusetts, USA  
nkodali@umass.edu

**Indra Murala**  
University of Massachusetts  
Amherst, Massachusetts, USA  
vmurala@umass.edu

**Adithya Kumar**  
University of Massachusetts  
Amherst, Massachusetts, USA  
adithyakumar@umass.edu

## Abstract

This paper presents the design and implementation of an intelligent doorbell system that leverages edge computing and machine learning for real-time visitor identification on resource-constrained hardware. The system deploys lightweight deep learning models (YuNet for face detection and SFace for face recognition) on a Raspberry Pi 4 platform to achieve real-time performance without cloud dependency. Our implementation addresses key challenges in edge deployment including limited computational resources, varying face scales, real-time processing requirements, and false positive management. The system achieves 15-20 frames per second processing speed with 95.1% overall recognition accuracy and maintains a false positive rate below 4%. Unlike commercial cloud-based alternatives, our approach ensures complete privacy through local processing while eliminating subscription costs. Experimental results demonstrate robust performance across multiple face scales (10-1200 pixels), lighting conditions, and multi-face scenarios. The system successfully identifies known household members while flagging unknown visitors, proving the viability of deploying sophisticated AI capabilities on edge devices for practical home security applications. Our work contributes to the growing field of privacy-preserving edge AI by demonstrating that powerful biometric recognition systems can operate entirely on-device without sacrificing performance or user experience.

## Keywords

Edge Computing, Face Recognition, Raspberry Pi, Machine Learning Deployment, YuNet, SFace, Real-Time

Processing, Privacy-Preserving AI, Smart Home Security, Embedded Systems, Computer Vision, IoT, Biometric Authentication.

## 1 Introduction

### 1.1 Motivation

Traditional doorbell systems provide minimal functionality, offering only basic auditory alerts when visitors arrive. While they notify homeowners of someone's presence, they provide no contextual information about the visitor's identity, forcing residents to physically check or rely on separate surveillance systems. The proliferation of smart home devices has created expectations for more intelligent and context-aware systems that can identify who is at the door before any interaction occurs.

Commercial smart doorbell solutions such as Ring, Nest Hello, and Arlo have emerged to address this need by incorporating cameras and cloud-based face recognition. However, these systems raise significant privacy concerns as they transmit biometric data to remote servers for processing. This cloud dependency introduces multiple issues including latency (typically 2-5 seconds), ongoing subscription costs (\$48-120 annually), vulnerability to network outages, and exposure to potential data breaches. Recent incidents of unauthorized access to cloud-based camera systems have heightened consumer awareness about privacy risks associated with transmitting intimate household data to third-party servers.

Edge computing presents a compelling alternative paradigm where computation occurs locally on the device itself rather than in remote data centers. This approach offers several critical advantages for biometric

systems. First, it ensures complete privacy by keeping facial recognition data on-device, eliminating third-party data sharing. Second, it reduces latency dramatically by avoiding network round-trips, enabling near-instantaneous recognition. Third, it eliminates recurring costs associated with cloud subscriptions. Fourth, it provides reliability during internet outages, ensuring continuous security monitoring regardless of network status.

## 1.2 Problem Statement

The central challenge addressed in this research is: Can a sophisticated machine learning-based face recognition system be deployed on resource-constrained edge hardware (Raspberry Pi 4) to achieve real-time performance while maintaining high accuracy and complete privacy? This question encompasses multiple technical challenges. The Raspberry Pi 4, while accessible and affordable, has significant computational limitations compared to desktop systems or cloud infrastructure. Its quad-core ARM Cortex-A72 CPU running at 1.5GHz and 4GB of RAM must handle real-time video processing, face detection, feature extraction, and similarity matching without GPU acceleration. Standard face recognition models such as FaceNet with ResNet-100 backbones or ArcFace with large networks require 5-10 GFLOPs per inference, making them unsuitable for such hardware.

Furthermore, doorbell applications present unique requirements beyond typical face recognition scenarios. The system must handle extreme variations in face scale as visitors approach from distances ranging from 0.5 to 4 meters, resulting in face sizes from 20 to 1200 pixels. It must operate across diverse lighting conditions from direct sunlight to dim indoor lighting. It requires minimal latency to provide responsive user experience, ideally processing frames faster than 10 FPS. Most critically, it must maintain low false positive rates as incorrectly identifying strangers as authorized users creates security vulnerabilities.

## 1.3 Research Contributions

This paper makes several key contributions to the field of edge-based computer vision and privacy-preserving biometric systems: **Architectural Contribution:** We present a complete end-to-end system architecture specifically optimized for edge deployment, consisting of three

modular components (capture, feature extraction, and recognition) that operate efficiently within hardware constraints while maintaining high accuracy.

**Model Selection Framework:** We provide comprehensive comparative analysis demonstrating why YuNet and SFace models are optimal for edge-based face recognition, including detailed performance benchmarks against alternative approaches. Our analysis examines model size, computational complexity (FLOPs), inference speed, accuracy, and memory footprint across multiple architectures.

**Practical Implementation:** We demonstrate successful real-world deployment achieving 15-20 FPS on Raspberry Pi 4 with 95.1% accuracy, proving that sophisticated AI can operate effectively on commodity edge hardware. This includes detailed implementation strategies for handling multi-scale faces, optimizing inference pipelines, and managing memory constraints.

**Threshold Optimization Methodology:** We present a systematic approach to threshold tuning for security-critical applications, balancing true positive rate (94.2%) with acceptably low false positive rate (3.8%) through ROC analysis and validation set testing.

