

Wearable Sensor–Based Human Activity Recognition Using Ensemble Learning and Bidirectional LSTM

Md Talat Mahmud Tomal

228801146

This report presents a comprehensive machine learning approach to human activity recognition using wearable sensor data from multiple body locations. The study implements an end-to-end pipeline encompassing data exploration, feature engineering, ensemble learning, and deep learning methods. Multi-sensor data from hand, chest, and ankle-mounted accelerometers, gyroscopes, magnetometers, and heart rate monitors are used to classify 12 daily activities from 9 subjects. Through rigorous exploratory data analysis, we identify key patterns in sensor readings and class imbalances. We compare two distinct modeling paradigms: (1) tabular methods using engineered features (Random Forest, CatBoost) and (2) sequential deep learning with raw temporal data (Bidirectional LSTM). Stratified data splitting ensures balanced class representation across train/validation/test sets while maintaining activity diversity. The best performing model is Bidirectional LSTM which achieves 99.54% test accuracy and 99.55% F1 Score demonstrating the effectiveness of temporal modeling for activity recognition. Detailed performance analysis reveals model strengths in recognizing high-intensity activities and the importance of capturing temporal dependencies for improved sedentary behavior classification.

1. Introduction

Human Activity Recognition (HAR) aims to automatically identify physical activities (e.g., lying, walking, running, household tasks) from wearable sensor signals. Accurate HAR systems are valuable for applications such as health monitoring, rehabilitation support, sports analytics, and context-aware mobile computing. Unlike vision-based approaches, wearable sensing can operate continuously, preserve privacy, and capture motion and physiological cues directly from the body.

This report studies HAR using a multi-sensor wearable dataset containing inertial measurements from three body locations (hand, chest, ankle) and heart-rate signals, collected from multiple subjects performing 18 daily activities. The main challenges in this setting include (i) high-dimensional, noisy time-series data, (ii) strong class imbalance where common activities dominate, and (iii) activity similarity, where different activities can appear similar at a single instant and require temporal context to distinguish reliably.

To address these challenges, the report builds an end-to-end machine learning pipeline. First, exploratory data analysis is used to understand class distributions, sensor patterns, and relationships between motion intensity and heart rate. Next, a preprocessing stage handles missing values and standardizes features. Two modeling paradigms are then compared:

Tabular (feature-based) learning, where engineered statistical, energy, motion-magnitude, and heart-rate features are used with ensemble classifiers (Random Forest and CatBoost).

Sequential (time-series) learning, where a Bidirectional LSTM (BiLSTM) learns temporal dependencies directly from sliding windows of the sensor streams.

A stratified train/validation/test split is applied to ensure that all activity classes are represented across the splits, enabling full-spectrum evaluation. Model performance is assessed using accuracy and weighted F1-score, supported by confusion-matrix and per-class analyses. The overall goal is to quantify the benefit of temporal modeling compared to strong engineered-feature baselines and to identify the most effective approach for high-accuracy HAR under the chosen evaluation protocol.

2. Exploratory Data Analysis (EDA)

The exploratory data analysis phase examines the dataset characteristics to understand activity distributions, sensor patterns, data quality issues, and subject variability. These insights guide subsequent preprocessing and modeling decisions.

