
A DUAL-USE FRAMEWORK FOR CLINICAL GAIT ANALYSIS: ATTENTION-BASED SENSOR OPTIMIZATION AND AUTOMATED DATASET AUDITING

Hamidreza Sadeghsalehi
Imperial College London

ABSTRACT

Background: Objective gait analysis using wearable sensors is critical for managing neurological and orthopedic conditions. However, optimizing sensor placement for specific clinical tasks and mitigating the impact of hidden dataset biases on artificial intelligence (AI) models remain significant challenges. **Methods:** A multi-stream, attention-based deep learning model was applied to the Voisard et al. (2025) multi-cohort gait dataset, addressing four binary classification tasks: Parkinson’s Disease (PD) screening, Osteoarthritis (OA) screening, post-stroke asymmetry detection (CVA), and a differential diagnosis (PD vs. CVA). We implemented strategies to mitigate severe class imbalance and evaluated the model using imbalance-robust metrics (e.g., Balanced Accuracy, MCC) and bootstrapped confidence intervals. The model’s primary outputs were classification performance and a quantitative map of learned sensor importance weights for four body locations: Head (HE), Lower Back (LB), Left Foot (LF), and Right Foot (RF). **Results:** The model’s attention mechanism identified distinct, task-specific sensor synergies, such as a novel Head-Right-Foot (HE-RF) combination for PD screening and a Head-Left-Foot (HE-LF) combination for the complex PD-vs-CVA differential diagnosis. Critically, the model functioned as an automated data auditor by discovering a severe laterality confound in the dataset. For both OA and CVA screening—tasks where bilateral information is clinically essential—the model assigned over 70% of its attention exclusively to the Right Foot sensor, with near-zero (e.g., $<0.1\%$) attention to the Left Foot. This counter-intuitive result, confirmed by statistically significant 95% CIs, was a direct reflection of a severe laterality bias in the patient cohorts (15/0 right-sided OA; 47/2/0 right-dominant CVA) and not a valid clinical finding. **Conclusion:** The attention-based framework is a powerful dual-use tool. Its primary methodological contribution is its function as an automated data auditor capable of discovering and quantifying hidden dataset biases, a critical step for responsible AI development. Concurrently, when applied to valid data subsets, it provides a data-driven method for identifying minimal, clinically relevant sensor sets (e.g., HE+RF for PD) to optimize future clinical monitoring protocols.

Keywords Digital Biomarkers · Gait Analysis · Wearable Sensors · Deep Learning · Attention Mechanism · Model Interpretability · Dataset Bias · Data Auditing

1 Introduction

Gait impairment is a cardinal symptom and a major contributor to disability across a wide spectrum of chronic health conditions, particularly in neurology and orthopedics [1–3]. In Parkinson’s Disease (PD), characteristics such as reduced walking speed, shortened step length, and increased gait variability are not merely symptoms but are core to the disease’s motor pathology and progression [1, 4]. Similarly, in conditions like osteoarthritis (OA) or following a cerebrovascular accident (CVA), alterations in gait symmetry and kinematics are primary indicators of disease severity and functional limitation [2, 3]. Traditional clinical assessment of gait, however, often relies on observational scoring and qualitative descriptions; while valuable, these methods can lack the sensitivity and objectivity required to track subtle changes over time, evaluate the efficacy of interventions, or develop precise diagnostic and prognostic models [5, 6]. This gap has created a clear clinical imperative for quantitative, objective, and scalable methods for gait assessment.

The advent of wearable sensors, particularly inertial measurement units (IMUs), has catalyzed a paradigm shift in movement analysis [7, 8]. These low-cost, unobtrusive devices can capture high-fidelity kinematic data—including acceleration and angular velocity—in ambulatory, real-world settings, overcoming the spatial and temporal limitations of traditional laboratory-based, optical motion capture systems [5]. The resulting high-dimensional time-series data streams offer a rich substrate for the development of digital biomarkers capable of objectively quantifying pathological gait patterns [9]. Large-scale, clinically annotated datasets, such as the comprehensive collection provided by Voisard et al. (2025), are instrumental in bridging the gap between raw sensor data and translational clinical research, providing the necessary scale and diversity to develop robust analytical models [10].

Despite these advances, a fundamental and often overlooked question persists: for a given clinical question, what is the minimal yet sufficient set of sensors required for an accurate assessment [11–13]? Current research and clinical protocols frequently employ either a large array of sensors, which can be burdensome for patients and computationally intensive, or rely on heuristic choices for sensor placement. This one-size-fits-all approach is suboptimal, as different pathologies manifest through distinct kinematic signatures across different body segments. For example, the systemic motor deficits of PD, which affect both axial trunk control and appendicular limb movement, are fundamentally different from the localized joint-level impairments of unilateral knee OA [1, 3]. A logical and more efficient paradigm would therefore be task-dependent, tailoring the sensor configuration to the specific clinical question being asked.