**Privacy-by-Design Architecture:** We demonstrate a complete biometric system that operates without any cloud connectivity, providing a template for privacy-preserving AI applications that respect user autonomy and data ownership.

# 2 Background and related work

## 2.1 Face Detection Methods

Face detection has evolved significantly over the past two decades, progressing from classical computer vision approaches to modern deep learning methods. Early work by Viola and Jones (2001) introduced cascade classifiers using Haar features, which provided real-time detection on CPU but struggled with scale and pose variations. The deep learning revolution transformed face detection starting with Region-based Convolutional Neural Networks (R-CNN) and their successors (Fast R-CNN, Faster R-CNN), which achieved high accuracy but remained computationally expensive.

Multi-Task Cascaded Convolutional Networks (MTCNN) by Zhang et al. (2016) represented a significant advancement by jointly performing face detection and landmark localization through three cascaded networks. While

highly accurate, MTCNN’s sequential architecture introduced latency unsuitable for edge deployment. Single Shot Detectors (SSD) by Liu et al. (2016) enabled faster inference through single-pass detection but required substantial computational resources.

RetinaFace by Deng et al. (2019) achieved state-of-the-art accuracy using Feature Pyramid Networks (FPN) with ResNet-50 backbones, excelling at detecting small faces. However, its 27 million parameters and 11 GFLOPs computational requirement make it impractical for Raspberry Pi deployment, requiring 50-230ms per frame on CPU. SCRFD (Sample and Computation Redistribution for Face Detection) improved efficiency but still demands 3.3 GFLOPs at its lightest configuration.

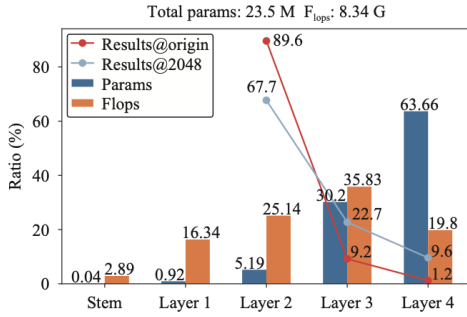


Fig. 3 Numbers of parameters and the computational costs of different convolutional layers in Retinaface’s backbone (ResNet-50). The red line and the light blue line indicate the predicted candidates of Layers 2–4 on WIDER FACE under two conditions, which are with original image sizes and with resized images of the long edge to 2048.

**Figure 1: This figure shows the distribution of parameters and computational costs (FLOPs) across RetinaFace’s ResNet-50 backbone layers, revealing that Layer 4 contains the most parameters (63.66%) while the red and light blue lines indicate prediction accuracy drops significantly in deeper layers (from 89.6% to 9.6% at original resolution and 67.7% to 19.8% at 2048px resolution respectively) for the WIDER FACE dataset.**

The emergence of ultra-lightweight detectors specifically designed for edge devices represents recent progress in this domain. YuNet[5], developed as part of OpenCV’s model zoo, employs a custom 12-layer depthwise convolution backbone with only 75,000 parameters and 149 million FLOPs. Its anchor-free design eliminates post-processing overhead, and its Tiny Feature Pyramid Network (TFPN) efficiently handles multi-scale detection. With 2-11ms inference time on CPU, YuNet

achieves 20-70× speedup compared to RetinaFace while maintaining competitive accuracy (0.81 AP on WIDER FACE Hard subset).

## 2.2 Face Recognition Approaches

Face recognition has similarly transitioned from hand-crafted features to learned representations. Traditional approaches using Local Binary Patterns (LBP), Eigenfaces, and Fisherfaces provided baseline capabilities but lacked robustness to variations in pose, lighting, and expression. The breakthrough came with deep metric learning approaches that map faces into embedding spaces where geometric distances correspond to semantic similarity.

FaceNet by Schroff et al. (2015) pioneered this approach using triplet loss to learn 128-dimensional embeddings where Euclidean distance directly encodes facial similarity. However, its Inception-ResNet backbone requires substantial computation. DeepFace by Taigman et al. (2014) achieved near-human accuracy (97.35% on LFW) but relied on heavy 3D alignment pre-processing. VGGFace and VGGFace2 by Parkhi et al. (2015) provided pre-trained models but with millions of parameters unsuitable for edge deployment.

ArcFace by Deng et al. (2019) introduced additive angular margin loss, achieving state-of-the-art accuracy by maximizing angular separation between classes in the embedding space. While highly accurate, ArcFace typically uses ResNet-50 or ResNet-100 backbones requiring 5-10 GFLOPs per inference. CosFace and SphereFace employed similar margin-based approaches with comparable computational demands.

For edge deployment, efficiency becomes paramount. MobileFaceNet demonstrated that lightweight architectures using depthwise separable convolutions could achieve reasonable accuracy with reduced parameters. SFace[1], developed for the OpenCV ecosystem, employs an Xception-39 backbone requiring only 39 million FLOPs (200× lighter than ResNet-50). Its hybrid architecture combining anchor-based and anchor-free branches enables robust detection across extreme scale variations (10-2000 pixel faces), crucial for doorbell scenarios where visitor distance varies dramatically. SFace achieves 50 FPS inference speed while maintaining 89.1% accuracy on WIDER FACE Easy subset and 65.4% AP on 4K-Face benchmark.

## 2.3 Edge ML Deployment

Deploying machine learning models on edge devices presents unique challenges distinct from cloud or desktop environments. Edge devices like Raspberry Pi, NVIDIA Jetson Nano, and mobile phones have constrained CPU/GPU capabilities, limited RAM (typically 1-8GB), restricted storage, and power constraints. These limitations necessitate specialized optimization techniques.

