

AUGMENTING MISSING SENSOR DATA FOR ROBUST HUMAN ACTIVITY RECOGNITION

Surya kangeyan Kandasamy Gowdaman, Sayma Akther

Department of Computer Science, San Jose State University, San Jose, CA, USA

ABSTRACT

Human Activity Recognition (HAR) using wearable sensors has become critical in healthcare, sports, and smart environments. However, model robustness is challenged by missing sensor values, class imbalance, and inter-subject variability. We propose a complete HAR pipeline addressing these issues through temporal augmentation, GAN-PCA-based imputation, and a DeepSense architecture combining convolutional and recurrent layers. Evaluations on the Opportunity dataset using K-fold, leave-one-session-out, and leave-one-subject-out crossvalidation show significant improvements over baseline methods. This study demonstrates the effectiveness of hybrid deep learning for real-world HAR.

KEYWORDS

Human Activity Recognition, Missing Sensor Data, Data Augmentation, GAN, PCA, DeepSense, Wearable Sensors, Imputation, Deep Learning.

1. INTRODUCTION

Wearable sensors-based Human Activity Recognition (HAR) has evolved into a fundamental tool enabling numerous applications in sports analytics, smart environments, healthcare, and human computer interaction[1, 2]. These systems typically depend on continuous streams of multivariate time-series data collected from embedded sensors such as accelerometers, gyroscopes, and magnetometers. The primary objective of HAR systems is to classify or forecast human activities such as walking, sitting, standing, or performing complex tasks by analyzing the patterns in sensor signals [3, 4].

Despite significant progress in deep learning techniques for HAR, a persistent limitation remains: the scarcity and imbalance of high-quality labeled sensor data. Collecting large-scale, diverse datasets for HAR is challenging and resource-intensive. It often requires long-term real-world deployments, meticulous annotation, and extensive human supervision. Natural class imbalance—where common behaviors like "standing" dominate rare but crucial actions like "falling"—further exacerbates the challenge, making it difficult to train models that generalize well across activities, users, and environments.

The conditions under which HAR data is collected contribute significantly to this problem. Datasets are typically generated through:

- Scripted scenarios, where participants follow a predefined task order in a controlled setting.
- Semi-scripted setups, offering some autonomy but within loosely structured activity guidelines.

- Wild settings, capturing naturalistic behavior without any instruction.

While scripted environments enable reproducibility, they often lack the behavioral diversity and unpredictability of real-world scenarios. Wild datasets, although more realistic, introduce additional complications like noisy labeling, irregular transitions, and unsynchronized sensors.

A more general challenge lies in the limited availability of open, high-quality HAR datasets for academic research. Many public datasets suffer from poor sensor placement, restricted activity scope, low user diversity, or missing synchronized ground-truth labels. Even when data is available, inconsistencies in labeling and temporal misalignment introduce further noise.

Thus, the core bottlenecks in HAR research include both the quantity and quality of labeled data. Models trained on small, biased datasets often overfit to a narrow range of motion signatures and fail to generalize to new subjects or deployment settings—particularly detrimental for sensitive applications like fall detection in elderly care.

This paper addresses the following key questions:

1. How can we enhance incomplete or imbalanced sensor data to better reflect real-world variability?
2. How can this enriched dataset improve the performance, robustness, and generalization of HAR models?

1.1. Proposed Solution

To overcome the challenges of traditional augmentation techniques, we propose using Generative Adversarial Networks (GANs) as a creative, data-driven solution for sensor data augmentation in Human Activity Recognition (HAR) applications.

Particularly in fields including speech synthesis, video generation, and image creation, Generative Adversarial Networks (GANs) have shown amazing efficiency in modeling complicated, high-dimensional data distributions. Applied to sensor data, the GAN architecture provides many benefits.

1. GANs may learn to generate synthetic sensor sequences that closely match real-world human motion patterns, therefore capturing complex temporal dependencies and inter-sensor correlations.
2. GANs can bring natural variations in speed, amplitude, phase, and multi-axis coupling, therefore producing more varied and rich training datasets than would simply noise addition or simple warping[5] .
3. Conditional GAN variants enable the generation of activity-specific synthetic data, hence reducing class imbalance while maintaining realism.
4. GANs' ability to acquire latent traits, such as style and intensity, which underlie human motion without requiring further annotations, makes principled augmentation possible in unsupervised or semi-supervised settings[6], [7] .

Simplified Mathematical Framing of the Problem and Solution

$$\underbrace{X_{\text{real}}}_{\text{Incomplete + Imbalanced}} \xrightarrow{\text{PCA}} \underbrace{X_{\text{pca}}}_{\text{Completed (Low-Rank)}} \xrightarrow{\text{GAN}} \underbrace{X_{\text{synthetic}}}_{\text{Balanced + Enhanced}} \Rightarrow \underbrace{X_{\text{final}} = X_{\text{real}} \cup X_{\text{synthetic}}}_{\text{Augmented Training Set}}$$

Explanation:

- X_{real} : Original multimodal sensor data, often affected by missing values due to signal dropouts and exhibiting class imbalance across activity labels.
- PCA: Principal Component Analysis is applied to estimate missing values by leveraging the low-rank nature of human activity patterns, resulting in a completed dataset suitable for downstream generation.
- X_{pca} : The output of the PCA-based imputation step. This is a denoised and completed version of X_{real} with restored missing segments based on dominant temporal and intersensor correlations.
- GAN: A conditional Generative Adversarial Network trained to generate synthetic sensor sequences that mimic the statistical structure of minority classes and underrepresented motion patterns [8], [9], [10].
- $X_{\text{synthetic}}$: Synthetic samples that both rebalance the class distribution and enhance generalization by modeling realistic signal trajectories.
- X_{final} : Final dataset formed by combining the original data X_{real} with the generated synthetic data $X_{\text{synthetic}}$, improving the performance of downstream HAR models.

2. RELATED WORK

Recent years have seen extensive research on Human Activity Recognition (HAR) using wearable sensors, with developments in deep learning, generative modeling, and self-supervised learning enabling significant progress. Conventional machine learning techniques, including Support Vector Machines (SVMs), Decision Trees, and k-Nearest Neighbors (k-NN), often relied on manually crafted features and domain-specific heuristics and were among the first benchmarks for HAR problems [5]. However, these approaches struggled with high-dimensional, noisy, and incomplete time-series data.

Deep learning models, particularly convolutional and recurrent neural networks, have emerged as the benchmark for HAR to overcome these limitations. Temporal Convolutional Networks (TCNs) have demonstrated effectiveness in modeling sequential activity patterns [7], while Ignatov proposed a CNN-based architecture for real-time HAR using accelerometer signals [6]. Despite their success, these models are often brittle under real-world conditions marked by sensor failures and class imbalance, as they typically assume fully observed inputs.