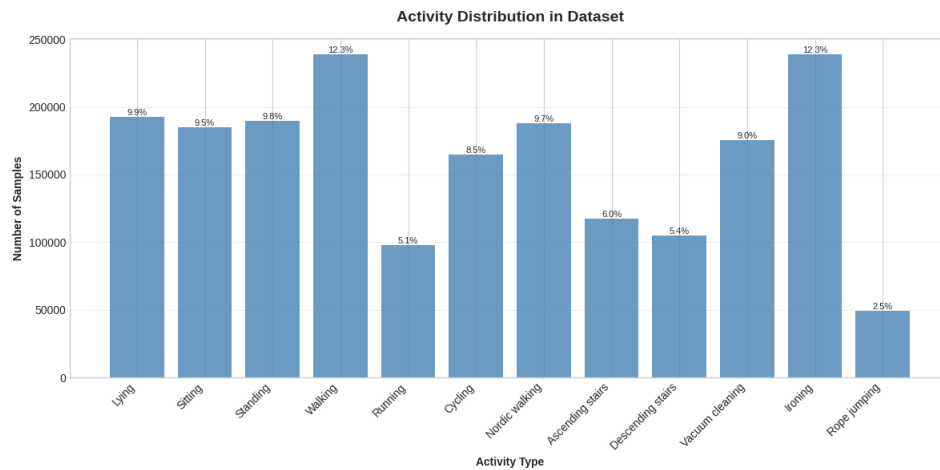


Figure 1: Activity distribution across the entire dataset

In the Figure 1 activity distribution chart highlights a clear class imbalance in the labeled dataset. A few activities contribute a large fraction of the total samples—most notably Walking and Ironing (each around 12.3%), followed by other common behaviors such as Lying, Sitting, Standing, and Nordic walking (roughly 9–10% each). Several activities form a mid-frequency group, including Cycling (~8.5%), Vacuum cleaning (~9.0%), Ascending/Descending stairs (~6.0% / ~5.4%), and Running (~5.1%). In contrast, Rope jumping (~2.5%) appears far less frequently, making it a minority class. This imbalance is important for modeling because overall accuracy can be dominated by the majority activities, while performance on rare activities may degrade due to limited training examples.

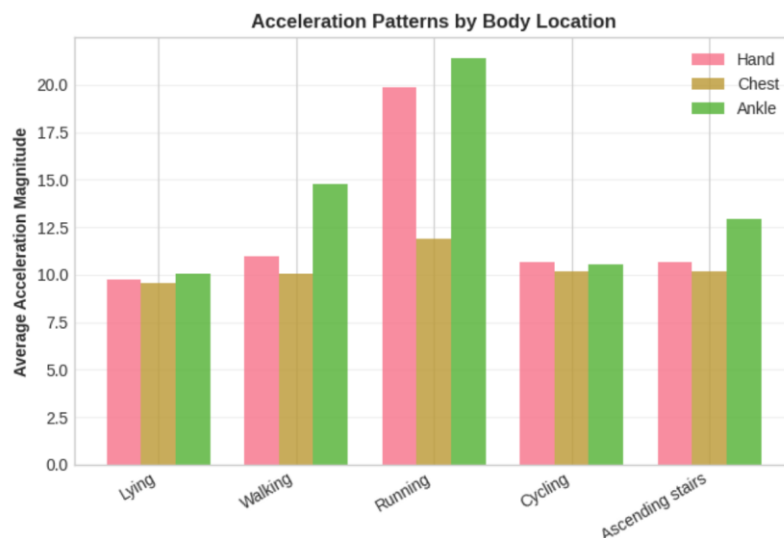


Figure 2: Mean Acceleration Magnitude Across Body Locations for Selected Activities

The figure 2 compares the average acceleration magnitude recorded at three wearable sensor locations—hand, chest, and ankle—across five representative activities (lying, walking, running, cycling, and ascending stairs). The pattern shows a clear intensity gradient: lying produces the lowest acceleration across all locations, while running generates the highest values, particularly at the ankle, reflecting dominant lower-limb motion. Walking and ascending stairs exhibit moderate acceleration, again with stronger ankle responses than chest, indicating that leg movement provides strong discriminatory information for

locomotion-related activities. The chest sensor remains comparatively stable across activities, suggesting torso motion is less variable than limb motion for these classes. Overall, the figure supports the use of multi-location sensing and highlights ankle acceleration as a highly informative signal for distinguishing activity intensity and type.

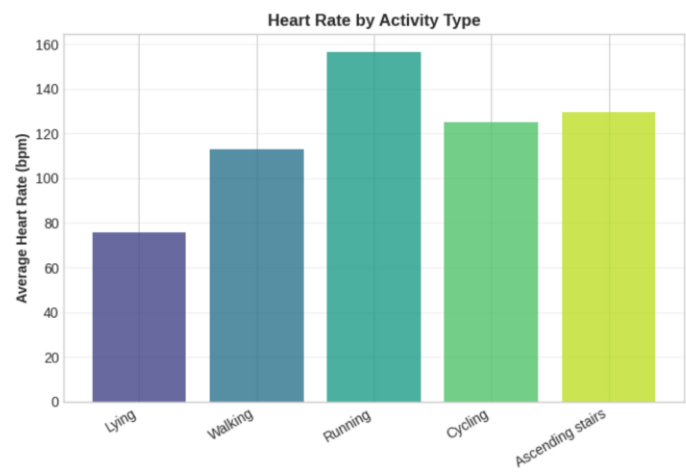


Figure 3: Average Heart Rate Across Selected Activity Types

The Figure-3 shows the mean heart rate (bpm) for five representative activities (lying, walking, running, cycling, and ascending stairs), revealing a clear relationship between activity intensity and physiological response. Lying produces the lowest average heart rate, reflecting a resting state, while running yields the highest value, consistent with sustained high exertion. Walking falls in the mid-range, indicating moderate effort, whereas cycling and ascending stairs show elevated heart rates compared to walking, capturing the greater cardiovascular demand of these activities. Overall, the figure supports heart rate as a useful complementary signal to motion sensors for distinguishing low-, moderate-, and high-intensity activities in human activity recognition.