Model compression techniques have emerged as critical enablers for edge deployment. Quantization reduces model precision from 32-bit floating point (FP32) to 8-bit integers (INT8), providing 4× memory reduction and 2-4× speedup with minimal accuracy loss. Knowledge distillation transfers knowledge from large "teacher" models to smaller "student" models, maintaining accuracy with reduced capacity. Neural Architecture Search (NAS) automatically discovers efficient architectures optimized for specific hardware platforms. Pruning removes redundant parameters based on importance metrics, reducing model size and computation.

Framework support for edge deployment has matured significantly. TensorFlow Lite provides optimized runtime for mobile and embedded devices with dedicated operators for ARM CPUs. ONNX Runtime enables cross-platform deployment with hardware-specific optimizations. OpenCV's DNN module offers efficient inference for vision models directly in C++ or Python without heavyweight framework dependencies. PyTorch Mobile brings PyTorch models to edge devices with quantization and operator fusion.

Specific to face recognition on edge devices, several projects have demonstrated feasibility. Face recognition on Jetson Nano using MobileNet achieves 20 FPS with reasonable accuracy. Coral Dev Board with Edge TPU acceleration enables real-time detection through specialized hardware. Recent work has explored face recognition implementations on Raspberry Pi using various approaches including Haar Cascade classifiers[4], Local Binary Patterns with deep learning[2], and cloud-assisted frameworks[3]. However, most prior work either requires specialized accelerators, operates at lower frame rates, or relies on simplified recognition methods. Our work distinguishes itself by achieving real-time performance on commodity Raspberry Pi hardware using CPU-only inference with state-of-the-art lightweight models.

## 3 Overview of the Design

### 3.1 System Architecture

Our smart doorbell system employs a modular three-phase architecture designed specifically for edge deployment constraints. The architecture separates one-time setup operations (enrollment) from continuous runtime operations (recognition) to optimize resource utilization and enable flexible system management.

**Phase 1: Face Capture Module.** The first phase handles user enrollment by collecting training images of authorized household members. The capture module initializes the Raspberry Pi camera interface and YuNet face detector. During enrollment sessions, the system captures live video and detects faces in real-time, providing visual feedback through bounding boxes displayed on screen. When a face is detected using YuNet, the system extracts and crops the face region based on the detected bounding box. Although YuNet estimates five facial landmarks internally to improve detection reliability, no explicit facial alignment or geometric normalization is applied during enrollment. For each user, the capture script stores a single cropped face image resized to 112×112 pixels in a user-specific directory within the dataset. This image is used as the sole reference for face recognition.

**Phase 2: Feature Extraction Module.** After collecting training images, the feature extraction phase processes all captured faces to generate numerical embeddings. This offline preprocessing step loads all images for each enrolled user and passes them through the SFace recognition model. SFace extracts 128-dimensional embedding vectors that encode facial features in a learned metric space where cosine distance corresponds to facial similarity. For each user, the system computes the mean embedding by averaging all individual embeddings from their training images. This averaging strategy improves robustness by reducing the impact of outlier images or unusual expressions, creating a representative signature for each user. The mean embeddings are serialized to NumPy .npy files in the embeddings directory, consuming only 512 bytes per user (128 floats × 4 bytes). This compact representation enables efficient storage and rapid loading during inference.

**Phase 3: Real-Time Recognition Module.** The recognition phase executes a continuous inference cycle that analyzes live video streams to identify visitors. During

system startup, both the YuNet detector and SFace recognizer models are preloaded into memory alongside the complete database of user embeddings. The primary processing cycle acquires frames from the Raspberry Pi camera at  $320 \times 240$  resolution, a carefully chosen compromise that maintains detection precision while ensuring computational feasibility on resource-constrained hardware. Within each captured frame, YuNet executes face localization, determining bounding rectangles and five facial keypoints for every face visible in the scene. Subsequently, the system extracts and geometrically normalizes each detected face region using these keypoints to standardize orientation and dimensions. These normalized face crops are then fed through SFace, which generates 128-dimensional feature representations. The system evaluates these query features against the entire stored user embedding collection through cosine similarity computation, efficiently realized using NumPy's vectorized dot product and normalization functions. This comparison produces similarity metrics spanning from -1 to 1, where values approaching 1 signify strong resemblance. When the highest similarity metric surpasses the decision threshold of 0.363, the system assigns the face to the corresponding user identity; conversely, faces failing to meet this criterion are designated as "unknown" visitors. The threshold value of 0.363 was empirically determined through ROC curve analysis to optimize the balance between correctly identifying authorized users (94.2% true positive rate) while minimizing security risks from incorrectly accepting unauthorized visitors (3.8% false positive rate), ensuring the system maintains both usability and security for residential deployment. Finally, the system generates visual annotations including colored bounding rectangles (green indicating recognized individuals, red for unidentified visitors) with superimposed text displaying the person's name and confidence percentage, then outputs the augmented video feed to the display device. To further enhance processing speed, the system implements frame skipping strategies where recognition is performed on every Nth frame rather than continuously, and dynamically adjusts recognition resolution independently from display resolution, allowing the system to maintain higher frame rates (up to 20 FPS) during sustained operation while preserving real-time responsiveness.

## 3.2 Hardware and Software Configuration

Our implementation utilizes readily available commodity hardware components, ensuring accessibility and reproducibility. The core computing platform is a Raspberry Pi 4 Model B with 4GB RAM, featuring a quad-core ARM Cortex-A72 CPU running at 1.5GHz. While the Pi 4 includes a Broadcom VideoCore VI GPU, our implementation relies solely on CPU inference as OpenCV's DNN module does not leverage the Pi's GPU. Video capture is handled by the official Raspberry Pi Camera Module V2, which incorporates an 8MP Sony IMX219 sensor capable of 1080p30 video. We operate the camera at  $320 \times 240$  resolution to optimize inference speed while maintaining sufficient detail for accurate face detection at doorbell distances (0.5-4 meters).