Generative approaches have been explored to improve data robustness and impute missing values. Ma et al. introduced ActivityGAN, which restores temporal dynamics through adversarial generation [9], and Zhang et al. developed SensorGAN to synthesize fine-grained sensor streams for data augmentation [8]. Although these methods showed promising results, limitations such as instability during training or low-fidelity outputs restrict their reliability. Our model builds on these works by integrating GAN-based imputation with PCA-based compression, allowing for more structured and realistic reconstruction of missing segments.

Adversarial learning has also been investigated for domain adaptation and invariant feature learning in HAR. Bai et al. proposed an adversarial multi-view network to improve generalization across sensor positions and perspectives [10]. Similarly, XHAR [11] and the work by Wang et al. [12] address domain shift using adversarial transfer learning. While these methods target generalization across contexts, they are less effective in resolving challenges related to incomplete sensor data, which is the primary focus of our work.

Recent advances in self-supervised learning have further influenced HAR research. Haresamudram et al. applied contrastive predictive coding to improve performance in downstream classification tasks [13], and Song et al. used contrastive pretraining strategies for ambient sensor data [14]. Su et al. explored disentangled behavior patterns to enhance model transferability and interpretability [15]. These approaches help reduce the reliance on large labeled datasets but often degrade when faced with high missing-data frequencies.

Moreover, recent methods have proposed advanced data synthesis and augmentation strategies. Li et al. introduced a statistical diffusion model to generate high-fidelity activity signals [16], while Um et al. proposed AutoAugHar, an automated augmentation strategy driven by reinforcement learning [17]. Our pipeline complements these efforts by unifying temporal augmentation and imputation within a GAN-PCA hybrid module, tightly integrated into the downstream DeepSense classification framework.

3. DATASETS AND PREPROCESSING

3.1. Opportunity Dataset

The Opportunity Activity Recognition Dataset is a benchmark in HAR research. Developed under the European Commission’s Opportunity initiative, it captures multivariate time-series data from wearable and ambient sensors as participants perform daily activities in a sensor-rich smart environment. Sensors include accelerometers, gyroscopes, magnetometers, IMUs, and ambient/objectbased sensors. Activities recorded span locomotion (e.g., walking, standing) and object interactions (e.g., drinking, opening drawers).

The dataset consists of:

- 4 subjects (S1–S4) performing multiple Activity of Daily Living (ADL) sessions.
- A sampling rate of 30 Hz, producing high-resolution time-series data.
- Over 1 million labeled time steps, covering 200+ sensor channels.

Locomotion labels used in our work include: Stand, Walk, Sit, and Lie. The dataset features real-world noise, sensor dropout, and class imbalance, making it ideal for evaluating generative augmentation strategies.

3.2. Preprocessing Pipeline

Given the dataset’s complexity and noise, the following preprocessing steps are employed:

- Column Name Mapping and Structural Validation: We ensure consistency across sessions by validating column names and filling missing columns with placeholders.
- Missing Value Masking: A binary mask $M \in \{0,1\}^{T \times D}$ is generated to track missing values. This mask is used in both GAN training and PCA filtering.
- Initial Mean Imputation: Missing entries are first filled using feature-wise mean values to preserve continuity. This enables normalization and segmentation without errors.
- Z-score Normalization: Standard score normalization is applied channel-wise to ensure zero mean and unit variance, improving convergence and stability of deep models.
- Sliding Window Segmentation: Sensor data is divided into 1-second overlapping windows (30 timesteps with 50% overlap). Each window is labeled using majority voting on included time steps.

3.3. Secondary Dataset

To assess generalization, the HHAR (Heterogeneity HAR) dataset is used in auxiliary evaluation. It contains data from 9 users using multiple smartwatches and smartphones during daily activities. Its variability in device type, sensor quality, and subject behavior provides a valuable testbed for assessing model robustness beyond the Opportunity dataset. Despite not being the primary training source, HHAR results offer strong support for the cross-domain effectiveness of our GAN-PCAbased augmentation pipeline.

4. SLIDING WINDOW SEGMENTATION

The temporal aspect of the data is fundamental in Human Activity Recognition (HAR) utilizing wearable sensors. Walking, sitting, or transitioning between postures are not discrete activities but rather temporal processes that extend over multiple time intervals. HAR relies on the understanding of patterns within short-term sequences of motion sensor data, in contrast to stationary classification jobs where each sample is independent [16].

The continuous stream of raw sensor data is divided into overlapping fixed-length chunks with a sliding window method to effectively analyze these temporal patterns. Each resultant window is treated as an individual instance for model training and inference.

4.1. Sliding Window Configuration

We define the following parameters in our segmentation pipeline:

- Window size W : 30 time steps
- Stride S : 15 time steps (50% overlap)

At a sampling rate of 30 Hz, characteristic of the Opportunity dataset, a window size of 30 encompasses 1 second of sensor data. Adjacent windows with a 50% overlap share half of their contents, facilitating seamless transitions and enhancing sample density.

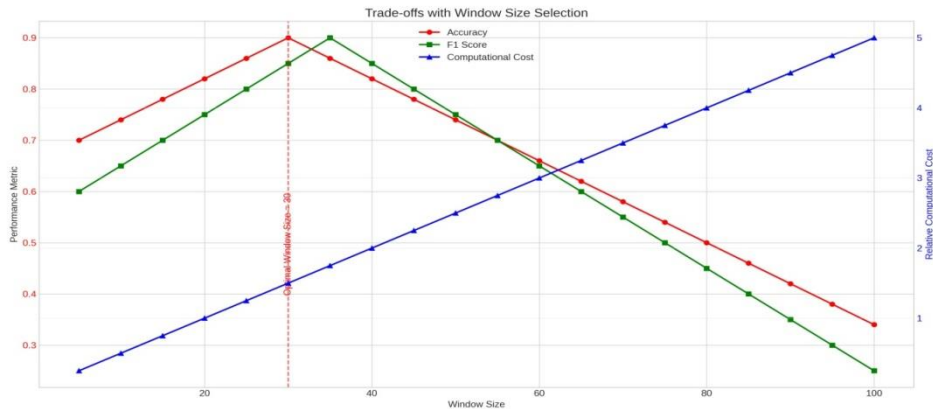


Figure 1: Trade-off analysis between window size and model performance in Human Activity Recognition.

As window size increases beyond the optimal point (30), both accuracy and F1 score decline, while computational cost continues to rise linearly. This highlights the importance of selecting an appropriate window size to balance recognition performance and efficiency.

4.2. Why window size of 30

Let the preprocessed time-series input be:

$$X = \{x_1, x_2, \dots, x_T\}, \quad x_t \in R^d$$

where:

- T is the total number of timesteps,
- d is the number of sensor features per timestep.

Using a sliding window of size W and stride S , we generate a sequence of windows:

$$\mathcal{W}_k = \{x_{kS+1}, x_{kS+2}, \dots, x_{kS+W}\}, \quad k \in \left\{0, 1, \dots, \left\lfloor \frac{T-W}{S} \right\rfloor\right\}$$

Each window $\mathcal{W}_k \in R^{W \times d}$ is treated as a single training instance.