Deep learning models, including Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), have emerged as state-of-the-art for analyzing complex time-series data from wearable sensors, demonstrating strong performance in tasks ranging from fall risk classification to gait recognition [14–16]. Within this domain, the attention mechanism, a technique originally developed for natural language processing, offers a powerful means of enhancing both model performance and interpretability [17]. By enabling a model to dynamically assign importance weights to different segments of its input, attention can be adapted for multi-sensor fusion, learning the relative contribution of each sensor’s data stream to a final prediction [18]. While the reliability of attention mechanisms for explaining individual predictions has been a subject of debate—as single-instance attention maps can be inconsistent across model initializations—their utility can be harnessed more robustly [19–22]. By aggregating attention weights across an entire cohort, the mechanism can be transformed from a potentially fragile explanator into a stable and powerful tool for discovering aggregate, population-level insights about feature importance for a given task.

Parallel to the challenge of optimizing sensor configurations is the pervasive and critical problem of dataset bias in medical AI [23–28]. AI models are susceptible to learning and amplifying biases present in their training data, which can lead to poor generalization, reinforcement of health inequities, and clinically harmful outcomes [23, 25]. Confounding variables—extraneous factors correlated with both input and outcome—can cause a model to learn spurious associations [29]. For instance, if a disease cohort is significantly older than a control cohort, a model may inadvertently become a highly accurate age detector rather than a true disease detector [25]. This underscores the urgent need for rigorous data auditing and governance as a foundational step in responsible AI development [27–30]. However, such auditing is typically a manual, resource-intensive process, ill-equipped to uncover subtle or unexpected confounders in large, complex datasets [29]. This contrasts with our work, which explores an emergent auditing capability embedded within the model itself.

This paper presents a dual contribution to the fields of digital health and medical AI. First, we introduce a novel, attention-based deep learning framework that provides a data-driven approach to identifying minimal, task-specific sensor configurations for clinical gait analysis. Second, and more significantly, we demonstrate that this same framework functions as a powerful tool for automated dataset auditing. We show how the model’s learned attention weights can quantitatively expose a critical, hidden laterality confound within a major public dataset, transforming the model from a simple classifier into an engine for data quality assurance and scientific discovery.

2 Methods

2.1 Study Dataset and Participants

This study utilized the “Dataset of Clinical Gait Signals with Wearable Sensors,” a publicly available, multi-pathology, and clinically annotated dataset curated by Voisard et al. (2025) [10]. The dataset comprises 1356 gait trials from a total of 260 participants, organized into three primary groups: Healthy Subjects (HS), patients with neurological conditions (Neuro), and patients with orthopedic conditions (Ortho) [10]. For the specific tasks defined in this study, participants were selected from four cohorts: Healthy Subjects (HS, $n=73$), Parkinson’s Disease (PD, $n=143$), Hip Osteoarthritis (HOA, $n=44$), and Cerebrovascular Accident (CVA, $n=47$) [10]. The original study was conducted in accordance with the Declaration of Helsinki, and all participants provided informed consent; full details regarding participant recruitment, clinical scoring, and ethical approvals are available in the original data descriptor [10].

2.2 Data Acquisition and Preprocessing

Gait data were collected using four IMU sensors (XSens or Technoconcept) placed on the Head (HE), Lower Back at the L5 vertebra level (LB), Left Foot (LF), and Right Foot (RF) [10]. Each sensor provided a 9-channel time-series data stream, consisting of 3D accelerometer signals, 3D gyroscope signals, and 3D gravity-corrected free acceleration signals [7]. All data were sampled at a frequency of 100 Hz. The standardized data collection protocol involved participants performing a 10-meter walk, followed by a 180-degree U-turn and a 10-meter walk back to the starting point [10]. For this analysis, the preprocessed data files provided by Voisard et al. were used; these contain synchronized, multi-sensor time-series data for each trial, along with metadata including detected gait events such as heel-strikes and toe-offs [10, 31].

2.3 Formulation of Clinical Classification Tasks

To evaluate the model’s ability to learn task-specific sensor configurations, four distinct binary classification tasks were formulated, each designed to simulate a relevant clinical question:

- **PD Screening:** Differentiating between participants with Parkinson’s Disease (PD) and Healthy Subjects (HS). This task simulates the use of gait analysis as a screening tool for a common neurodegenerative disorder [11].
- **OA Screening:** Differentiating between participants with Hip Osteoarthritis (HOA) and Healthy Subjects (HS). This task models the identification of gait patterns associated with a prevalent orthopedic condition [3].
- **Asymmetry Detection:** Differentiating between participants post-Cerebrovascular Accident (CVA) and Healthy Subjects (HS). Given the typically unilateral motor deficits following a stroke, this task was designed to test the model’s ability to identify asymmetrical gait [2].
- **Differential Diagnosis:** Differentiating between participants with Parkinson’s Disease (PD) and those with Cerebrovascular Accident (CVA). This represents a more complex clinical challenge, requiring the model to distinguish between two distinct pathologies that can both affect gait.