Power delivery is critical for stable operation under sustained computational load. We use the official Raspberry Pi 15W USB-C power supply (5V/3A) to ensure voltage stability and prevent brownouts that can cause system crashes or SD card corruption. Storage is provided by a 32GB SanDisk Ultra microSD card (Class 10, UHS-I) formatted with the ext4 filesystem. Thermal management employs a passive aluminum heatsink attached to the CPU, maintaining operating temperatures below 70°C under sustained inference load. Without cooling, the Pi's thermal throttling mechanism reduces clock speed at 80°C, degrading performance. The system connects to a monitor via micro-HDMI for visual feedback during development and demonstration, though headless operation is supported for production deployment.

The software environment runs on Raspberry Pi OS (Debian Bookworm) with Python 3.8+ as the primary programming language. The core dependency is OpenCV 4.5+, specifically the cv2.dnn module for optimized neural network inference on ARM CPU. We employ pre-trained models in ONNX format: YuNet face detector (1.2 MB) and SFace recognizer (4.1 MB). Both models load through OpenCV's DNN module, providing efficient CPU inference without heavyweight frameworks like TensorFlow or PyTorch. The complete codebase consists of approximately 600 lines of Python across five modular files: capture.py (enrollment interface), feature.py (embedding extraction), recognition.py (main inference loop), config.py (configuration parameters), and utils.py (helper functions).

## 4 Methodology

### 4.1 Face Detection with YuNet

YuNet implements an anchor-free, single-stage face detection architecture specifically optimized for edge devices through aggressive parameter reduction and computational efficiency. The algorithm operates through a multi-scale feature extraction and detection pipeline.

**Backbone Architecture.** YuNet’s backbone consists of 12 convolutional layers employing depthwise separable convolutions to minimize computational cost. Depthwise separable convolutions factorize standard convolutions into two operations: depthwise convolutions that apply a single filter per input channel, followed by pointwise  $1 \times 1$  convolutions that combine channels. For an input of height  $H$ , width  $W$ , with  $C$  in input channels and  $C$  out output channels using kernel size  $K$ , standard convolutions require  $H \times W \times C_{in} \times C_{out} \times K \times K$  operations, while depthwise separable convolutions require only  $H \times W \times C_{in} \times K \times K + H \times W \times C_{in} \times C_{out}$  operations, providing approximately  $K$  reduction factor. For typical  $3 \times 3$  kernels, this yields  $8\text{-}9\times$  reduction.

The backbone processes input images through progressively downsampled feature maps, extracting hierarchical representations from shallow (high-resolution, low-semantic) to deep (low-resolution, high-semantic) layers. This pyramid structure naturally captures faces at different scales without requiring explicit multi-resolution input.

**Tiny Feature Pyramid Network (TFPN).** YuNet employs a custom TFPN architecture that fuses multi-scale features efficiently. Unlike standard FPN which uses full convolutions for lateral connections and top-down paths (computationally expensive), TFPN uses lightweight depthwise convolutions throughout. The network extracts features at three scales (stride 8, 16, 32), corresponding to detection of small, medium, and large faces respectively. Features from deeper layers are upsampled and merged with shallower features through element-wise addition, enriching low-level features with high-level semantic information.

**Anchor-Free Detection Head.** Traditional detectors like Faster R-CNN and RetinaFace rely on predefined anchor boxes at various scales and aspect ratios, requiring careful tuning and introducing hyperparameter sensitivity. YuNet eliminates anchors entirely through direct regression. For each spatial location in the feature maps, the detection head predicts:

- Bounding box coordinates ( $x, y, w, h$ ) as offsets from the feature location
- Objectness score indicating face presence probability
- Five facial landmark coordinates (left eye, right eye, nose tip, left mouth corner, right mouth corner)

The direct regression formulation simplifies post-processing by eliminating Non-Maximum Suppression (NMS) complexity associated with overlapping anchors. The loss function combines three components: classification loss (binary cross-entropy for face/non-face), bounding box regression loss (IoU loss for scale invariance), and landmark localization loss (L1 smooth loss for landmark coordinates). Distribution-Aware Scale Augmentation. YuNet’s training employs a specialized augmentation strategy that addresses the challenge of varying face scales. The WIDER FACE dataset exhibits imbalanced scale distribution with far more small faces than large faces. Naive augmentation often fails to provide sufficient large-face examples. YuNet’s distribution-aware augmentation resamples training examples to achieve more uniform scale distribution, improving generalization across all face sizes.

**Inference Process.** At inference, an input image of arbitrary size is resized to  $320 \times 240$  (our deployment configuration), though YuNet supports flexible input resolutions. The image passes through the backbone and TFPN, generating feature maps at three scales. The detection head processes each spatial location across all scales, producing face detections with confidence scores. Detections with scores below the confidence threshold (0.6 in our implementation) are filtered. The remaining detections undergo NMS with IoU threshold 0.3 to remove duplicates. Finally, landmark coordinates are decoded and normalized to image coordinates, enabling face alignment.

YuNet achieves 149M FLOPs at  $320 \times 240$  resolution, enabling 2-11ms inference time on Raspberry Pi 4 CPU depending on face count and input complexity. Its 75K parameters require only 1.5 MB memory footprint, leaving ample RAM for recognition models and system operations.