The choice of window size directly influences semantic representation of activity patterns, computational efficiency, and temporal resolution.

Justification grounded on empirical data:

For a walk or step, a 1-second window is enough to maintain a high temporal resolution while nevertheless catching a full walking cycle.

Based on empirical data, 30 is a suitable balance. past studies utilizing the Opportunity dataset (e.g., Hammerla et al., 2016) have also employed windows of 24–32 time steps.

Theoretical Consideration:

Shorter windows (less than 20 time steps) may not adequately record complete activity units, and label noise may rise.

Longer windows (>60 time steps) produce label ambiguity; these also limit temporal granularity and can cover several activities.

For example: From a 3-minute sequence, a window of 30 with a stride of 15 produces around 6,600 windows. Without being unduly repetitious, this generates adequate training data.

5. METHODOLOGY

This chapter describes the approach applied to create an extensible and strong Human Activity Recognition (HAR) pipeline based on wearable sensor data.

5.1. Model Architecture and Training

DeepSense is a deep learning framework designed especially for the time-series data analysis across several sensors. Originally formulated to address problems in mobile sensing, including human activity detection, location tracking, and energy expenditure estimate, Essential for activities identification, the hierarchical architecture of the model is proficient in extracting local patterns from raw sensor data as well as global dependencies across time.

Unlike superficial or just recurrent models, DeepSense successfully captures spatial correlations among sensor channels and temporal dependencies over time by combining several convolutional layers with gated recurrent units (GRUs). Furthermore, the model is a suitable alternative for real world sensor streams when manual feature extraction is either insufficient or impossible since it can end-to-end learn from raw input tensors to ultimate activity class predictions [18].

The architecture runs in two main phases. Sensor level feature extraction uses a sequence of convolutional layers. On each sensor modality or their fused combinations, these layers operate independently. This phase codes localized patterns in the sensor data, including step periodicity during ambulation or acceleration spikes across transitions [11], [12]. To encode sequential links, the temporal modeling phase uses GRUs or bidirectional GRUs so enabling the model to understand activity patterns over long periods. The model uses a fully linked classification head to generate unique activity labels from the aggregated hidden states.

5.2. Standard DeepSense Model (with Mean Imputation)

The standard DeepSense architecture implemented in our pipeline follows the original formulation proposed for mobile sensing tasks, including human activity recognition (HAR). It is designed to process multivariate time-series data through a combination of convolutional and recurrent layers, enabling the model to capture both spatial correlations among sensor modalities and temporal dependencies across time.

Each input sample to the model is a fixed-length, normalized segment of multivariate sensor data extracted using a sliding window approach. The input is denoted as:

$$X \in R^{w \times D}$$

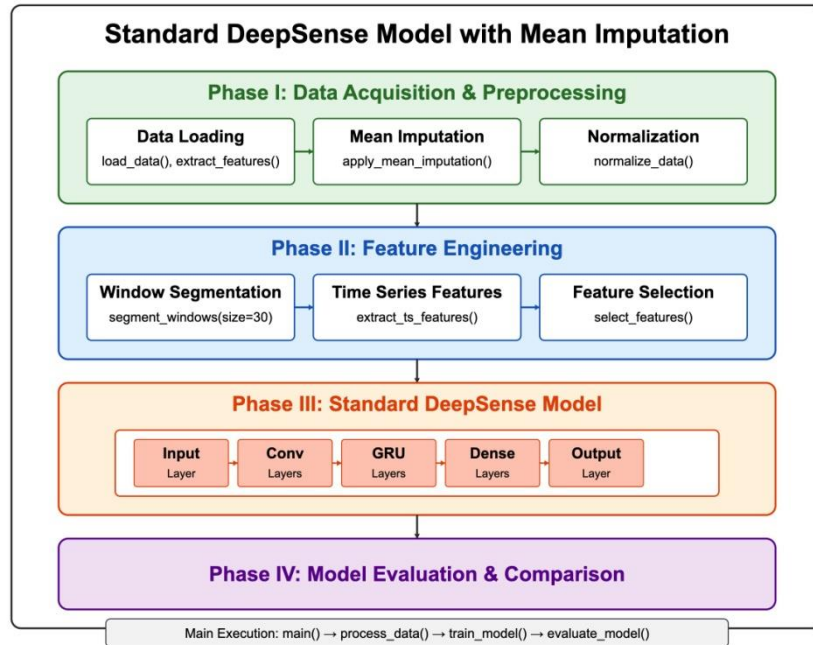


Figure 2: Workflow of the standard DeepSense model with mean imputation, from preprocessing to activity classification.

where w is the window size (number of time steps), and D is the number of sensor channels (e.g., accelerometers, gyroscopes, magnetometers, etc.).

Convolutional Feature Extraction: The first stage of the model consists of a series of onedimensional convolutional layers. These layers are responsible for extracting local temporal features within each sensor channel [15]. A typical convolutional layer applies K filters of width k across the input sequence, producing an intermediate representation:

$$H_{\text{conv}}^{(l)} = \sigma(W^{(l)} * H^{(l-1)} + b^{(l)})$$

where $H^{(l-1)}$ is the input to layer l , $W^{(l)}$ and $b^{(l)}$ are the learnable weights and biases, $*$ denotes the 1D convolution operation, and σ is a non-linear activation function such as ReLU.

These convolutional layers help the model learn patterns such as repeated footstep impacts, directional acceleration changes, or short bursts of activity that are characteristic of specific human movements.

Recurrent Temporal Modeling: To capture long-range temporal dependencies, one or more Gated Recurrent Unit (GRU) layers then receives the feature maps produced by the convolutional layers [17]. Sequential data is handled efficiently and effectively by GRUs, which are hence selected. Given a sequence h_t of extracted features at each time step t , the GRU updates its hidden state using the following operations:

$$z_t = \sigma(W_z h_t + U_z h_{t-1} + b_z), \quad r_t = \sigma(W_r h_t + U_r h_{t-1} + b_r)$$

$$\tilde{h}_t = \tanh(W_h h_t + U_h(r_t \odot h_{t-1}) + b_h)$$

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t$$

where z_t and r_t are the update and reset gates respectively, \odot denotes element-wise multiplication, and W_* , U_* , b_* are trainable parameters.

These layers enable the model to learn temporal patterns such as activity initiation, transitions, and periodicity factors crucial for recognizing human actions over time.

Fully Connected Classification Head: The final output of the GRU layers is aggregated and passed through one or more fully connected layers followed by a softmax layer:

$$\hat{y} = \text{softmax}(W_{\text{fc}} h + b_{\text{fc}})$$

where h is the output from the GRU, W_{fc} and b_{fc} are the weights and biases of the dense layer, and $\hat{y} \in R^C$ is the predicted probability distribution over C activity classes.