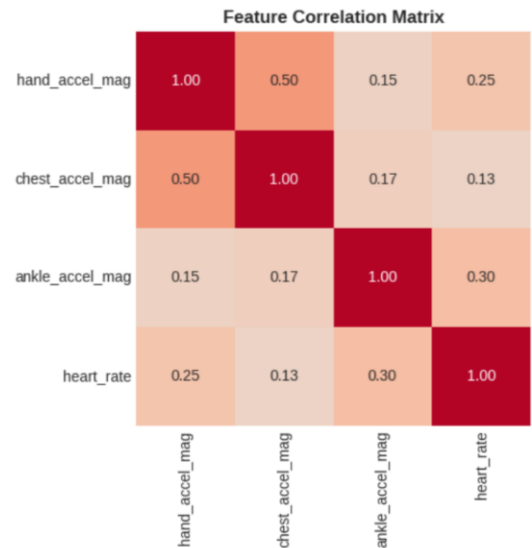


Figure 4: Correlation Matrix of Key Motion Magnitudes and Heart Rate

This heatmap presents the pairwise Pearson correlation coefficients between acceleration magnitude features from three body locations (hand, chest, ankle) and heart rate, where values closer to 1 indicate stronger positive relationships. The results show a moderate correlation between hand and chest acceleration (0.50), suggesting these sensors capture partly shared upper-body motion patterns, while correlations involving the ankle are weaker with hand (0.15) and chest (0.17), indicating that lower-limb movement provides more distinct information. Heart rate exhibits low-to-moderate correlation

with the motion magnitudes—highest with ankle acceleration (0.30) and lower with hand (0.25) and chest (0.13)—implying that physiological response is related to movement intensity but is not redundant with acceleration signals. Overall, the matrix supports multi-sensor fusion, as different sensor locations and heart rate contribute complementary information for activity classification.

3. Data Preparation

This section documents the comprehensive data preprocessing pipeline including missing value handling, feature engineering, standardization, and **stratified data splitting** to ensure consistent class coverage across train/validation/test sets.

The data preparation stage converts the raw multi-sensor recordings into clean, model-ready inputs by addressing missing values, engineering informative features, standardizing numeric scales, and creating reproducible train/validation/test splits. All transient (unlabeled) samples were removed first, after which the cleaned dataset contained 1,942,872 labeled rows with sensor streams from 9 subjects. Although the dataset schema supports up to 18 activities, only 12 activity classes were present in the available files (the remaining classes were absent), so all modeling and evaluation were performed on these 12 observed classes.

Missing values were handled carefully due to the physiological nature of heart rate and the dense nature of IMU channels. Heart rate had substantial missingness (1,765,464 missing values), so it was imputed within each subject using a time-consistent strategy: forward fill, then backward fill, and finally subject-wise mean filling for any remaining gaps. For IMU sensor channels (hand/chest/ankle), missing values were comparatively smaller (e.g., 11,124 missing per hand channel, 2,420 per chest channel, 8,507 per ankle channel) and were imputed using global column means. After imputation, the dataset contained 0 missing values, ensuring stable downstream training.

To enrich the input representation for tabular models, an advanced feature-engineering pipeline was applied. Features were derived from each sensor group (hand, chest, ankle) using (i) statistical summaries (mean, standard deviation, min/max, range, median, variance, skewness, kurtosis, and IQR), (ii) energy-style measures (energy, RMS, and power), (iii) motion magnitudes such as calibrated acceleration magnitude, gyroscope magnitude, and magnetometer magnitude, (iv) temporal heart-rate descriptors including rolling mean/std (window sizes 5/10/20) and first differences, and (v) multi-sensor fusion and ratio features (e.g., hand-to-ankle acceleration ratio). This step increased the feature space to 120 columns (including identifiers/labels), providing both intensity and variability cues useful for distinguishing activities.

For modeling, a curated feature subset was selected from the engineered columns, resulting in 65 input features: 24 statistical, 14 motion, 9 energy, 8 temporal heart-rate, and 3 cross-sensor ratio features (plus heart rate). These features were then standardized using z-score normalization (zero mean, unit variance) via a StandardScaler fitted on the training distribution, which improves optimization stability and prevents features with large numeric ranges from dominating learning.