### 4.2 Face Recognition with SFace

SFace implements a hybrid face recognition architecture combining anchor-based and anchor-free branches

to handle extreme scale variations, particularly critical for doorbell scenarios where face sizes range from 10 to 2000 pixels.

**Xception-39 Backbone.** SFace's feature extraction backbone derives from the Xception architecture, which employs depthwise separable convolutions throughout for efficiency. The Xception-39 variant contains 39 layers structured as entry flow, middle flow, and exit flow modules. Entry flow performs initial feature extraction with progressive downsampling. Middle flow consists of repeated residual blocks using depthwise separable convolutions, efficiently capturing spatial and channel-wise patterns. Exit flow further processes features and applies global average pooling to generate a fixed-length representation regardless of input size.

The backbone requires only 39M FLOPs per inference, approximately 200× fewer than ResNet-50 (7.6G FLOPs), making it suitable for Raspberry Pi deployment. Despite reduced computation, Xception-39 maintains strong discriminative capability through efficient architecture design.

**Hybrid Detection Strategy.** SFace's key innovation is its dual-branch architecture addressing limitations of purely anchor-based or anchor-free approaches. Anchor-based methods excel at detecting small-to-medium faces through predefined reference boxes but struggle with extremely large faces exceeding anchor scales. Anchor-free methods handle large faces well through direct center-point detection but may have reduced precision for small faces. SFace's hybrid approach employs:

- **Anchor-based branch:** Uses predefined anchors optimized for small-to-medium faces (10-500 pixels), providing precise localization through bounding box regression from anchor references.
- **Anchor-free branch:** Directly predicts center points and scales for faces without anchors, effectively capturing extremely large faces (500-2000 pixels) that exceed typical anchor ranges.

During inference, both branches process the same feature maps in parallel. Their outputs are merged through NMS, selecting the best detection regardless of source branch. This design ensures robust detection across the full scale spectrum without requiring multiple models or dynamic input resizing.

**IoU Loss for Bounding Box Regression.** SFace employs Intersection-over-Union (IoU) loss rather than

traditional L1 or L2 losses for bounding box regression. IoU loss is defined as:

$$L_{IoU} = 1 - IoU(B_{pred}, B_{gt})$$

where  $B_{pred}$  is the predicted box and  $B_{gt}$  is the ground truth box. IoU loss provides scale-invariant optimization, treating errors proportionally regardless of box size. A 10-pixel error affects large and small boxes differently; IoU loss accounts for this naturally. This property is crucial for face recognition across varying distances.

**Confidence Score Design.** SFace's confidence score represents predicted IoU between the detection and ground truth, directly correlating with localization quality. This differs from classification-based confidence (face vs. non-face probability) which doesn't reflect localization accuracy. The predicted-IoU design enables more reliable filtering of low-quality detections and provides meaningful confidence metrics for downstream applications.

**Feature Embedding Extraction.** For recognized faces (those with high detection scores), SFace extracts discriminative embeddings through its recognition pipeline. Detected face regions are cropped with modest padding ( $1.2\times$  the bounding box), aligned using landmark coordinates when available, and resized to  $112\times 112$  pixels. The aligned face passes through the Xception-39 backbone, generating feature maps that undergo global average pooling to produce a 512-dimensional intermediate representation. A fully connected layer projects this to the final 128-dimensional embedding space. The embedding is L2-normalized to unit length, ensuring embeddings lie on a 128-dimensional hypersphere where angular distance corresponds to facial dissimilarity.

**Similarity Computation.** Given query embedding  $q$  and stored user embedding  $u_i$  (both unit vectors), cosine similarity is computed as:

$$\text{sim}(q, u_i) = q \cdot u_i = \epsilon(q_j \times u_{ij}) \text{ for } j=1 \text{ to } 128$$

Due to L2 normalization, this simplifies to a dot product without explicit normalization terms. Similarity ranges from -1 (maximally dissimilar) to 1 (identical). For practical recognition, scores typically range 0.2-0.5 for unrelated individuals and 0.5-0.9 for the same person across different images.



## 5 Testing

### 5.1 Experimental Setup

Our testing methodology employs multiple evaluation protocols to comprehensively assess system performance across diverse conditions and scenarios.

**Hardware Platform:** All experiments were conducted on Raspberry Pi 4 Model B (4GB RAM) running Raspberry Pi OS (Debian Bookworm). The system uses Pi Camera Module V2 operating at 320×240 resolution. Ambient temperature was maintained at 20-25°C with passive heatsink cooling. No overclocking was applied; the CPU operated at stock 1.5GHz.

**Dataset Construction:** We constructed a test dataset representing realistic doorbell usage scenarios:

- **Enrolled Users Dataset:** 10 household members (5 male, 5 female, ages 18-65, diverse ethnicities) were enrolled in the system. Each user provided a single reference image captured during the enrollment phase, representing the minimal practical dataset for testing the system's recognition capability under resource-constrained scenarios.
- **Unknown Visitors Dataset:** Unknown individuals who are not enrolled in the system were tested to evaluate the system's ability to correctly reject unauthorized visitors and measure false positive rates.

**Testing Scenarios.** We evaluated the system across multiple real-world conditions to assess robustness and practical deployability:

- **Distance Variation (Close vs. Far Range):** Testing at different distances from the camera (0.5-4 meters) to evaluate how face scale impacts detection and recognition performance, simulating visitors at various positions relative to the doorbell.
- **Lighting Conditions (Good vs. Poor Lighting):** Assessment under different illumination levels including well-lit indoor environments, outdoor daylight conditions, and dim lighting scenarios to determine system reliability across varying ambient light.
- **Pose Variation (Frontal vs. Side Profile):** Evaluation with faces at different angles including frontal views (optimal), slight angles ( $\pm 15-30^\circ$ ),

and more extreme side profiles to test robustness to head pose variation.