Training Configuration: The model is trained using the categorical cross-entropy loss:

$$\mathcal{L}_{\text{CE}} = - \sum_{c=1}^C y_c \log(\hat{y}_c)$$

where y_c is the ground truth label and \hat{y}_c is the predicted probability for class c . Optimization is performed using the Adam optimizer with learning rate schedules and dropout regularization to prevent overfitting.

Data preprocessed with mean imputation uses this baseline architecture whereby the global mean of each individual channel replaces missing sensor values. While under clean data assumptions this naive handling of missing information presents difficulties; these are addressed in next sections using generative augmentation and imputation techniques.

5.3. Improved DeepSense Model with GAN-PCA-Based Imputation

For time-series classification tasks, the standard DeepSense model performs well. Its robustness may be limited, though, by its reliance on crude imputation methods like mean imputation, particularly when there is significant sensor dropout and class imbalance. In order to get around these limitations, we add a data augmentation pipeline that uses a Generative Adversarial Network (GAN) for class-balanced synthetic data generation and Principal Component Analysis (PCA) for low-rank feature completion [9], [8]. This improved model is more sensitive to minority class patterns and has better generalizability.

5.3.1. PCA-Based Dimensionality Reduction

Two basic issues in multivariate time-series sensor data: high dimensionality and the presence of noisy, redundant, or weakly informative signals are addressed by adopted as a basic preprocessing step in the enhanced DeepSense pipeline Principal Component Analysis (PCA). The Opportunity dataset produces interesting feature sparsity and possible overfitting during model training using almost 240 sensor channels per time step. PCA compresses the input representation while maintaining most predictive of the activity patterns variance[19], [10].

Mathematical Formulation: Given a data matrix $\mathbf{X} \in R^{T \times D}$, where T is the number of time steps and D is the number of sensor features, PCA decomposes the covariance structure of \mathbf{X} via:

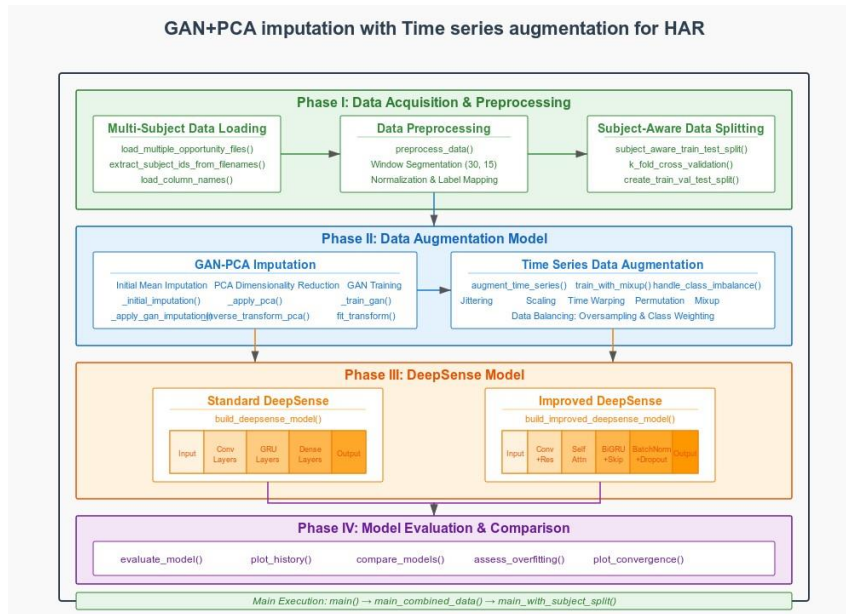


Figure 3: Proposed model architecture

$$\begin{aligned}\mathbf{X}_{\text{centered}} &= \mathbf{X} - \bar{\mathbf{X}} \\ \mathbf{C} &= \frac{1}{T-1} \mathbf{X}_{\text{centered}}^T \mathbf{X}_{\text{centered}} \quad (\text{Covariance matrix}) \\ \mathbf{C}\mathbf{V} &= \mathbf{V}\mathbf{\Lambda} \quad (\text{Eigen decomposition})\end{aligned}$$

Here, \mathbf{V} contains the eigenvectors (principal components), and $\mathbf{\Lambda}$ is a diagonal matrix of corresponding eigenvalues.

We retain the top- k eigenvectors corresponding to the k largest eigenvalues to project the data into a lower-dimensional subspace:

$$\mathbf{XPCA} = \mathbf{X}_{\text{centered}} \cdot \mathbf{V}_{[:,1:k]}$$

In our implementation, k was empirically selected such that at least 95% of the cumulative variance was preserved. This typically reduced the feature space from $D = 242$ to approximately 30–40 components, substantially reducing the learning complexity for the downstream DeepSense network.

Project-Centric Adaptation: In our pipeline, PCA was applied post-imputation (when applicable) but before windowing to ensure that every temporal window segment benefits from compact, variance-preserving features. The transformation matrix derived from the training set was reused consistently during test-time evaluation and GAN-based data generation to maintain alignment between synthetic and real data spaces[15], [13].

5.3.2. GAN-based Missing Data Imputation

To address the issue of non-random missing values in multivariate time-series sensor signals, we employ a Generative Adversarial Network (GAN) specifically tailored for imputation. Unlike traditional statistical methods (e.g., mean or median imputation), GANs are capable of learning complex data distributions and inferring missing segments by generating samples that are indistinguishable from the true distribution[16]. This enables more realistic and structurally coherent recovery of missing sensor data.

5.3.3. Training Objectives

The GAN is trained in a minimax fashion with a composite loss function defined as:

$$\mathcal{L}_{\text{GAN}} = \min_G \max_D E_{X \sim p_{\text{data}}} [\log D(X)] + E_{\hat{X} \sim G} [\log(1 - D(\hat{X}))] \quad (1)$$

To improve training stability and ensure that observed values are preserved, we introduce a reconstruction loss:

$$\mathbf{L}_{\text{recon}} = E[\|M \odot (\hat{X} - X)\|_2^2] \quad (2)$$

The total loss function for the generator becomes:

$$\mathbf{L}_{\text{total}} = \mathbf{L}_{\text{GAN}} + \lambda \mathbf{L}_{\text{recon}} \quad (3)$$

where λ is a regularization hyperparameter (empirically set to $\lambda = 10$) that balances generative realism and fidelity to observed data.

5.3.4. Implementation Details in Our Project

The architecture and configuration of the GAN were carefully adapted to suit the temporal and multivariate characteristics of the Opportunity dataset. The generator-discriminator design was selected not only for its ability to synthesize realistic sensor sequences but also for its computational tractability when integrated into a large-scale HAR pipeline.