2.4 Attention-Based Multi-Stream Neural Network Architecture

A multi-stream neural network incorporating a sensor-level attention mechanism was designed to process the multi-modal sensor data [17, 18]. The architecture consists of three main components:

- **Modality-Specific Feature Extraction:** The model employs four parallel and independent branches, one for each sensor location (HE, LB, LF, RF). Each branch is a 1D Convolutional Neural Network (1D-CNN) that processes the 9-channel time-series input from its corresponding sensor. These CNNs are designed to learn high-level, abstract feature representations from the raw sensor signals. The output of each branch is a fixed-length feature vector, $v_i \in \mathbb{R}^{128}$, where $i \in \{\text{HE, LB, LF, RF}\}$.
Each 1D-CNN branch consists of three sequential convolutional layers with filter counts of 32, 64, and 128, respectively. All layers use a `kernel_size=15` and `padding=7`. Each convolutional layer is followed by a ReLU activation and a `MaxPool1d` layer with `kernel_size=2`. A final `AdaptiveAvgPool1d(1)` layer collapses the time dimension, producing the 128-dimensional feature vector for each sensor.
- **Sensor-Level Attention Mechanism:** The four feature vectors $\{v_{\text{HE}}, v_{\text{LB}}, v_{\text{LF}}, v_{\text{RF}}\}$ are fed into an attention module. This module consists of a single linear layer that learns to compute four unnormalized scalar importance scores, e_i . These scores are then passed through a softmax function to generate the final attention weights, α_i , which represent the model’s learned judgment of each sensor’s relative importance for the given task. The weights are positive and sum to one.

$$\alpha_i = \frac{\exp(e_i)}{\sum_{j \in \{\text{HE, LB, LF, RF}\}} \exp(e_j)}$$

- **Fusion and Classification:** The four feature vectors are fused into a single context vector via a weighted summation using their learned attention weights: $c = \sum_i \alpha_i v_i$. This context vector, which represents a task-optimized summary of the information from all four sensors, is then passed through a final 2-layer fully-connected classifier (with ReLU activation and Dropout (0.5)) to produce a single logit for binary class prediction.

Table 1: Model classification performance and mean sensor attention weights on the test set. Performance is reported with metrics robust to class imbalance. Attention weights are shown as mean % (95% Confidence Interval).

Clinical Task	N (Pat./Ctrl.)	ROC-AUC	PR-AUC	Bal. Acc.	MCC	Sens. (Rec.)	Spec.	HE (%)	LB (%)	LF (%)	RF (%)
PD Screening	15 (4 Pos/11 Neg)	0.821	0.671	0.639	0.369	0.333	0.945	32.9 (28.3–37.4)	3.4 (3.1–3.8)	11.2 (10.3–12.4)	52.5 (48.4–56.9)
OA Screening	14 (2 Pos/12 Neg)	0.990	0.906	0.942	0.682	1.000	0.885	6.3 (4.9–7.7)	22.7 (18.2–27.5)	0.1 (0.0–0.1)	71.0 (65.7–76.0)
Asymmetry (CVA)	19 (8 Pos/11 Neg)	0.950	0.753	0.747	0.521	0.563	0.932	22.6 (17.1–28.1)	0.0 (0.0–0.0)	0.0 (0.0–0.0)	77.4 (71.8–82.9)
Differential (PD vs. CVA)	11 (4 Pos/7 Neg)	0.657	0.821	0.607	0.202	0.500	0.714	51.5 (41.8–60.6)	0.2 (0.1–0.2)	46.8 (37.9–56.4)	1.5 (1.1–2.0)

2.5 Experimental Design and Statistical Analysis

To ensure robust and generalizable evaluation, the dataset was split at the patient level into training (70%), validation (15%), and test (15%) sets, a critical methodological step to prevent data leakage [29, 32]. The exact class distributions (number of patients vs. controls) for each test set are reported in Table 1.

For each of the four clinical tasks, a separate model was trained from a random initialization. Given the significant class imbalance observed in the data splits, models were trained using a weighted binary cross-entropy loss function (`nn.BCEWithLogitsLoss` with `pos_weight`), where the positive class weight was set to $\frac{\text{negative_samples}}{\text{positive_samples}}$ to penalize misclassification of the minority class.

Training was performed using the Adam optimizer with a learning rate of $1e-4$ for 50 epochs, a `batch_size` of 32, and early stopping with a `patience` of 5 epochs based on validation loss.

Model performance was evaluated on the held-out test set using metrics robust to class imbalance, including: Area Under the Receiver Operating Characteristic Curve (ROC-AUC), Precision-Recall AUC (PR-AUC), Balanced Accuracy, Matthew’s Correlation Coefficient (MCC), Sensitivity (Recall), and Specificity. Ninety-five percent confidence intervals for these metrics were calculated using 1000 bootstrap samples.

The primary explanatory output of the model, the “Sensor Importance Map,” was derived for each task by calculating the mean of the learned attention weights ($\alpha_{HE}, \alpha_{LB}, \alpha_{LF}, \alpha_{RF}$) across all samples in the test cohort. We calculated 95% confidence intervals for these mean attention weights using 1000 bootstrap samples to provide a measure of statistical uncertainty in the model’s learned strategy.