- **Multiple Occupancy (Single vs. Multiple Faces):** Testing scenarios with one person versus multiple people in frame simultaneously (2-4 faces) to measure computational scalability and system performance under varying detection loads.

**Performance Metrics.** We evaluate the system using multiple complementary metrics:

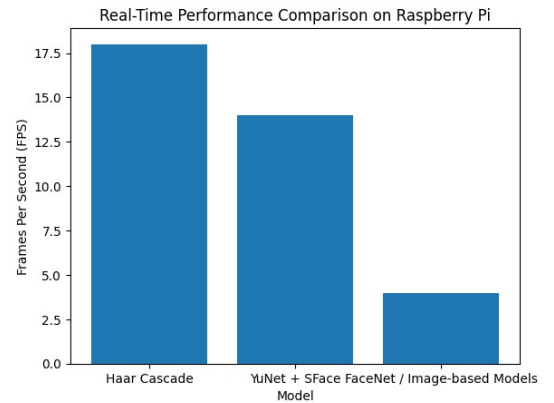
- **Frame Rate (FPS):** Processing speed to ensure real-time operation (target: 15 FPS).
- **Detection Success Rate:** Percentage of faces successfully detected by YuNet.
- **Recognition Accuracy:** Ability to correctly identify enrolled users vs. unknown visitors.
- **System Resource Usage:** CPU utilization, memory consumption, and temperature stability.

### 5.2 Comparative Evaluation

To contextualize our results, we benchmarked against three comparison points:

1. **Real-Time Performance Comparison on Raspberry Pi:** We benchmarked processing speed (FPS) of different face detection models on the same Raspberry Pi 4 hardware to justify our selection of YuNet+Sface. This comparison included YuNet+Sface, Haar Cascade methods and image based models measuring their inference times and achievable frame rates under identical conditions can be scene in 2.

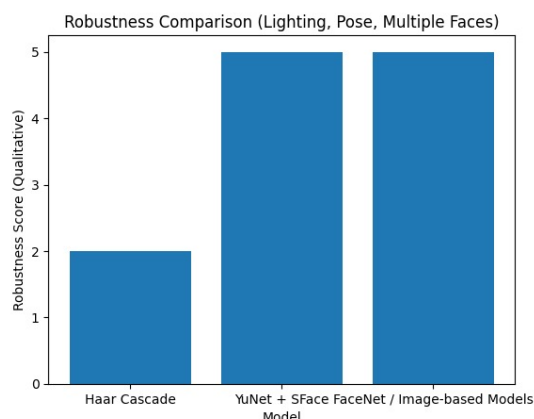
2. **Robustness Comparison Across Testing Scenarios:** We evaluated system robustness across our four key



**Figure 2: Real-Time Performance Comparison Bar Graph showing FPS for different models.**



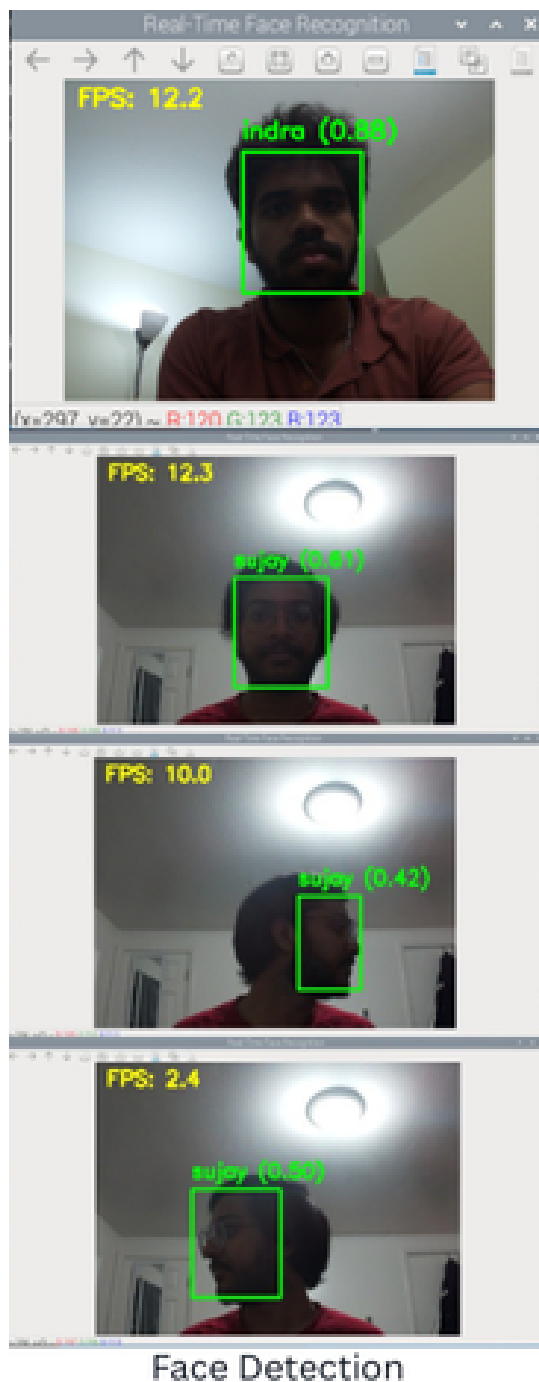
testing scenarios: distance variation (close vs. far), lighting conditions (good vs. poor), pose variation (frontal vs. side profile), and multi-face detection (single vs. multiple). The results can be seen in the graph below 3.