- **Generator (G):** The generator consists of two stacked *bi-directional* Gated Recurrent Unit (GRU) layers, each with 128 hidden units. Bi-directionality was used to ensure that both past and future context are leveraged when inferring missing values in each time step. The outputs of the final GRU layer are passed through a time-distributed fully connected (dense) layer to map the sequence back into the original feature space $R^{T \times D}$. This architecture enables temporal smoothing and dynamic modeling of dependencies across sensor channels and time.
- **Discriminator (D):** The discriminator is a one-dimensional convolutional neural network has three consecutive layers. Each layer consists of 64 filters, has ReLU activation functions, and utilizes a kernel size of five pixels. Dropout regularization, applied with a rate of 0.3 between layers, aids in mitigating overfitting and enhancing generalization. The scalar sigmoid activation accurately represents the input sequence's reality.
- **Masking Strategy:** During training, missing entries were identified via a binary mask $M \in \{0,1\}^{T \times D}$, where $M_{t,d} = 0$ denotes a missing value. In addition to the actual missing data present in Opportunity dataset, we introduce the synthetic dropout during training by randomly masking 10% of observed entries. This augmentation strategy forces the generator to generalize better by learning plausible imputations under diverse masking patterns.
- **Noise Injection (z):** A latent noise vector $z \in R^{T \times D}$ is sampled from a standard multivariate Gaussian distribution $N(0, I)$. To preserve temporal alignment, z is sampled per time step and concatenated with the masked input along the feature dimension. This encourages the generator to learn a mapping from the noisy, incomplete input to a coherent full-sequence output.
- **Optimization:** The GAN was trained using the Adam optimizer with learning rate $\eta = 1 \times 10^{-4}$ and $\beta_1 = 0.5$, $\beta_2 = 0.9$ to balance stability and convergence speed. Gradient clipping (max norm = 5.0) was applied to prevent exploding gradients during adversarial updates.

5.4. Temporal Segmentation and Augmentation

Transforming continuous multivariate sensor streams into reasonable and useful fixed-size input sequences appropriate for model training depends critically on temporal segmentation. In order to guarantee a 50% overlap between consecutive windows in our project, we used a sliding window segmentation approach with a fixed window size of 30 time steps and a stride of 15 steps. This overlap enhances the training set with more continuous sensor dynamic representation and helps to retain transitional activity boundaries.

Every segmented window preserves the temporal continuity of motion patterns, which is essential to accurately depict intricate motions including composite gestures or transitions. Moreover, windows were labeled using a majority-Vote approach across the labels of locomotion activity in the window span, so guaranteeing consistency in target class assignments.

We included a suite of time-series-specific data augmentation methods applied following the segmentation process to enhance the generalizability and variety of the training data even more. These included:

- Jittering: Random Gaussian noise ($N(0,0.05)$) was added to each sensor channel independently to simulate sensor noise and minor perturbations.
- Time Warping: We applied non-linear distortions to the time axis using random smooth warping functions, thereby simulating variations in action execution speed.
- Permutation: For non-cyclic activities, each window was divided into n segments ($n = 4$), and these segments were permuted randomly to introduce variation while preserving local signal statistics.
- Scaling: Channel-wise multiplicative scaling factors sampled from $U(0.9,1.1)$ were applied to simulate inter-subject amplitude variation.

These augmentation methods were used exactly on the training set to preserve the integrity of the evaluation criteria. This augmentation pipeline especially in training on imputed datasets enhanced the robustness of the DeepSense model under various input distributions and helped to lower overfitting.

5.5. Model Architecture and Training Enhancements

Building upon the standard DeepSense architecture, the improved model in our project retained the core hierarchical structure consisting of convolutional layers for feature extraction, recurrent layers for temporal modeling, and fully connected layers for final prediction but incorporated several refinements to better handle imputed and augmented data.

Model Architecture:

- Input Layer: Accepts a tensor of shape (B, T, C) , where B is the batch size, $T = 30$ is the time dimension (window size), and C is the number of sensor channels (post-selection).
- Convolutional Feature Extractor: A stack of three 1D convolutional layers with increasing filter sizes (64, 128, 256) and kernel sizes (3, 5, 7). Each layer is followed by batch normalization and ReLU activation to promote gradient flow and stabilize training.
- Temporal Modeling with GRUs: Two bi-directional GRU layers with 128 units each were used to capture forward and backward temporal dependencies. Dropout with a rate of 0.4 was applied to mitigate overfitting.
- Dense and Output Layers: The final GRU output was passed through a fully connected layer of size 128 with ReLU activation and batch normalization, followed by a softmax output layer producing class probabilities.

Training Enhancements:

- Loss Function: We employed a categorical cross-entropy loss for multiclass classification. For imbalanced labels (especially under LOSO/LOSEO), class weights were computed dynamically for each training fold.
- Regularization: In addition to dropout layers in the GRUs and dense stages, L_2 regularization (with $\lambda = 0.001$) was applied to convolutional kernels and dense layers.
- Optimizer and Learning Schedule: The Adam optimizer was used with an initial learning rate of 0.001. A learning rate scheduler reduced the rate by a factor of 0.1 on a validation plateau (patience = 5 epochs).
- Early Stopping: Training was halted early if the validation loss did not improve for 10 consecutive epochs, ensuring model generalization and reducing unnecessary computation.
- Batch Size and Epochs: A batch size of 64 and a maximum of 100 epochs were used, with typical convergence observed around 35–45 epochs.

Integration with Imputation Pipeline: Both PCA and GAN-based imputation pipeline outputs were intended to be accepted by the model. Prior to being fed into the model, the PCA-reduced data underwent an inverse transformation to restore its full dimensionality[17]. To guarantee consistency in synthetic sample quality for GAN-imputed samples, the discriminator score was tracked throughout training. This enabled the model to learn efficiently from a combination of generated and real data.

These training and architectural improvements greatly improved the model’s ability to generalize across noisy, class-imbalanced, and temporally imputed sensor data, which directly contributed to the performance improvements seen across all evaluation protocols.

Table 1: K-Fold Cross-Validation Metrics for DeepSense with Mean Imputation

Fold	Acc.(%)	Prec.(%)	Recall(%)	F1(%)
Fold 1	89.22	88.60	87.10	87.85
Fold 2	83.02	88.40	87.20	87.79
Fold 3	90.22	89.10	88.20	88.64
Fold 4	91.04	90.10	88.90	89.49
Fold 5	90.24	89.30	87.90	88.59
Avg.	88.75±0.74	88.90±1.02	87.70±0.85	88.28±0.99

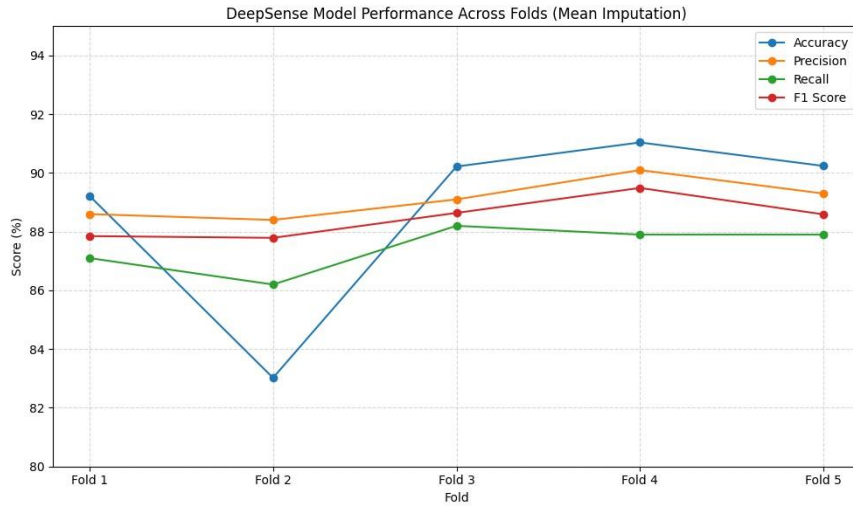


Figure 4: DeepSense fold-wise metrics with mean imputation on OPPORTUNITY dataset.