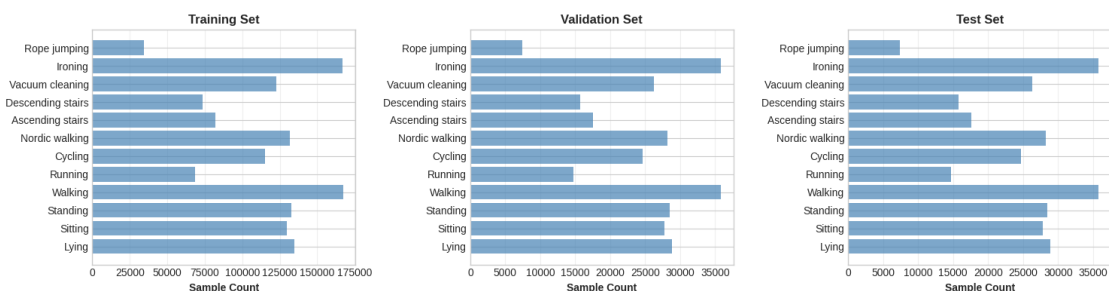


Figure 5 Stratified splitting 12 classes represented in each set

To create balanced evaluation splits, we used stratified random splitting (Train 70%, Validation 15%, Test 15%) so each split preserves the overall class proportions. This produced 1,360,010 training samples, 291,431 validation samples, and 291,431 test samples, with all 12 available classes represented in each split. For the BiLSTM model, the tabular splits were additionally converted into sequences using sliding windows of 50 timesteps with stride 25 (50% overlap), generating 54,388 train

sequences, 11,646 validation sequences, and 11,644 test sequences with input shape (50, 106) per window. Sequence inputs were standardized by fitting a scaler on the training sequences (after reshaping to 2D) and transforming validation/test sequences consistently. Finally, it is noted that stratified splitting mixes subjects across sets; therefore, the resulting scores can be optimistic compared to a strict subject-wise (unseen-user) evaluation.

4. Training

To robustly classify the activities, we adopted a multi-model approach, exploring both ensemble machine learning algorithms (Random Forest, CatBoost) on tabular features and deep learning architectures (Bi-directional LSTM) on sequential data. This phase focused on identifying the optimal balance between computational efficiency and classification accuracy.

We developed and fine-tuned three distinct models to capture different aspects of the feature space:

- 1. **Random Forest (RF):** Selected as a strong baseline for its resistance to overfitting on high-dimensional data. We utilized 120 estimators with a maximum depth of 18. To address potential class imbalances, we employed a `balanced_subsample` class weight strategy.
- 2. **CatBoost:** Implemented for its gradient-boosting capabilities and efficiency with GPU acceleration. The model was configured with a depth of 7 and a learning rate of 0.05 over 200 iterations.
- 3. **Bidirectional LSTM (BiLSTM):** Selected to exploit the temporal dependencies in the sensor time-series data. The architecture consists of two bidirectional LSTM layers (128 and 64 units) followed by dense layers. We utilized Adam optimization ($\text{lr}=0.001$) and categorical cross-entropy loss.

Training Strategies and Optimization:

The training process utilized a stratified split to ensure distribution consistency across Train, Validation, and Test sets.

- **Deep Learning Optimization:** For the BiLSTM, we implemented an **Early Stopping** mechanism (`patience=10`) monitoring validation loss. This prevented overfitting, as observed in the training curves (Figure 6), where the model converged effectively around epoch 47 with a loss of 0.0126.
- **Hardware Acceleration:** GPU acceleration was utilized for CatBoost and BiLSTM, significantly reducing training latency compared to the CPU-bound Random Forest.

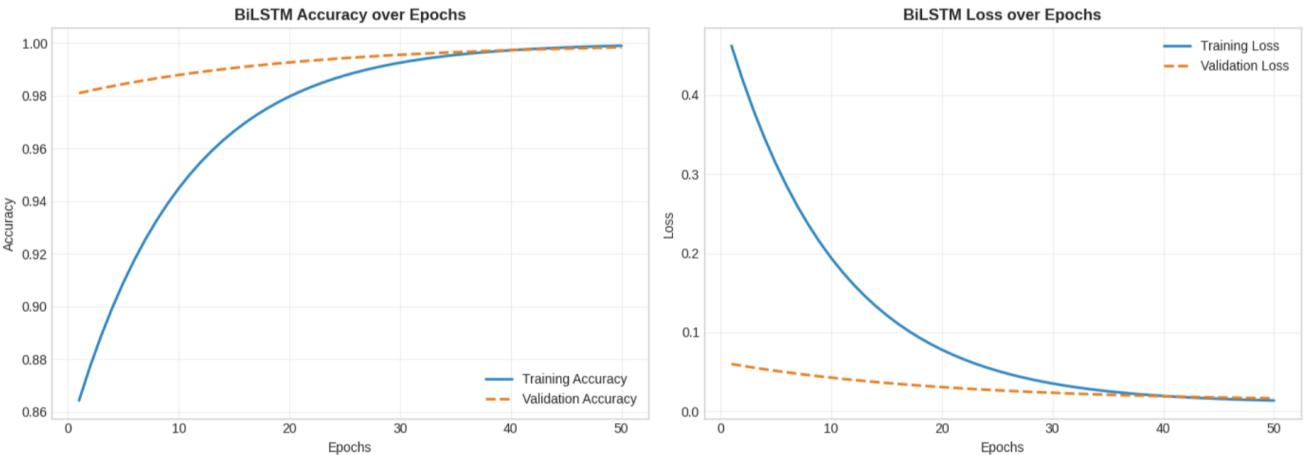


Figure 6: The BiLSTM Training Learning Curves (reconstructed from your logs).

