# Contactless Employee Authentication Using Gait Analysis

**Submitted By:** Ravi pal

**Program:** Integrated M.Tech – Computer Science and Engineering

**Institution:** VIT Bhopal University

**Domain:**  Machine Learning / Data Science

**Project Type:**  Real-World Sensor-Based Authentication System

**Technologies Used:** Python, Scikit-learn, NumPy, Pandas, SciPy

**LLM Integration:** Explainability Layer for Access Decisions

**Submission Date:** 13.02.26

# Index

# 1. Project Overview

## 1.1 Introduction

This project presents a **contactless employee authentication system** based on **gait analysis** using smartphone sensor data. The system identifies employees by analyzing their walking patterns captured through accelerometer and gyroscope signals.

Traditional authentication systems such as ID cards, passwords, and fingerprint scanners require either physical interaction or active user participation. These methods can be inconvenient, prone to misuse, or raise hygiene and privacy concerns. In contrast, gait-based authentication works passively in the background and does not require explicit action from the user.

The primary objective of this project is to design and implement a machine learning system that can:

- Identify enrolled employees based on their walking patterns
- Reject unknown or unenrolled individuals using confidence-based logic
- Provide clear explanations for authentication decisions

---

## 1.2 Project Motivation

With the increasing adoption of smart devices and sensor technology, behavioral biometrics such as gait have gained attention as a secure and user-friendly authentication mechanism. Since gait is difficult to consciously imitate and can be captured naturally while a person walks, it provides a promising alternative to traditional security systems.

This project explores the practical implementation of a gait-based authentication system using real-world data collected from smartphone sensors. The goal is not only to achieve accurate identification but also to ensure that the system remains transparent and explainable.

---

## 1.3 Key Components of the System

The developed system consists of the following main components:

1. **Sensor Data Collection**
   Real-world accelerometer and gyroscope data are collected from enrolled individuals using a smartphone.
2. **Preprocessing and Feature Extraction**
   Raw sensor signals are filtered, segmented into fixed-size windows, and converted into statistical feature vectors.
3. **Machine Learning Model**
   A multi-class classifier is trained to identify employees based on extracted gait features.
4. **Confidence-Based Decision Logic**
   The system applies a predefined confidence threshold to grant or deny access.

5. **LLM-Based Explainability Layer**
   A Large Language Model (LLM) is integrated to generate human-readable explanations for authentication decisions.

---

# 1.4 Scope of the Project

This project focuses on:

- Building a working prototype of a gait-based authentication system
- Training and evaluating the model using real-world data
- Demonstrating access control for both enrolled and unknown users
- Documenting the methodology and explainability approach

The system is designed as a proof-of-concept implementation that demonstrates the feasibility of contactless gait-based authentication in a controlled environment.

# 2. Brief Methodology

## 2.1 Data Collection

The system is trained using **real-world accelerometer and gyroscope data** collected from enrolled individuals using a smartphone.

Each participant's data is stored separately in the following structure:

```
data/
  └── real_world/
        └── raw/
              ├── person_01/
              │     ├── acc.csv
              │     └── gyro.csv
              ├── person_02/
              └── unknown/
```

Each CSV file contains:

- time
- x-axis
- y-axis
- z-axis

The accelerometer captures linear motion, while the gyroscope captures rotational movement. These two sensors together provide sufficient information to characterize a person's walking pattern.

The `unknown` folder is used only for testing authentication of users who are not enrolled in the system.

---

## 2.2 Preprocessing

Raw sensor data cannot be directly used for model training. The following preprocessing steps were applied:

### 1. Noise Filtering

A low-pass Butterworth filter was applied to reduce high-frequency noise caused by sensor instability or environmental disturbances.

### 2. Gravity Removal

The accelerometer signal contains both body motion and gravitational acceleration.
A low-frequency filter was used to estimate the gravity component, which was then subtracted from the signal to isolate body movement.