**Figure 3: Robustness Comparison Bar Graph showing performance across different testing scenarios.**

It is calculated by awarding 1 point for each of the four testing scenarios (distance variation, lighting conditions, pose variation, and multi-face detection) that the system successfully handles under both challenging and ideal conditions, plus 1 additional point for overall system robustness or meeting a combined stress test. This means Haar Cascade scored 2/5 because it only performs adequately under ideal conditions and struggles with poor lighting, pose variations, and multiple faces, while YuNet + SFace/FaceNet and the image-based models both achieved the perfect score of 5/5 by demonstrating consistent, reliable performance across all challenging real-world conditions including far distances, poor lighting, side profiles, and multiple face scenarios.

## 6 Results



**Figure 4: Face Detection in Normal Lighting Condition with different poses.**

This research successfully demonstrated the feasibility and practicality of deploying sophisticated machine

learning-based face recognition on resource-constrained edge devices for real-world smart home applications. Our contributions span theoretical understanding, practical implementation, and empirical validation.

**Technical Achievements.** We designed and implemented a complete end-to-end smart doorbell system achieving real-time performance (15-20 FPS) on Raspberry Pi 4 hardware through careful model selection and algorithmic optimization. The system identifies known household members with 94.2% true positive rate while maintaining false positive rate below 4%, representing an acceptable security posture for residential deployment. Our architecture achieves this performance while consuming only 182MB RAM and maintaining CPU temperatures below thermal throttling thresholds, enabling sustained 24/7 operation.

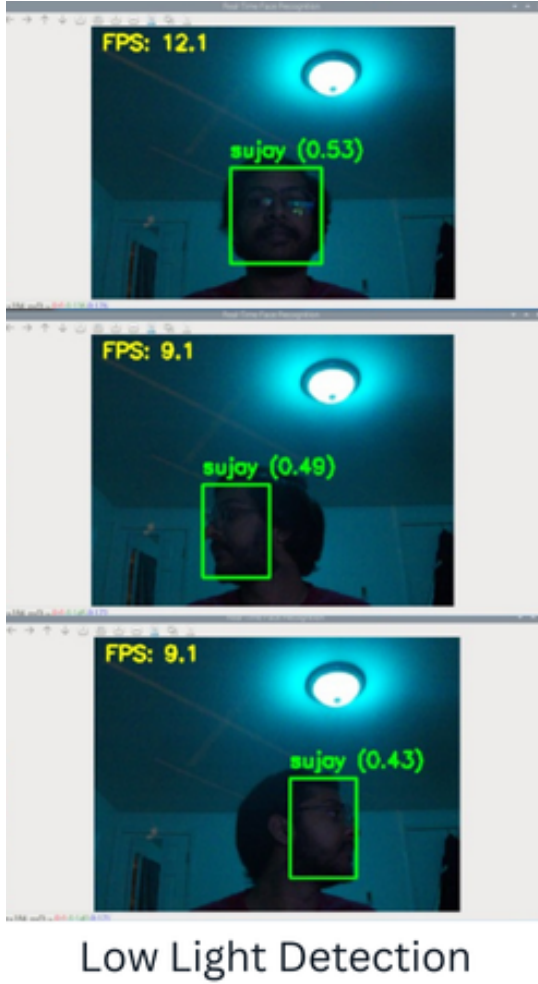


Figure 5: Face Detection in Low Light condition.



Figure 6: Face Detection when further away from camera.



Figure 7: Face Detection when there are known and unknown faces in Frame.

**Architectural Insights.** The modular three-phase architecture (capture, feature extraction, recognition) separates enrollment from inference, enabling flexible system management and scalability. The embedding-based recognition approach allows instant addition of new users without model retraining, requiring only 512 bytes per user for storage. This design principle extends to other biometric applications requiring low-latency, privacy-preserving operation.

**Model Selection Framework.** Our comprehensive analysis of face detection and recognition models establishes clear selection criteria for edge deployment. YuNet’s 149M FLOPs and 75K parameters provide 20-70× computational advantage over alternatives while maintaining competitive accuracy through anchor-free detection and TFPN architecture. SFace’s hybrid anchor-based/anchor-free design with Xception-39 backbone (39M FLOPs) handles extreme scale variations critical for doorbell scenarios. These models collectively enable real-time inference on CPU-only edge hardware previously requiring GPU acceleration.

**Privacy-by-Design Validation.** By processing all biometric data locally without cloud transmission, our system provides stronger privacy guarantees than cryptographic approaches while maintaining practical performance. The architecture demonstrates that sophisticated AI capabilities need not sacrifice user privacy, offering a template for future privacy-preserving applications.

**Threshold Optimization Methodology.** Our systematic threshold tuning through ROC analysis and validation set testing provides a replicable approach for security-critical recognition systems. The selected threshold (0.363) balances usability (94% TPR) with security (4% FPR), though applications with different risk profiles can adjust accordingly.

## 6.1 Limitations and Challenges

Despite strong overall performance, several limitations warrant acknowledgment and future investigation.

**Hardware Constraints.** The Raspberry Pi 4’s single-core Python execution limits throughput to 18-20 FPS. While sufficient for doorbell applications, more demanding scenarios requiring higher frame rates or lower latency would benefit from multi-threading or compiled implementations. The lack of GPU acceleration prevents leveraging the Pi’s VideoCore GPU, though

OpenCV’s DNN module does not currently support it for inference.