3 Results

3.1 Performance on Imbalanced Clinical Tasks

Our initial analysis revealed significant class imbalance across all tasks. Consequently, we adopted a weighted loss function and evaluated performance using imbalance-robust metrics, which are presented in Table 1. The model demonstrated strong discriminative ability on the OA Screening (ROC-AUC 0.990) and Asymmetry Detection (ROC-AUC 0.950) tasks. Performance on the PD Screening (ROC-AUC 0.821) and the highly complex Differential Diagnosis (ROC-AUC 0.657) tasks was more modest, reflecting the clinical difficulty of these classifications.

The metrics highlight the challenge of class imbalance. For PD Screening, despite a high Specificity (0.945), the Sensitivity (0.333) and Balanced Accuracy (0.639) indicate the model still struggles to identify the minority PD class, even with re-weighting. Conversely, the model achieved perfect Sensitivity (1.000) for OA Screening, correctly identifying all positive-class patients in the test set. These results underscore the necessity of using metrics beyond simple accuracy to interpret model performance in clinical datasets.

3.2 Attention Mechanism Discovers Task-Specific Sensor Importance Maps

The primary goal of the study was to determine if the model could learn to prioritize different sensors based on the clinical question. The aggregated attention weights, summarized in Table 1 and visualized in Figures 1–4, revealed distinct and complex sensor importance maps for each task.

PD Screening (Figure 1): For the PD Screening task, the model learned a novel synergy, allocating the most attention to the Right Foot (RF) sensor (52.5%) and the Head (HE) sensor (32.9%). This data-driven finding suggests a strategy of integrating information from both appendicular (foot) and axial (head) segments. This aligns well with the known pathophysiology of Parkinson’s Disease, which is characterized by both appendicular motor deficits (e.g., shuffling, reduced step length) and axial symptoms (e.g., postural instability, head tremor) [1, 4]. This HE-RF combination is a plausible, data-driven hypothesis for an efficient PD screening protocol.

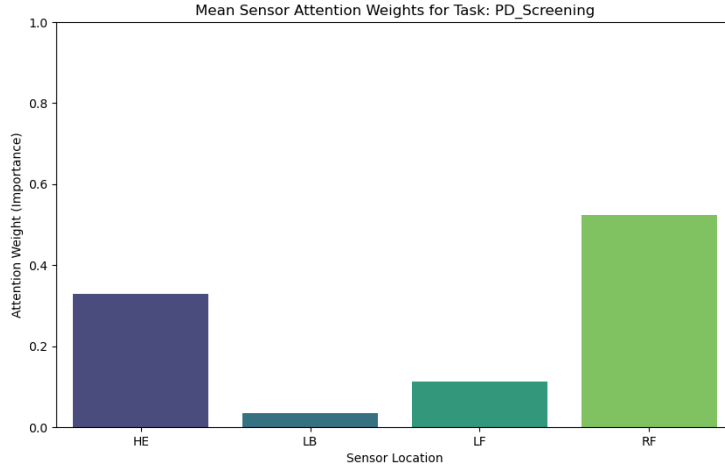


Figure 1: Mean Sensor Attention Weights for Parkinson’s Disease (PD) Screening. The model learned to prioritize a novel synergy between the Right Foot (RF) (52.5%) and Head (HE) (32.9%) sensors. This data-driven hypothesis aligns with the clinical presentation of PD, which involves both appendicular (gait) and axial (postural) motor deficits. Error bars represent 95% confidence intervals calculated from 1000 bootstrap samples, indicating the statistical uncertainty in the estimate of the mean attention.

OA and Asymmetry Screening (Figure 2 & 3): In stark contrast, the results for OA Screening and Asymmetry Detection (CVA) are not clinical findings, but rather a clear validation of our auditing methodology. For both tasks, the model developed an overwhelming and statistically significant reliance on a single sensor: the Right Foot (RF). In the OA task, the RF sensor received 71.0% of the attention, and in the CVA task, it received 77.4%. Critically, the Left Foot (LF) sensor was actively and confidently ignored in both scenarios, with mean attention weights of 0.1% and 0.0%, respectively. The 95% confidence intervals for the LF sensor attention (e.g., [0.0–0.1]) confirm that this is not an unstable result but a deliberate and stable strategy. This finding is highly counter-intuitive, as a clinical diagnosis of asymmetry or OA fundamentally relies on comparing bilateral limb information [2, 3]. The model’s categorical rejection of the left foot strongly suggested it had discovered a dataset confound, which we confirm in Section 3.3.

Differential Diagnosis (Figure 4): The most sophisticated sensor utilization strategy was observed in the Differential Diagnosis task. To distinguish between PD and CVA, the model learned a unique, non-obvious synergy, focusing almost equally on the Head (HE) sensor (51.5%) and the Left Foot (LF) sensor (46.8%). The Right Foot, which dominated the CVA screening task, was now almost completely ignored (1.5%). This demonstrates a remarkable level of adaptive learning. The model appears to have learned that the RF sensor is a confounder for CVA and thus provides no utility in differentiating it from PD. Instead, it reasons that it must compare central postural control (Head, often affected in PD) with a non-confounded limb (Left Foot) to solve the task. The wide, overlapping confidence intervals for HE and LF weights also suggest a high-variance strategy, reflecting the nuanced clinical challenge of this task.

