

Multimodal Indoor Localisation for Measuring Mobility in Parkinson’s Disease using Transformers

Ferdian Jovan¹, Ryan McConville¹, Catherine Morgan¹, Emma Tonkin¹, Alan Whone¹, and Ian Craddock¹

University of Bristol, BS8 1TH, United Kingdom
ferdian.jovan@bristol.ac.uk

Abstract. Parkinson’s disease (PD) is a slowly progressive debilitating neurodegenerative disease which is prominently characterised by motor symptoms (e.g. bradykinesia, rigidity and tremor). Indoor localisation, including number and speed of room to room transitions, provides a proxy outcome which represents mobility and could be used as a digital biomarker to quantify how mobility changes as this disease progresses. We use data collected from 10 people with Parkinson’s, and 10 controls, each of whom lived for five days in a smart home with various sensors. In order to more effectively localise them indoors, we propose a transformer-based approach utilizing two data modalities, Received Signal Strength Indicator (RSSI) and accelerometer data from wearable devices, which provide complementary views of movement. Our approach makes asymmetric and dynamic correlations by a) learning temporal correlations at different scales and levels, and b) utilizing various gating mechanisms to select relevant features within modality and suppress unnecessary modalities. On a dataset with real patients, we demonstrate that our proposed method gives an average accuracy of 89.9%, outperforming competitors. We also show that our model is able to better predict in-home mobility for people with Parkinson’s with an average offset of 1.13 seconds to ground truth.

Keywords: Parkinson’s disease · Indoor localisation · Multimodal data · Transformer · Time series.

1 Introduction

Parkinson’s disease (PD) is a chronic, debilitating neurodegenerative disease which is characterised by a variety of motor symptoms, such as slowness of movement, rigidity, tremor, and gait dysfunction [11]. It is a slowly progressive disease - however, for each individual, the symptoms can fluctuate hour-by-hour, related to medication intake timing, stress and other factors. A particular problem is that patients “wear off” from their medications and experience a worsening of symptoms prior to the next medication dose. Very frequent symptom evaluation is needed to give an accurate evaluation of how severe the symptoms and

their fluctuations are for an individual; additionally, constant monitoring could also capture the range of symptoms over time so that the slow disease progression can be sensitively detected and quantified. For example, as PD progresses, motor symptoms become more severe which hinder the subject’s gait and movement around their own house. As a result, the subject is more likely to stay in one room; once they move, they typically need more time to transition between rooms. Up until now, a PD evaluation has often been done in an artificial environment like a clinic or laboratory where only a snapshot view of the individual’s motor function and mobility can be captured. Given the known symptom fluctuations experienced by most patients, measuring a snapshot of the function and mobility may lead to an incorrect conclusion about individual’s PD progression.

Enabled by an IoT-based platform with multimodal devices designed to allow continuous, unobtrusive sensing, we make progress towards autonomous PD evaluation and monitoring in home environments by providing continuous indoor localisation of people with PD. Indoor localisation, including the number and speed of room-to-room transitions can enhance snapshot clinical assessments by objectively and unobtrusively capturing real-world function and behaviour, and it could be used to monitor PD progression, as a proxy outcome digital biomarker, to quantify how mobility changes as the disease progresses. Specifically, use wearable inertial measurement unit (IMU) sensors to collect Received Signal Strength Indicator (RSSI) and accelerometer data from PD and healthy control (HC) subjects living daily life in a home environment. Although both RSSI and accelerometer data come from wearable devices, those data represent different contexts. RSSI data are typically used for estimating location, while accelerometer data can be used to differentiate activities. As some activities are specific to particular locations or rooms (e.g. stirring a pan on the hob must be in a kitchen), accelerometer data may complement RSSI in separating adjacent rooms, which RSSI alone may struggle with.

Using both RSSI and accelerometer data, we propose a deep learning approach with dual modalities that encode temporal room signatures to perform indoor localisation, in particular room-level classification. We cover one challenge, particularly for PD, that is faced by any machine learning technique with several modalities. As PD is a heterogeneous disease, the severity of symptoms varies from one patient to another [6]. Severe symptoms, such as tremor, may affect the generalisation of accelerometer data and combining it with RSSI data may, in fact, worsen the performance of indoor localisation. Furthermore, the challenge is magnified by the free-living environment where the movements and mobility are greatly varied and unstructured. Our proposed architecture, Dual Context Modality Network (DCMN) for room-level classification, is based on the fusion of several layer neural networks that are designed to (1) capture temporal room signatures both local (i.e. within few time steps) and global (i.e. across all time steps) and (2) adaptively choose features and modality based on their importance. We show that such fusion of various layers helps in dealing with the challenges mentioned, as DCMN intelligently encodes the temporal room signatures while adaptively deciding whether some inputs (or modality) have discrim-

inative information in both modalities. Our evaluation on our PD dataset, which includes subjects with and without PD, shows that DCMN achieves the state-of-the-art accuracy for indoor localisation. We also demonstrate the effectiveness of DCMN in predicting in-home mobility (i.e. number of daily transition, room-to-room transition duration) for PD subjects with an average offset closest to ground truth.