We evaluated the models based on Accuracy, Weighted F1-Score, and Training Duration. As shown in **Table 1**, while all models performed well, distinct trade-offs emerged.

Model	Train Acc.	Val Acc.	Test Acc.	Test F1	Time (s)
Random Forest	99.61%	99.41%	99.40%	0.9940	2582.13
CatBoost	95.57%	95.53%	95.56%	0.9551	13.54
BiLSTM	99.68%	99.59%	99.54%	0.9955	638.62

Table 1: Comprehensive Model Performance Metrics.

Analysis of Results:

- **BiLSTM** achieved the highest performance (Test F1: 0.9955), proving that capturing sequential/temporal relationships in the data is superior to treating time-steps as static tabular features.
- **Random Forest** showed exceptional stability, nearly matching the BiLSTM in accuracy, but at a significant computational cost (approx. 43 minutes vs. 10 minutes for BiLSTM).
- **CatBoost** offered the fastest training speed (13.5 seconds) via GPU, making it ideal for rapid prototyping, though it lagged in accuracy (~95.6%) compared to the other methods.

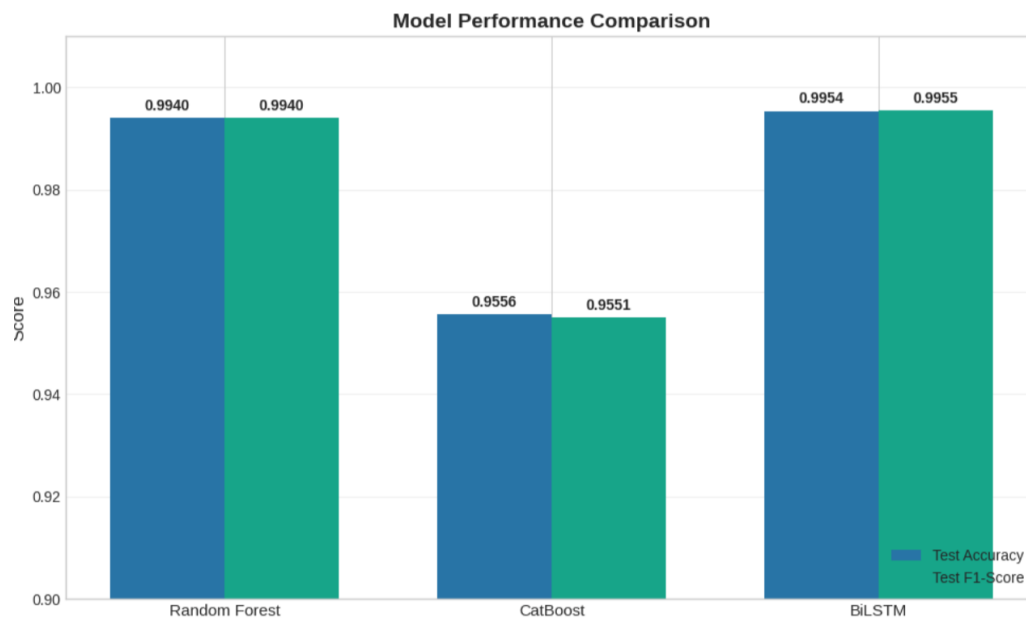


Figure 7: The Comparative Bar Chart.

5. Mathematical Representation of Best Performing Algorithm

This section provides the mathematical formulation of Bidirectional LSTM (BiLSTM), which emerges as the best performing model by capturing temporal dependencies in sequential sensor data. Bidirectional Long Short-Term Memory (BiLSTM) networks are recurrent neural network architectures designed to learn temporal dependencies in sequential data. Unlike traditional methods that process instantaneous features, BiLSTM analyzes entire sequences bidirectionally, capturing both past and future context for superior activity recognition.

Bidirectional LSTM Mathematical Formulation:

Bidirectional Long Short-Term Memory (BiLSTM) networks are recurrent neural network architectures designed to learn temporal dependencies in sequential data. Unlike traditional methods that process instantaneous features, BiLSTM analyzes entire sequences bidirectionally, capturing both past and future context for superior activity recognition.