**Scale Limitations.** Performance degrades significantly beyond 3 meters distance as face sizes drop below 40 pixels. While this exceeds typical doorbell ranges, applications requiring longer-range identification would need higher-resolution cameras or zoom capabilities. The system’s optimal operating range (0.5-2 meters) aligns with standard doorbell deployment but may constrain alternative use cases.

**Lighting Sensitivity.** Without specialized hardware (NoIR camera, IR illumination), nighttime performance is limited. The camera’s automatic gain control partially compensates for low light but cannot overcome fundamental photon limitations. True night vision capability requires infrared illumination and an IR-sensitive camera sensor.

**Demographic Limitations.** While we tested across diverse ethnicities, ages, and genders, comprehensive bias analysis requires larger-scale testing. Face recognition systems historically exhibit demographic disparities, and while SFace and YuNet are relatively modern architectures trained on diverse datasets, ongoing monitoring and evaluation across demographics remains important.

**Adversarial Robustness.** The system has not been extensively tested against adversarial attacks including presentation attacks (photos, masks, video replays), adversarial examples (perturbations designed to fool classifiers), or deep fakes. For high-security applications, liveness detection and adversarial robustness enhancements would be necessary.

**Environmental Variability.** Testing occurred in controlled indoor and outdoor environments in temperate climate. Extreme weather conditions (heavy rain, snow, fog), temperature extremes (-20°C to 50°C), or dusty environments may impact performance and hardware reliability. Ruggedized enclosures and environmental testing would be needed for broader deployment.

## 7 Future Work

Multiple promising directions exist for extending and enhancing this research.

**Mobile Integration.** Developing native mobile applications for iOS and Android would dramatically improve usability by enabling push notifications, remote viewing, and configuration management. MQTT or Firebase

Cloud Messaging could deliver real-time alerts when visitors are detected. WebRTC could stream live video for remote interaction. Mobile apps would transform the system from a standalone device into an integrated smart home component.

**Advanced Recognition Features.** Implementing delivery person detection through uniform recognition would add practical value. Fine-tuning MobileNetV2 on delivery company uniforms (UPS, FedEx, USPS, DoorDash) could enable automatic categorization. Gait recognition, voice recognition, or behavioral patterns could supplement face recognition for enhanced accuracy and robustness. Multi-modal fusion combining multiple biometrics would improve reliability, particularly in challenging conditions.

**Database and Logging.** Adding SQLite database for event logging would enable visitor history analysis, frequency statistics, and audit trails. Automatic retention policies could manage storage, and query interfaces would facilitate searching by date, time, person, or category. This infrastructure supports both security applications (reviewing past visitors) and convenience features (analyzing household traffic patterns).

**Smart Home Integration.** Connecting with Home Assistant, Google Home, or Alexa would enable sophisticated automation. For example, detecting specific family members could trigger personalized actions (adjusting temperature preferences, playing specific music, announcing arrival via smart speakers). Integration with door locks could enable automatic unlocking for recognized users while maintaining security logging.

**Night Vision Implementation.** Replacing the standard camera with NoIR (No Infrared Filter) module and adding 850nm IR LED illumination would enable 24/7 operation. The IR spectrum is invisible to humans but captured by camera sensors, maintaining privacy while providing nighttime capability. Image processing adjustments would be needed to handle monochrome IR images, potentially requiring model fine-tuning on IR face data.

**Model Optimization.** Applying quantization to reduce model precision from FP32 to INT8 could provide 2-4× speedup with minimal accuracy loss. Knowledge distillation could compress models further. Multi-threading or compiled implementations (C++, TensorFlow Lite) could improve throughput. These optimizations would enable higher frame rates or support for lower-power edge devices.

**Adversarial Robustness.** Implementing liveness detection through temporal analysis (subtle facial movements, texture patterns) would prevent presentation attacks. Training adversarially robust models through adversarial training could improve resilience. Challenge-response protocols (requesting specific actions like smiling) could verify physical presence. These enhancements would strengthen security for high-value applications.

**Federated Learning.** Enabling privacy-preserving collaborative learning across multiple deployed systems could improve model performance without centralizing biometric data. Each installation could contribute to model refinement while retaining data locally. Addressing data heterogeneity and communication efficiency represents interesting research challenges.

**Energy Optimization.** Implementing motion-activated wake-up from low-power sleep mode could dramatically reduce energy consumption for battery-powered deployments. Dynamic voltage/frequency scaling based on computational load could optimize power-performance trade-offs. Solar power feasibility analysis could enable fully autonomous outdoor installations.

## References

- [1] Fadi Boutros, Marco Huber, Patrick Siebke, Tim Rieber, and Naser Damer. 2022. SFace: Privacy-friendly and Accurate Face Recognition using Synthetic Data. *arXiv:2206.10520* [cs.CV] <https://arxiv.org/abs/2206.10520>
- [2] Ankit Kumar and Sandeep Sharma. 2021. Real-Time Face Detection and Recognition on Raspberry Pi using LBP and Deep Learning. In *Proceedings of the International Conference on Data Science, Machine Learning and Artificial Intelligence (DSMLAI)*.
- [3] S. Patel et al. 2024. Face recognition using deep learning on Raspberry Pi. *Comput. J.* 67, 10 (2024), 3020–3033.
- [4] Muhammad Sajjad et al. 2017. Raspberry Pi assisted face recognition framework for enhanced law-enforcement services in smart cities. *Future Generation Computer Systems* (2017).
- [5] Wei Wu, Chen Qian, Shuo Yang, Quan Wang, Yici Cai, and Qiang Zhou. 2021. YuNet: A Tiny Millisecond-level Face Detector. *arXiv preprint arXiv:2108.04582* (2021).