3.3 Automated Discovery of a Critical Dataset Confounder

The anomalous and counter-intuitive findings from the OA and CVA screening tasks prompted a deeper investigation into the dataset itself. The model’s extreme and consistent preference for the right foot sensor was not a model failure but a successful and quantitative discovery of a significant sampling bias in the underlying Voisard et al. (2025) dataset.

A review of the cohort demographics detailed in the original data descriptor publication revealed a critical confound related to pathology laterality [10]:

- **HOA Cohort:** The laterality of hip osteoarthritis was predominantly right-sided. The paper reports a right/left laterality of 15/0, meaning all patients with unilateral hip pathology had it on the right side.
- **CVA Cohort:** The cohort of stroke survivors was overwhelmingly composed of individuals with right-sided motor deficits. The reported laterality (right/left/ambidextrous) was 47/2/0, indicating that nearly all unilateral cases involved the right side of the body.

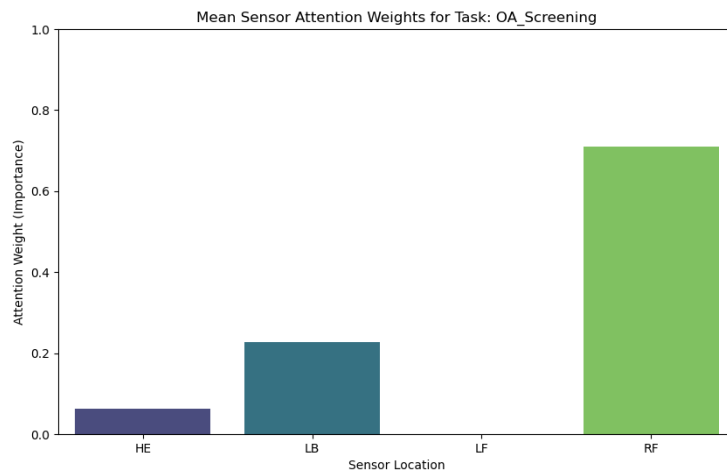


Figure 2: Mean Sensor Attention Weights for Osteoarthritis (OA) Screening. This result demonstrates the model’s function as a data auditor. It assigned over 70% of its attention to the Right Foot (RF) sensor while statistically ignoring the Left Foot (LF) (0.1%). This counter-intuitive pattern reflects the model’s discovery of a dataset confound (a 15/0 right-sided laterality bias in the HOA cohort), not a generalizable clinical signature of the disease. Error bars represent 95% confidence intervals (1000 bootstrap samples).

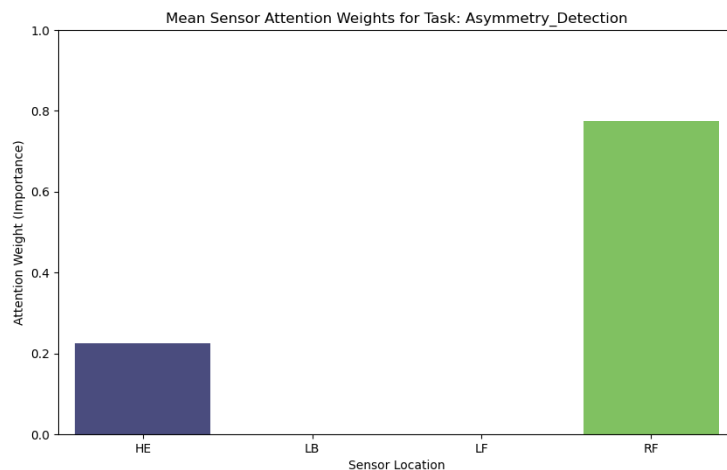


Figure 3: Mean Sensor Attention Weights for Asymmetry Detection (CVA). Similar to the OA task, the model relied almost exclusively on the Right Foot (RF) sensor (77.4%). The lack of attention to the Left Foot (LF) (0.0%) for a task predicated on asymmetry demonstrates the exploitation of a dataset confounder (a 47/2/0 right-dominant laterality bias in the CVA cohort), validating the attention mechanism’s ability to flag hidden biases. Error bars represent 95% confidence intervals (1000 bootstrap samples).

The model, as an optimization algorithm, learned the most parsimonious path to a correct classification. In this biased dataset, the most predictive feature for identifying a CVA or HOA patient was not a complex measure of “asymmetry” but simply the presence of an “abnormality in the right foot’s signal.” The model correctly deduced that the Left Foot sensor provided redundant or non-informative data for these specific tasks and therefore assigned it a near-zero attention weight. This finding powerfully demonstrates the capacity of an interpretable model to function as an automated data auditor, flagging hidden confounders that could severely limit the generalizability of any model trained on such data.