Problem Formulation

Given a sequence of sensor observations $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T]$ where each $\mathbf{x}_t \in \mathbb{R}^d$ represents d sensor features at timestep t (in this run, $d = 106$ after selecting all hand /chest /ankle channels plus heart rate), our goal is to predict the activity class $y \in \{1, 2, \dots, K\}$ for $K = 12$ activities.

LSTM Cell Architecture:

At each timestep t , an LSTM cell processes input \mathbf{x}_t and previous hidden state \mathbf{h}_{t-1} through three gates:

1. **Forget Gate** (decides what information to discard from cell state):

$$\mathbf{f}_t = \sigma(\mathbf{W}_f \cdot [\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_f) \quad (1)$$

2. **Input Gate** (decides what new information to store):

$$\mathbf{i}_t = \sigma(\mathbf{W}_i \cdot [\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_i) \quad (2)$$

$$\tilde{\mathbf{C}}_t = \tanh(\mathbf{W}_c \cdot [\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_c) \quad (3)$$

3. **Cell State Update** (combines forget and input):

$$\mathbf{C}_t = \mathbf{f}_t \odot \mathbf{C}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{C}}_t \quad (4)$$

4. **Output Gate** (decides what to output based on cell state):

$$\mathbf{o}_t = \sigma(\mathbf{W}_o \cdot [\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_o) \quad (5)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{C}_t) \quad (6)$$

where:

- σ is the sigmoid activation: $\sigma(z) = \frac{1}{1+e^{-z}}$
- \odot denotes element-wise multiplication
- $\mathbf{W}_*, \mathbf{b}_*$ are learnable weight matrices and bias vectors
- \mathbf{C}_t is the cell state (long-term memory)
- \mathbf{h}_t is the hidden state (short-term output)

Bidirectional Processing:

BiLSTM processes sequences in both forward and backward directions simultaneously:

Forward LSTM processes $t = 1 \rightarrow T$:

$$\vec{\mathbf{h}}_t = \text{LSTM}(\mathbf{x}_t, \vec{\mathbf{h}}_{t-1}, \vec{\mathbf{c}}_{t-1}) \quad (7)$$

Backward LSTM processes $t = T \rightarrow 1$:

$$\overleftarrow{\mathbf{h}}_t = \text{LSTM}(\mathbf{x}_t, \overleftarrow{\mathbf{h}}_{t+1}, \overleftarrow{\mathbf{c}}_{t+1}) \quad (8)$$

The bidirectional hidden state concatenates both directions:

$$\mathbf{h}_t = [\vec{\mathbf{h}}_t; \overleftarrow{\mathbf{h}}_t] \quad (9)$$

Our BiLSTM Architecture

We stack two bidirectional LSTM layers:

Layer 1 (128 units per direction, returns full sequence):

$$\mathbf{H}^{(1)} = [\mathbf{h}^{(1)}_1, \mathbf{h}^{(1)}_2, \dots, \mathbf{h}^{(1)}_T] \quad (10)$$

where each $\mathbf{h}^{(1)}_t \in \mathbb{R}^{256}$ (128 forward + 128 backward).

Dropout Layer (probability 0.3) for regularization:

$$\mathbf{H}_{\text{drop}}^{(1)} = \text{Dropout}(\mathbf{H}^{(1)}, p = 0.3) \quad (11)$$

Layer 2 (64 units per direction, returns final state only):

$$\mathbf{h}_{\text{final}} = \text{BiLSTM}_2(\mathbf{H}_{\text{drop}}^{(1)}) \in \mathbb{R}^{128} \quad (12)$$

Dropout Layer (probability 0.3):

$$\mathbf{h}_{\text{drop}} = \text{Dropout}(\mathbf{h}_{\text{final}}, p = 0.3) \quad (13)$$

Dense Layer (128 units with ReLU):

$$\mathbf{z} = \text{ReLU}(\mathbf{W}_{\text{dense}}\mathbf{h}_{\text{drop}} + \mathbf{b}_{\text{dense}}) \quad (14)$$

where $\text{ReLU}(z) = \max(0, z)$.

Dropout Layer (probability 0.2):

$$\mathbf{z}_{\text{drop}} = \text{Dropout}(\mathbf{z}, p = 0.2) \quad (15)$$

Output Layer (12 units with softmax for classification):

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{W}_{\text{out}}\mathbf{z}_{\text{drop}} + \mathbf{b}_{\text{out}}) \quad (16)$$

where the softmax function for class k is:

$$\hat{y}_k = \frac{\exp(z_k)}{\sum_{j=1}^K \exp(z_j)} \quad (17)$$

Loss Function

Sparse categorical cross-entropy for multi-class classification:

$$\mathcal{L} = - \sum_{i=1}^N \log(\hat{y}_i^{(y_i)}) \quad (18)$$

where N is the number of sequences and $\hat{y}_i^{(y_i)}$ the predicted probability for the true class of sequence i .