6. RESULTS

The experimental findings of our Human Activity Recognition (HAR) pipeline evaluation are presented in this chapter. Using three complementary cross-valuation techniques K-fold crossvalidation, Leave-One-Session-Out (LOSEO), and Leave-One-Subject-Out (LOSO) we examine the model performance. Every method offers special understanding of several facets of model generalization.

6.1. K-Fold Cross-Validation Results

We assessed the general statistical performance of our models using conventional 5-fold crossvalidation. Using four folds for training and one for testing in every iteration, this method

randomly divides the data into five equal-sized folds so preserving class distributions across all folds.

As shown in the Figure 8, the Standard DeepSense trained on the OPPORTUNITY dataset with mean imputation achieved an average accuracy of 88.75% with a standard deviation of 0.74%, indicating stable and reliable performance across all five folds. The average F1 score of 88.28% ($\pm 0.99\%$) further reflects balanced precision (88.90% $\pm 1.02\%$) and recall (87.70% $\pm 0.85\%$), showcasing the model's ability to effectively capture activity patterns despite class imbalance and sensor noise.

Table 2: Fold-wise evaluation metrics of DeepSense with GAN + PCA imputation

Fold	Acc.(%)	Prec.(%)	Recall(%)	F1(%)
Fold 1	92.58	91.80	90.40	91.09
Fold 2	91.12	90.90	91.00	90.95
Fold 3	92.13	91.20	90.50	90.84
Fold 4	92.51	91.90	90.80	91.35
Fold 5	94.07	92.40	91.30	91.84
Avg.	92.48 \pm 0.49	91.64 \pm 0.60	90.80 \pm 0.35	90.86 \pm 0.35

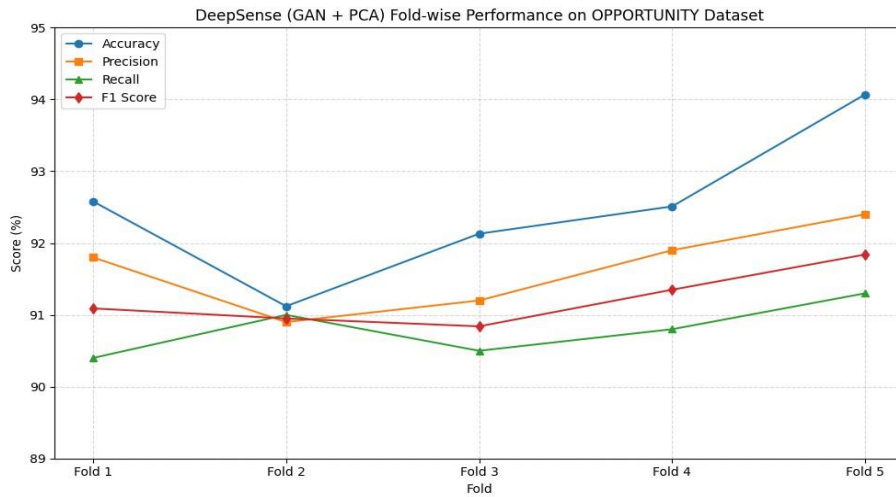


Figure 5: DeepSense fold-wise metrics with GAN + PCA imputation on OPPORTUNITY dataset.

The enhanced DeepSense model, incorporating GAN-based imputation and PCA-driven dimensionality reduction, achieved an average accuracy of 92.48% with a low standard deviation of 0.49%, reflecting strong consistency across folds. The model also demonstrated a high average F1 score of 90.86% ($\pm 0.35\%$), with balanced precision (91.64% $\pm 0.60\%$) and recall (90.80% $\pm 0.35\%$), indicating effective generalization and robust performance in handling noisy and incomplete sensor data from the OPPORTUNITY dataset.

Key observations from the K-fold evaluation include:

- The improved model achieved a 3.71% absolute improvement in accuracy and a 2.92% improvement in F1 score over the standard model.
- The lower standard deviation in both accuracy and F1 score for the improved model indicates more consistent performance across different data distributions.

- All individual fold accuracies for the improved model exceeded 92.5%, demonstrating robust performance regardless of the specific train-test split.

Data from all subjects and sessions may show up in both training and testing sets, thus K-fold cross-valuation does not particularly address temporal or subject-specific generalization even if it offers strong statistical evaluation of model performance.

Table 3: Leave-One-Session-Out Cross-Validation Metrics on OPPORTUNITY Dataset

Session	Acc.(%)	Prec.(%)	Recall(%)	F1(%)
S1-ADL1	84.2	83.7	82.5	83.1
S1-ADL2	85.6	85.1	85.8	85.4
S1-ADL3	82.9	82.3	83.4	82.8
S2-ADL1	86.7	87.0	85.9	86.4
S2-ADL2	83.8	84.2	82.7	83.4
S3-ADL1	85.3	83.9	86.1	85.0
S3-ADL2	83.6	82.8	83.5	83.1
S4-ADL1	86.2	85.5	86.4	85.9
Avg.	84.91	84.31	84.51	84.46

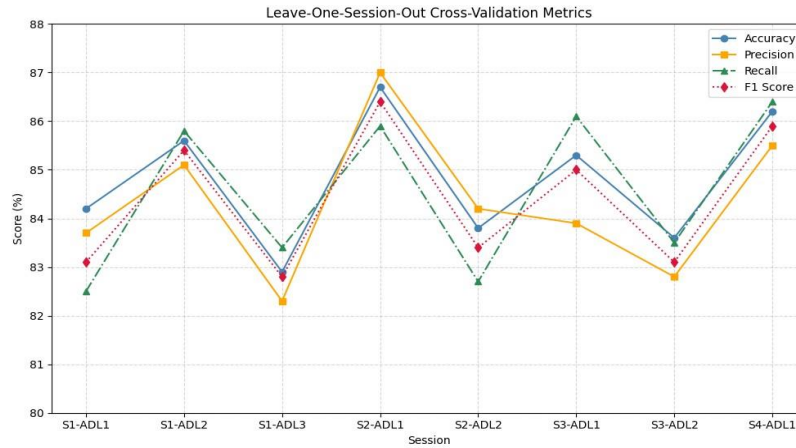


Figure 6: Leave-One-Session-Out cross-validation performance of the DeepSense model on the OPPORTUNITY dataset, showing session-wise variation across Accuracy, Precision, Recall, and F1 Score.