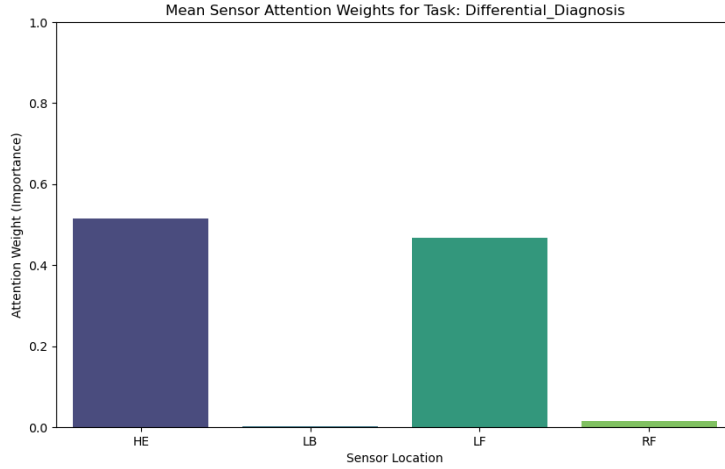


Figure 4: Mean Sensor Attention Weights for Differential Diagnosis (PD vs. CVA). To solve this complex task, the model learned a sophisticated synergy, allocating importance to the Head (HE) sensor (51.5%) and the Left Foot (LF) sensor (46.8%). It learned to ignore the confounded Right Foot sensor (1.5%), suggesting a strategy of comparing central postural control (Head) with a non-confounded limb (Left Foot). Error bars represent 95% confidence intervals (1000 bootstrap samples).

4 Discussion

This study demonstrates the dual utility of an attention-based deep learning framework for wearable sensor-based gait analysis. The principal findings are twofold: first, the model successfully identified minimal, task-specific sensor configurations, providing data-driven hypotheses for optimizing clinical monitoring protocols. Second, and more significantly, the model’s interpretability layer served as an effective, automated tool for discovering and quantifying a severe laterality confound in a major public dataset, highlighting a novel application for such architectures in data quality assurance and scientific discovery.

4.1 Clinical Implications for Optimized Gait Analysis Protocols

The sensor importance maps generated by our model offer actionable hypotheses for designing more efficient and less burdensome clinical gait analysis protocols. For PD screening, the model’s reliance on a Head-Right-Foot (HE-RF) synergy suggests that a minimal two-sensor setup, capturing both axial and appendicular features, could be a powerful data-driven hypothesis for future investigation [4, 11]. This has significant implications for remote patient monitoring and decentralized clinical trials, where patient adherence and ease of use are paramount [9]. For more complex tasks like differential diagnosis, the model’s learned HE-LF synergy suggests that capturing information from anatomically distant body segments, reflecting both central and peripheral nervous system function, may be necessary. These data-driven insights move beyond heuristic sensor placement, offering a principled approach to protocol design that is tailored to the specific clinical question at hand.

4.2 Interpretable AI as an Engine for Automated Scientific Discovery and Data Auditing

The most profound contribution of this work is methodological. The discovery of the laterality confound was not a limitation of our study but its most powerful result. It showcases a new paradigm for the use of interpretable AI: not merely as a tool for prediction, but as an active instrument for probing, validating, and understanding the datasets upon which our models are built. While the broader literature on responsible AI rightly focuses on mitigating demographic and societal biases to ensure fairness and equity [23, 24, 26–28, 30], our work demonstrates that these same principles and tools can be extended to uncover clinical and sampling biases that are equally detrimental to a model’s scientific validity and generalizability. This automated approach represents a scalable and powerful alternative to the often manual and time-consuming processes of data auditing [29]. An interpretable model can serve as a “canary in the coal mine,” quantitatively flagging potential confounders that might be missed by human reviewers examining summary statistics.

By training models on various sub-tasks within a dataset and analyzing their learned attention patterns, researchers can systematically screen for hidden biases, thereby strengthening the foundations of evidence-based digital medicine.

4.3 Limitations and Considerations

It is crucial to acknowledge the limitations of this study. First, the specific sensor importance maps for the OA and CVA tasks are intrinsically linked to the Voisard et al. (2025) dataset and its identified confounder. They should not be interpreted as generalizable gait signatures for these conditions but rather as signatures of the dataset’s specific bias. Second, the compelling sensor importance maps derived for PD Screening and Differential Diagnosis must be treated as data-driven hypotheses that require urgent validation on independent, multi-centric datasets. It is paramount that these validation datasets are themselves screened for similar biases (e.g., laterality) and other potential confounders not analyzed here, such as systematic differences in age, disease severity, or walking speed between cohorts, which could also be learned by the model.

Furthermore, our methodology must be contextualized within the active ‘attention as explanation’ debate [19–22]. We acknowledge that instance-level attention weights can be unstable and may not reliably reflect a model’s decision-making process for a single prediction. However, our approach deliberately avoids this contentious use case. Instead, we use cohort-level aggregation—averaging the attention weights across all samples in the test set—to identify stable, population-level importance maps. We contend this aggregation effectively averages out instance-level noise and instability, revealing a robust and meaningful signal about the model’s learned strategy for the task as a whole. This aggregated map, supported by bootstrapped confidence intervals, reflects the underlying data distribution and feature importance, transforming attention from a fragile ‘explainer’ into a reliable tool for dataset auditing.