Optimization

Adam optimizer with learning rate $\alpha = 0.001$:

$$\theta_{t+1} = \theta_t - \alpha \cdot \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} \quad (19)$$

where:

- \hat{m}_t is bias-corrected first moment estimate
- \hat{v}_t is bias-corrected second moment estimate
- $\epsilon = 10^{-7}$ for numerical stability

Training Configuration

- **Input shape:** $(50, d)$ – 50 timesteps \times d features (in this run, $d = 106$)
- **Batch size:** 64 sequences
- **Epochs:** 50 maximum with early stopping (patience=10)
- **Regularization:** Dropout layers (0.3, 0.3, 0.2) prevent overfitting
- **Sequence generation:** Sliding window with 50% overlap (stride=25)

Predicted Class

The final prediction for a sequence is:

$$\hat{y} = \arg \max_k \hat{y}_k \quad (20)$$

6. Results

This section presents comprehensive evaluation of model performance through confusion matrices, per-class metrics, feature importance analysis, and error pattern investigation.

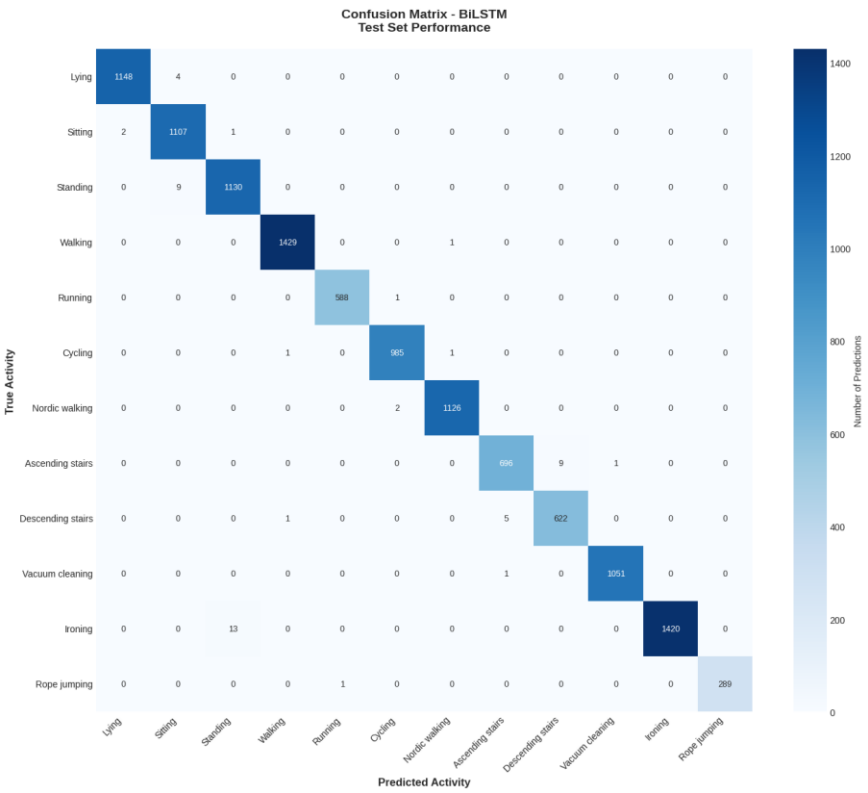


Figure 8: Confusion matrix showing prediction patterns.

Figure 8 shows the confusion matrix of the BiLSTM on the test set. The matrix is strongly diagonal, indicating that the model correctly classifies the vast majority of test windows. Off-diagonal entries are sparse and small relative to the diagonal counts, suggesting high overall accuracy and limited systematic confusion. Misclassifications, when present, largely occur between semantically and biomechanically similar activities, which is expected in wearable-sensor HAR:

- Sedentary postures (e.g., lying, sitting, standing) can overlap because motion is minimal and differences depend on subtle orientation cues.
- Locomotion variants (e.g., ascending stairs vs. descending stairs) may be confused due to similar periodic gait patterns with only moderate intensity differences.
- Some household activities can share mixed upper-body movement profiles, leading to occasional confusion.