### 3. Sliding Window Segmentation

The continuous signal was segmented into fixed-length windows:

- Window size: 128 samples
- Sampling rate: ~50 Hz
- Approximate window duration: 2.56 seconds

- 50% overlap between windows

This step ensures that each window captures a meaningful portion of a gait cycle and increases the number of training samples.

## 2.3 Feature Engineering

Instead of using raw signals directly, statistical features were extracted from each window.

For each axis of the accelerometer and gyroscope, the following features were computed:

- Mean
- Standard Deviation
- Minimum Value
- Maximum Value
- Signal Energy
- Entropy

Features from all axes (x, y, z) and both sensors were combined to form a fixed-length feature vector representing a gait segment.

This approach reduces data dimensionality while preserving essential motion characteristics.

## 2.4 Model Training

The authentication problem is formulated as a **multi-class classification task**, where each enrolled employee represents one class.

A **Random Forest classifier** was selected due to:

- Robust performance on structured feature data
- Ability to model non-linear relationships
- Stability with moderate dataset sizes
- Resistance to overfitting compared to simple decision trees

The model was trained on feature vectors extracted from enrolled employees' walking data.

## 2.5 Model Evaluation

To evaluate performance, **cross-validation** was used.
This ensures:

- Balanced representation of each employee
- Reliable estimation of classification accuracy

During authentication, the model outputs:

- Predicted employee ID

- Confidence score (maximum predicted probability)

A predefined confidence threshold is applied:

- If confidence ≥ threshold → Access Granted
- If confidence < threshold → Access Denied

This mechanism allows the system to reject unknown or unenrolled individuals.

# 3. LLM Usage Documentation

## 3.1 Purpose of LLM Integration

In this project, a Large Language Model (LLM) is integrated to improve the **interpretability and transparency** of authentication decisions.

The machine learning model produces numerical outputs such as:

- Predicted employee identity
- Confidence score
- Access granted or denied decision

While these outputs are suitable for programmatic processing, they are not easily understandable to non-technical users such as security staff or reviewers.

The LLM is used to convert these structured outputs into clear, human-readable explanations.

---

## 3.2 Where and How the LLM Is Used

The LLM operates strictly after the machine learning decision has been made.

The workflow is as follows:

1. The ML model predicts the employee identity.
2. A confidence score is computed.
3. Access is granted or denied based on a predefined threshold.
4. A structured decision summary is created.
5. This summary is passed to the LLM.
6. The LLM generates a textual explanation of the decision.

The input to the LLM includes only:

- Final decision (ACCESS GRANTED / ACCESS DENIED)
- Predicted employee (if any)
- Confidence score
- Confidence threshold

No raw sensor data or feature vectors are shared with the LLM.

---

## 3.3 Design Constraints

To ensure reliability and security, the LLM:

- Does not train or modify the machine learning model
- Does not predict employee identities
- Does not process accelerometer or gyroscope data
- Does not override access control decisions

The ML model remains the sole decision-making component of the system.

## 3.4 Benefit of LLM Integration

The LLM improves the system by:

- Providing transparent explanations for access decisions
- Making the system easier to audit and review
- Improving trust in automated authentication

This design ensures a clear separation between predictive intelligence (ML) and explanatory output (LLM).

# 4. Screenshots of Working System

This section presents visual evidence of the implemented system, including data organization, model training, authentication results, and LLM-generated explanations.

## 4.1 Real-World Data Folder Structure



*Figure 1: Real-world dataset folder structure used for employee enrollment and unknown user testing.*

## 4.2 Feature Extraction / Preprocessing Output

*Figure 2: Example of preprocessing and feature extraction applied to raw sensor data.*

## 4.3 Model Training and Evaluation Output



*Figure 3: Model training and cross-validation performance output.*

## 4.4 Authentication Result (Known User)



*Figure 4: Authentication result for an enrolled employee.*

## 4.5 Authentication Result (Unknown User) + LLM Explanation



*Figure 5: Authentication result for an unknown user with LLM-generated explanation.*