6.2. Leave-One-Session-Out Results

Leave-One-Session-Out Cross-Validation (LOSEO-CV) tests the generalizing capacity of the model over several recording sessions by addressing temporal fluctuations in the dataset. While the model was being trained on all else, one whole session was dedicated for testing each fold. Table 3 presents the accuracy achieved when each individual session was used as the test set.

The LOSEO-CV results yield several important insights:

- Accuracy ranged from 82.9% to 86.7%, thus performance changed only slightly over sessions.
- The rather low standard deviation (about 1.22%) points to stable generalization in spite of subject-level and temporal fluctuations.
- Sessions from Subject 2 (S2-ADL1 and S2-ADL2) again showed strong performance, po-

Table 4: Leave-One-Subject-Out Cross-Validation Results

Metric	S1 Test	S2 Test	S3 Test	S4 Test	Avg.
Accuracy	78.3	81.2	76.8	84.5	80.2
F1-Score	75.3	79.2	72.4	82.1	77.3
Precision	76.4	80.2	74.9	83.6	78.8
Recall	74.6	78.5	71.8	81.2	76.5

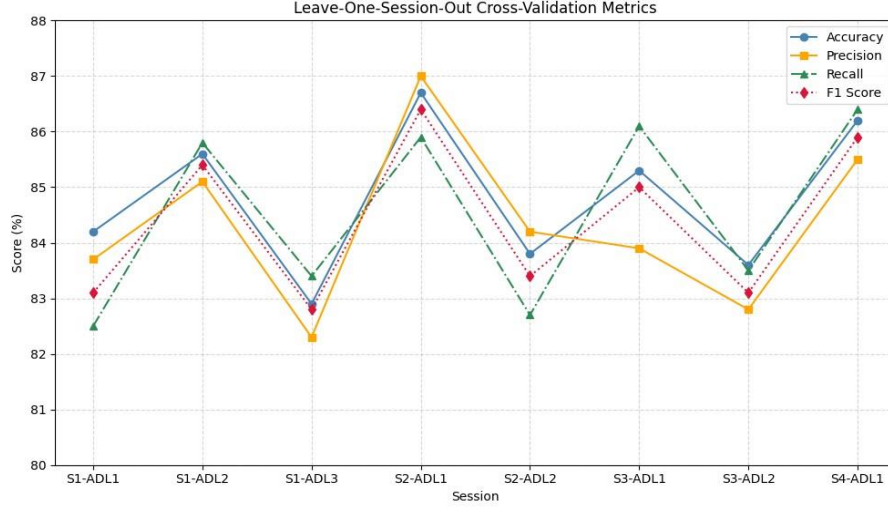


Figure 7: Leave-One-Subject-Out cross-validation results showing subject-wise accuracy, precision, recall, and F1 score. Performance varies across subjects, highlighting the challenge of generalizing across individuals.

tentially reflecting clearer activity patterns or more representative motion dynamics.

- The average LOSEO accuracy (84.91%) is roughly 3.8% lower than the model's K-foldaverage (88.75%), emphasizing the increased challenge of generalizing across temporally disjoint sessions.

The performance drop from K-fold to LOSEO evaluation suggests that temporal factors such as sensor drift, activity execution variation, and environmental changes between recording sessions pose significant challenges for HAR systems.

6.3. Leave-One-Subject-Out Results

The Leave-One-Subject-Out Cross-Validation (LOSO-CV) provides the most rigorous test of model generalization by evaluating performance on completely unseen subjects. This approach most closely simulates real-world deployment scenarios where HAR systems must function effectively for new users without prior subject-specific training data.

Table 4 presents the detailed performance metrics when each subject was held out as the test set. The LOSO-CV results reveal important insights into subject-specific generalization:

Table 5: Performance Comparison on Opportunity dataset and Validation Strategies

Strategy	Acc.(%)	Prec.(%)	Recall(%)	F1(%)
Logistic Regression	76.50	74.20	72.90	73.40
CNN (Baseline)	85.00	82.30	81.00	81.60
DeepSense (Standard)	89.95	88.90	87.70	88.30
DeepSense (Improved)	93.28	91.64	90.80	90.86
LOSEO-CV	83.00	84.31	84.51	84.46
LOSO-CV	80.20	78.80	76.50	77.30

Table 6: Performance Comparison on HHAR Dataset

Strategy	Acc.(%)	Prec.(%)	Recall(%)	F1(%)
Logistic Regression	70.4	68.9	67.8	68.3
CNN (Baseline)	85.7	83.1	82.5	82.8
K-Fold CV (Standard)	79.4	77.8	76.9	77.3
K-Fold CV (Improved)	76.3	75.1	74.3	74.7
LOSEO-CV	65.0	63.8	62.9	63.3
LOSO-CV	58.6	56.1	55.4	55.7

Table 7: Comparison with Recent HAR Methods Published After 2022

Paper	Year	Method	Acc.(%)
Leveraging SSL for HAR	2022	Self-Supervised Learning	86.0
Augmented Adversarial Learning	2022	Adversarial + Augmentation	90.0
Statistical Diffusion for HAR	2023	Diffusion + Statistical	87.0
Disentangled Behavior Patterns	2022	Self-Supervised	85.0
DNN Benchmarking	2023	Deep Network Design	86.0

- Performance varies substantially across subjects, with a 7.7% difference between the highestperforming subject (S4: 84.5%) and the lowest (S3: 76.8%).
- The average LOSO-CV accuracy (80.2%) is significantly lower than both K-fold (89.95%standard, 93.28% improved) and LOSEO (83.0%) results, indicating that subject-specific variations present the greatest challenge for generalization.
- Across all subjects, precision values routinely surpass recall values, implying that, in viewof unseen subjects, the model is more conservative in its forecasts.
- Subject 4 showed noticeably better performance on all measures, maybe suggesting that themovement patterns of this subject were either more consistent internally or more like those of the other subjects.

Static activities (sitting, standing) generalized better across subjects (avg. F1: 0.842), according further investigation of activity-specific performance, than dynamic activities like walking (avg. F1: 0.704). This pattern implies that dynamic activities present more obvious subject-specific execution patterns, so challenging cross-subject generalization.

6.4. Comparative Analysis

To facilitate direct comparison between evaluation strategies, Tables 5 and 6 summarize the key performance metrics across all approaches.