Finally, our choice of a CNN-based architecture is one of many valid approaches for time-series analysis; other architectures, such as Transformers, might yield different performance or insights [33, 34].

4.4 Future Directions

This work opens several avenues for future research. The immediate next step is the prospective validation of the hypothesized minimal sensor sets (e.g., the HE-RF configuration for PD monitoring) in a real-world clinical study to confirm their efficacy and clinical utility. Beyond this, the methodology can be expanded into a formal, generalized framework for AI-driven dataset auditing. Such a framework could be deployed as a standard pre-analysis step in digital health research to automatically screen for potential confounders.

Once a confounder is identified, future work should focus on its mitigation. This could involve integrating advanced techniques designed to train confounder-free models, such as the adversarial de-confounding methods cited in [35], to explicitly remove the influence of the discovered laterality bias from the learned feature representations. The long-term vision is the integration of these validated, optimized, and bias-aware models into clinical decision support systems and remote patient monitoring platforms, enabling more precise, efficient, and equitable care.

5 Conclusion

The attention-based deep learning framework presented in this study provides a robust, data-driven method for generating hypotheses about sensor optimization in clinical gait analysis. However, its true novelty and most significant contribution lie in its emergent capacity as an automated data auditor. By revealing hidden scientific insights and critical dataset flaws through its interpretable architecture, this methodology demonstrates a powerful dual-use capability. It represents a significant step forward in our ability to build more efficient, reliable, and responsible AI systems for the future of digital health.

6 Data Availability

The “Dataset of Clinical Gait Signals with Wearable Sensors” is publicly available at the Scientific Data page: <https://www.nature.com/articles/s41597-025-05959-w>. For citation and resolution, use the DOI link: <https://doi.org/10.1038/s41597-025-05959-w>.

7 Code Availability

The complete code base used for model training, evaluation, and generation of all results presented in this paper is available at: <https://github.com/hamidreza-s-salehi/GaitSensorAttention>.

References

- [1] Anat Mirelman, Paolo Bonato, Richard Camicioli, Terry D Ellis, Nir Giladi, Jamie L Hamilton, Chris J Hass, Jeffrey M Hausdorff, Elisa Pelosin, and Quincy J Almeida. Gait impairments in parkinson’s disease. *The Lancet Neurology*, 18(7):697–708, 2019.
- [2] Kara K Patterson, William H Gage, Dina Brooks, Sandra E Black, and William E McIlroy. Evaluation of gait symmetry after stroke: a comparison of current methods and recommendations for standardization. *Gait & posture*, 31(2):241–246, 2010.
- [3] Kathryn Mills, Michael A Hunt, and Reed Ferber. Biomechanical deviations during level walking associated with knee osteoarthritis: a systematic review and meta-analysis. *Arthritis care & research*, 65(10):1643–1665, 2013.
- [4] Silvia Del Din, Alan Godfrey, Brook Galna, Sue Lord, and Lynn Rochester. Free-living gait characteristics in ageing and parkinson’s disease: impact of environment and ambulatory bout length. *Journal of neuroengineering and rehabilitation*, 13(1):46, 2016.
- [5] Richard Baker. The history of gait analysis before the advent of modern computers. *Gait & posture*, 26(3):331–342, 2007.
- [6] Martina Mancini and Fay B Horak. The relevance of clinical balance assessment tools to differentiate balance deficits. *European journal of physical and rehabilitation medicine*, 46(2):239, 2010.
- [7] Weijun Tao, Tao Liu, Rencheng Zheng, and Hutian Feng. Gait analysis using wearable sensors. *Sensors*, 12(2):2255–2283, 2012.
- [8] ACRMDOG Godfrey, Richard Conway, David Meagher, and Gearoid ÓLaighin. Direct measurement of human movement by accelerometry. *Medical engineering & physics*, 30(10):1364–1386, 2008.
- [9] Andrea Coravos, Sean Khozin, and Kenneth D Mandl. Developing and adopting safe and effective digital biomarkers to improve patient outcomes. *NPJ digital medicine*, 2(1):14, 2019.
- [10] Cyril Voisard, Rémi Barrois, Nicolas de l’Escalopier, Nicolas Vayatis, Pierre-Paul Vidal, Alain Yelnik, Damien Ricard, and Laurent Oudre. A dataset of clinical gait signals with wearable sensors from healthy, neurological, and orthopedic cohorts. *Scientific Data*, 12(1):1674, 2025.
- [11] Carlotta Caramia, Diego Torricelli, Maurizio Schmid, Adriana Muñoz-Gonzalez, Jose Gonzalez-Vargas, Francisco Grandas, and Jose L Pons. Imu-based classification of parkinson’s disease from gait: A sensitivity analysis on sensor location and feature selection. *IEEE journal of biomedical and health informatics*, 22(6):1765–1774, 2018.
- [12] Fabio A Storm, Christopher J Buckley, and Claudia Mazzà. Gait event detection in laboratory and real life settings: Accuracy of ankle and waist sensor based methods. *Gait & posture*, 50:42–46, 2016.
- [13] Matthias Van Parijs and Leen Boonen. Overview of measurement methods of inter-limb coordination. 2019.
- [14] Nils Y Hammerla, Shane Halloran, and Thomas Plötz. Deep, convolutional, and recurrent models for human activity recognition using wearables. *arXiv preprint arXiv:1604.08880*, 2016.
- [15] Francisco Javier Ordóñez and Daniel Roggen. Deep convolutional and lstm recurrent neural networks for multimodal wearable activity recognition. *Sensors*, 16(1):115, 2016.
- [16] Brett M Meyer, Lindsey J Tulipani, Reed D Gurchiek, Dakota A Allen, Lukas Adamowicz, Dale Larie, Andrew J Solomon, Nick Cheney, and Ryan S McGinnis. Wearables and deep learning classify fall risk from gait in multiple sclerosis. *IEEE journal of biomedical and health informatics*, 25(5):1824–1831, 2020.
- [17] Dzmitry Bahdanau. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [18] Haojie Ma, Wenzhong Li, Xiao Zhang, Songcheng Gao, and Sanglu Lu. Attnsense: Multi-level attention mechanism for multimodal human activity recognition. In *IJCAI*, pages 3109–3115, 2019.
- [19] Sarthak Jain and Byron C Wallace. Attention is not explanation. *arXiv preprint arXiv:1902.10186*, 2019.
- [20] Sofia Serrano and Noah A Smith. Is attention interpretable? *arXiv preprint arXiv:1906.03731*, 2019.
- [21] Sarah Wiegreffe and Yuval Pinter. Attention is not not explanation. *arXiv preprint arXiv:1908.04626*, 2019.
- [22] Shashank Yadav and Vignesh Subbian. When attention fails: Pitfalls of attention-based model interpretability for high-dimensional clinical time-series. In *Conference on Health, Inference, and Learning*, pages 289–305. PMLR, 2025.
- [23] Ziad Obermeyer, Brian Powers, Christine Vogeli, and Sendhil Mullainathan. Dissecting racial bias in an algorithm used to manage the health of populations. *Science*, 366(6464):447–453, 2019.