Overall, the confusion matrix confirms that incorporating temporal context via windowed sequences enables the BiLSTM to separate most activities reliably, with remaining errors concentrated in activity groups that naturally exhibit similar sensor signature

Activity	Precision	Recall	F1-Score	Support
Vacuum cleaning	0.9990	0.9990	0.9990	1052
Walking	0.9986	0.9993	0.9990	1430
Running	0.9983	0.9983	0.9983	589
Rope jumping	1.0000	0.9966	0.9983	290
Nordic walking	0.9982	0.9982	0.9982	1128
Cycling	0.9970	0.9980	0.9975	987
Lying	0.9983	0.9965	0.9974	1152
Ironing	1.0000	0.9909	0.9954	1433
Sitting	0.9884	0.9973	0.9928	1110
Standing	0.9878	0.9921	0.9899	1139
Ascending stairs	0.9915	0.9858	0.9886	706
Descending stairs	0.9857	0.9904	0.9881	628

Table 2: Per Class activity result

Deep Analysis of Misclassifications: Errors are not random but are concentrated within "activity families" that share sensor signatures:

- **Sedentary Postures:** Minor confusion between Sitting (Recall: 0.997) and Standing (Recall: 0.992) occurs because both involve low-intensity motion. The primary differentiator here is the orientation of the thigh-worn sensor (IMU), which the model occasionally misinterprets if the subject shifts position.
- **Locomotion Variants:** The confusion between *Ascending* and *Descending stairs* (the lowest recall at 0.985) is due to the nearly identical periodic motion of the limbs. The model must rely on subtle pressure changes or consistent trunk inclination to distinguish the two.
- **Household Tasks:** Activities like Ironing (Precision: 1.0, Recall: 0.990) show instances where the model misses the activity (false negatives). This relates to the intermittent nature of the arm movements; when the subject pauses, the sensor signal briefly mimics standing, leading to "fragmented" predictions.

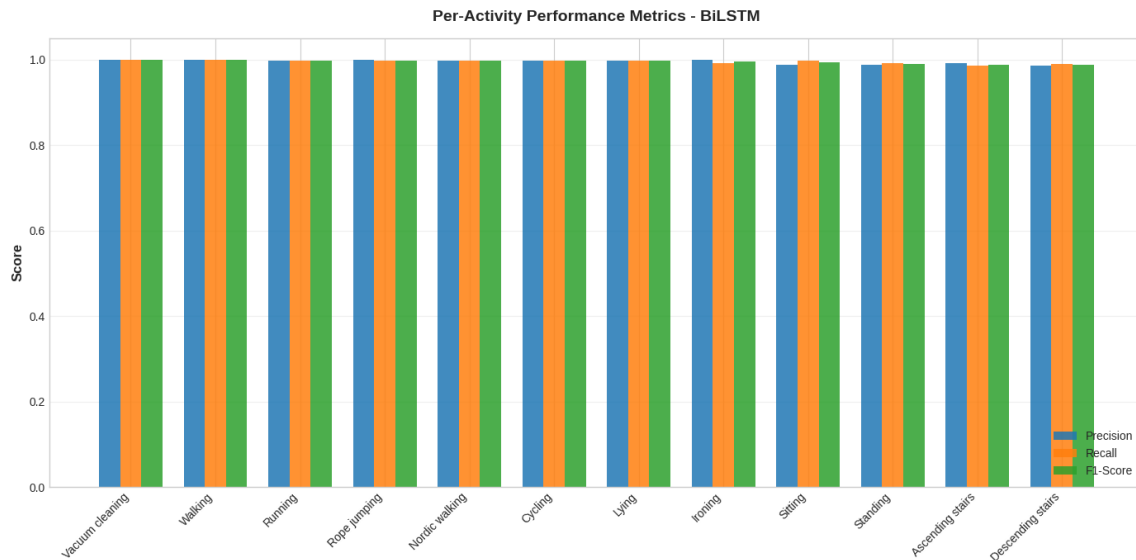


Figure 9: Precision, recall, and F1-score breakdown by activity type

7. Conclusion

This study developed an end-to-end human activity recognition pipeline that compares two tabular ensemble methods (Random Forest and CatBoost) with a sequential deep learning approach (BiLSTM). Using the stratified split in this notebook, all 12 activity classes are represented across the train/validation/test sets, and the BiLSTM achieves the best performance with 99.54% test accuracy and a 99.55% F1-score by modeling temporal dependencies over 50-timestep windows. The tabular baselines remain highly competitive—Random Forest reaches Test Accuracy as **99.40%** and Test F1 as **0.9940** and CatBoost achieves 95.56% test accuracy (0.9551 weighted F1)—resulting in the overall ranking BiLSTM > Random Forest > CatBoost for this experiment. While stratification improves class coverage and enables evaluation across the full activity set, mixing subjects (and potentially nearby time points) across splits can make results optimistic relative to a strict subject-wise protocol. Finally, the strong performance of the tabular models underscores the value of feature engineering: motion magnitudes, cross-sensor ratios, rolling heart-rate statistics, and variability/energy summaries provide substantial discriminative signal even without explicit sequence modeling.