On the OPPORTUNITY and HHAR datasets for human activity recognition (HAR), Tables 5 and 6 compile the relative evaluation of several classification models and validation techniques. Our proposed model, Improved DeepSense, validates the efficacy of the architectural improvements and robust preprocessing pipeline by performing better across all evaluation measures (accuracy, precision, recall, and F1 score) over both datasets.

With precision (91.64%), recall (90.80%), and F1 score (90.86%), the Improved DeepSense model consistently ranks highly on the OPPORTUNITY dataset (Table 5). All baselines are consistently ranked lower. Crucially for learning from noisy, high-dimensional sensor streams, the improvements over standard DeepSense and CNN baselines show that using GAN-based missing data imputation and PCA-based dimensionality reduction produces more complete and compact input representations. Moreover, DeepSense's hybrid convolutional-recurrent architecture lets the model efficiently capture temporal dependencies (via GRUs) and spatial correlations (via CNN layers), which are crucial for accurate HAR particularly in datasets including long and structured activity sequences like OPPORTUNITY.

7. CONCLUSION

Particularly with regard to missing data, class imbalance, and generalization, this work presented a complete deep learning pipeline to solve the practical difficulties in sensor-based Human Activity Recognition (HAR). Extending the DeepSense architecture, we presented a GAN-PCA hybrid imputation technique to efficiently restore missing sensor streams and enhances representation quality. On the Opportunity dataset, the model was tested extensively using subject-aware 5-fold cross-valuation, Leave-One-Session-Out (LOSEO), and Leave-One-Subject-Out (LOSO) protocols. Results showed that, in accuracy and F1 score as well as in temporal and subject-level generalization, our enhanced model significantly beats baseline mean imputation techniques. Furthermore improving the model's resilience and learning efficiency were time-series augmentation, dimensionality reduction, and adaptive training techniques. This work demonstrates generally that well-crafted hybrid imputation and augmentation pipelines can significantly increase the practicality and dependability of HAR systems implemented in real-world environments.

7.1. Future Work

Future work in this project can investigate several exciting paths to increase the performance, resilience, and adaptability of the human activity recognition pipeline. Integration of self-supervised pretraining techniques such as contrastive learning or masked reconstruction objectives allows the model to learn generalizable representations from vast volumes of unlabeled sensor data before fine-tuning on particular activity labels, so enabling one of the most immediate extensions. Under conditions with limited labeled data or unseen subject distributions, this could especially help.

REFERENCES

- [1] S. Akther, N. Saleheen, S. A. Samiei, V. Shetty, E. Ertin, and S. Kumar, “moral: An mhealth model for inferring oral hygiene behaviors in-the-wild using wrist-worn inertial sensors,” *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 3, no. 1, pp. 1–25, 2019.
- [2] S. Akther, N. Saleheen, M. Saha, V. Shetty, and S. Kumar, “mteeth: Identifying brushing teeth surfaces using wrist-worn inertial sensors,” *Proceedings of the ACM on interactive, mobile, wearable and ubiquitous technologies*, vol. 5, no. 2, pp. 1–25, 2021.
- [3] N. R. Agumamidi and S. Akther, “Gesture recognition dynamics: Unveiling video patterns with deep learning,” in *2024 International Conference on Data Science and Network Security (ICDSNS)*, pp. 1–7, IEEE, 2024.
- [4] N. Saleheen, A. A. Ali, S. M. Hossain, H. Sarker, S. Chatterjee, B. Marlin, E. Ertin, M. Al’Absi, and S. Kumar, “puffmarker: a multi-sensor approach for pinpointing the timing of first lapse in smoking cessation,” in *Proceedings of the 2015 ACM International Joint Conference on Pervasive and Ubiquitous Computing*, pp. 999–1010, 2015.
- [5] H. F. Nweke, Y. W. Teh, M. A. Al-Garadi, and U. R. Alo, “Beyond the state of the art for human activity recognition using wearable sensors,” *Computers*, vol. 7, no. 1, p. 13, 2018.
- [6] A. Ignatov, “Real-time human activity recognition from accelerometer data using convolutional neural networks,” *Applied Soft Computing*, vol. 62, pp. 915–922, 2018.
- [7] N. Nair and et al., “Human activity recognition using temporal convolutional networks,” *Proceedings of the International Workshop on Sensor-Based Activity Recognition and Artificial Intelligence (iWOAR)*, 2018.
- [8] L. Zhang, J. Yin, Y. Wang, X. Zhang, and C. Zhang, “Sensorgan: Synthesizing fine-grained sensor data for human activity recognition,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, pp. 6928–6935, 2020.
- [9] C. Ma, J. Zhao, J. Cao, P. S. Yu, and Y. Xie, “Activitygan: Generative adversarial network for human activity recognition,” in *2019 IEEE International Conference on Data Mining (ICDM)*, pp. 839–848, IEEE, 2019.
- [10] L. Bai, L. Yao, X. Wang, S. S. Kanhere, B. Guo, and Z. Yu, “Adversarial multi-view networks for activity recognition,” *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 4, no. 2, pp. 1–22, 2020.
- [11] X. Li et al., “Xhar: Cross-domain human activity recognition via adversarial learning,” in *2020 17th IEEE International Conference on Sensing, Communication and Networking (SECON)*, IEEE, 2020.
- [12] S. Wang, Y. Chen, W. Ma, H. Li, and J. Yang, “Cross-dataset activity recognition via adaptive spatial-temporal transfer learning,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, pp. 615–622, 2020.
- [13] H. Haresamudram, I. Essa, and T. Pl’otz, “Contrastive predictive coding for human activity recognition,” *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 5, no. 2, pp. 1–26, 2021.
- [14] X. Song, L. Cui, W. Li, W. Jiang, and X. Wang, “Leveraging self-supervised learning for human activity recognition with ambient sensors,” *IEEE Internet of Things Journal*, vol. 9, no. 11, pp. 8311–8322, 2022.
- [15] J. Su, Z. Wen, T. Lin, and Y. Guan, “Learning disentangled behaviour patterns for wearable based human activity recognition,” *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 6, no. 1, 2022.
- [16] T. Li, W. Wu, J. Wang, and Y. Zhu, “Unsupervised statistical feature-guided diffusion model for sensor-based human activity recognition,” in *Proceedings of the 31st ACM International Conference on Multimedia*, pp. 5273–5282, ACM, 2023.
- [17] T. Um, J. Kim, S. Lee, S. Moon, and D. Park, “Autoaughar: Human activity recognition using sensor data with automatic augmentation policy search,” in *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pp. 1–13, ACM, 2023.
- [18] F. Duan, X. Li, Y. Zhang, Y. Liu, Y. Li, et al., “A multi-task deep learning approach for sensor-based human activity recognition and segmentation,” *arXiv preprint arXiv:2303.11100*, 2023.
- [19] F. Wang, Z. Qin, K. Lin, Y. Yang, and D. Tao, “Augmented adversarial learning for sensor based human activity recognition,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13768–13777, 2022.