- [24] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020.
- [25] John R Zech, Marcus A Badgeley, Manway Liu, Anthony B Costa, Joseph J Titano, and Eric Karl Oermann. Variable generalization performance of a deep learning model to detect pneumonia in chest radiographs: a cross-sectional study. *PLoS medicine*, 15(11):e1002683, 2018.
- [26] Effy Vayena, Alessandro Blasimme, and I Glenn Cohen. Machine learning in medicine: addressing ethical challenges. *PLoS medicine*, 15(11):e1002689, 2018.
- [27] Alvin Rajkomar, Jeffrey Dean, and Isaac Kohane. Machine learning in medicine. *New England Journal of Medicine*, 380(14):1347–1358, 2019.
- [28] Christopher J Kelly, Alan Karthikesalingam, Mustafa Suleyman, Greg Corrado, and Dominic King. Key challenges for delivering clinical impact with artificial intelligence. *BMC medicine*, 17(1):195, 2019.
- [29] Michael Roberts, Derek Driggs, Matthew Thorpe, Julian Gilbey, Michael Yeung, Stephan Ursprung, Angelica I Aviles-Rivero, Christian Etzmann, Cathal McCague, Lucian Beer, et al. Common pitfalls and recommendations for using machine learning to detect and prognosticate for covid-19 using chest radiographs and ct scans. *Nature Machine Intelligence*, 3(3):199–217, 2021.
- [30] Kaustubh Chakradeo, Inchuen Huynh, Sedrah B Balaganeshan, Ole L Dollerup, Hjørdis Gade-Jørgensen, Susanne K Laupstad, Mikkel Malham, Tri-Long Nguyen, Adam Hulman, and Tibor V Varga. Navigating fairness aspects of clinical prediction models. *BMC medicine*, 23(1):567, 2025.
- [31] Benoit Mariani, Hossein Rouhani, Xavier Crevoisier, and Kamiar Aminian. Quantitative estimation of foot-flat and stance phase of gait using foot-worn inertial sensors. *Gait & posture*, 37(2):229–234, 2013.
- [32] Shachar Kaufman, Saharon Rosset, Claudia Perlich, and Ori Stitelman. Leakage in data mining: Formulation, detection, and avoidance. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 6(4):1–21, 2012.
- [33] George Zerveas, Srideepika Jayaraman, Dhaval Patel, Anuradha Bhamidipaty, and Carsten Eickhoff. A transformer-based framework for multivariate time series representation learning. In *Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining*, pages 2114–2124, 2021.
- [34] Yanbo Xu, Shangqing Xu, Manav Ramprasad, Alexey Tumanov, and Chao Zhang. Transehr: self-supervised transformer for clinical time series data. In *Machine Learning for Health (ML4H)*, pages 623–635. PMLR, 2023.
- [35] Qingyu Zhao, Ehsan Adeli, and Kilian M Pohl. Training confounder-free deep learning models for medical applications. *Nature communications*, 11(1):6010, 2020.