2 Related Work

Advancement in machine learning has motivated research on automatic PD monitoring and evaluation. Early research started with simple PD classification [5] or easy-to-distinguish symptoms identification [1,4]. Much of this research relies on accelerometer data from smart phones or wearable devices as their main data source. Alternatively, there are some other methods using vision sensors [13]. Although raw data can be used for a simple PD classification, some researches do feature extraction before applying any classification methods. For example, in [4], restricted Boltzmann machines are trained using features extracted from wrist-worn accelerometer data in a home environment to predict PD state. Similarly, [20] use convolutional neural networks (CNN) on augmented accelerometer data to classify PD motor state. Li et al. [13] use CNNs on RGB data to first estimate human pose and then extract features from trajectories of joints movements. RF is finally used to classify PD vs. non-PD symptoms and measure their severity.

While much of the work in PD evaluation report high performance for their learning methods, they tend to use single modality for homogeneity, which raises the question of whether additional modalities can further improve performance. There is limited research in PD utilizing multiple sensors for a better prediction and evaluation. In other healthcare applications, several works have started using multiple modalities to improve their performance and robustness. [16] proposes a network, called CaloriNet, for fusing accelerometer and silhouette data to estimate the calorie expenditure of subjects. Heidarvinchek et al. [7] proposed MCPD-Net deep network that uses two data modalities, acquired from vision and accelerometer sensors in a home environment for a PD classification. They minimise the difference between the latent spaces corresponding to the two data modalities before the final representation is fed up to a linear layer for PD classification. Masullo et al. [17] match video sequences of silhouettes to accelerations from wearable sensors for a person re-identification in a home environment. Their application is used to identify members of a household while respecting their privacy. All of this research uses vision as their main data source; while vision has proved to be a powerful modality for PD monitoring, privacy issues in home settings has limited research on RGB data.

A promising direction for more privacy-friendly PD monitoring in home environments is via indoor localisation. Indoor localisation typically uses fingerprinting to collect data to train a machine learning model. Fingerprinting uses either classification or regression methods to estimate the location of wearable devices

by exploiting signal sources present in the environment. Utilising RSSI effectively is challenging. A significant challenge is due to the random fluctuations in RSSI values due to shadowing, fading and multi-path effects. However, in recent years, many techniques have been proposed to tackle the RSSI fluctuations and, indirectly, improve the localisation accuracy. Zhang et al. in [27] proposed a 4-layer deep neural network (DNN) that generates coarse positioning estimates, which is then refined to produce a final location estimate by a hidden Markov model (HMM). To further improve location accuracy for different buildings and floors, Ibrahim et al. in [9] exploit a time-series of RSSI data from access points (AP) to estimate room locations. A CNN is used to build localisation models to further leverage the temporal dependencies across time-series readings.

Even though RSSI data has shown a promising result for indoor localisation, relying on RSSI data alone is not enough to tackle home environments for PD subjects due to shadowing and rooms with tight separation. As we aim to measure the PD progression through in-home mobility within naturalistic home environments, we propose to use multiple modalities, i.e. RSSI (which can estimate location) and accelerometers (which can measure movement), to expand our input domain and capture a wider range of features, and, in turn, improve the localisation accuracy.

3 Proposed Framework: DCMN

We introduce Dual Context Modality Network (DCMN) for indoor localisation, a novel method for improving the accuracy of room-level classification by merging various modalities and learning their temporal room signatures. There are several challenges that arise from it:

1. **Considering multivariate features.** RSSI signals are commonly used for room level localisation [9]. These signals, that are measured at multiple access points within a home, are transmitted by a wearable. However, most previous approaches have not considered the use of additional features to more effectively enrich the RSSI signals [10,19]. Naturally, the wearable also produces acceleration measurements, and thus we can explore if accelerometer data will enrich the RSSI signals, in particular to help distinguish adjacent rooms, which RSSI only systems typically struggle with. If it will, how can we incorporate these extra features into the existing features for accurate room predictions, particularly within the context of PD where the acceleration signal may be significantly impacted by PD itself.
2. **Capturing feature and multimodal importance.** RSSI signals have been widely used for indoor localisation, typically using a fingerprinting technique that produces a ground truth radio map of a home. Similar to RSSI signals, accelerometers data can be used to identify typical activities performed in a specific room, and in turn, help identify in which room a person is. However, identifying which of, and when, these features become important is a challenging problem. Can we identify which features are important?

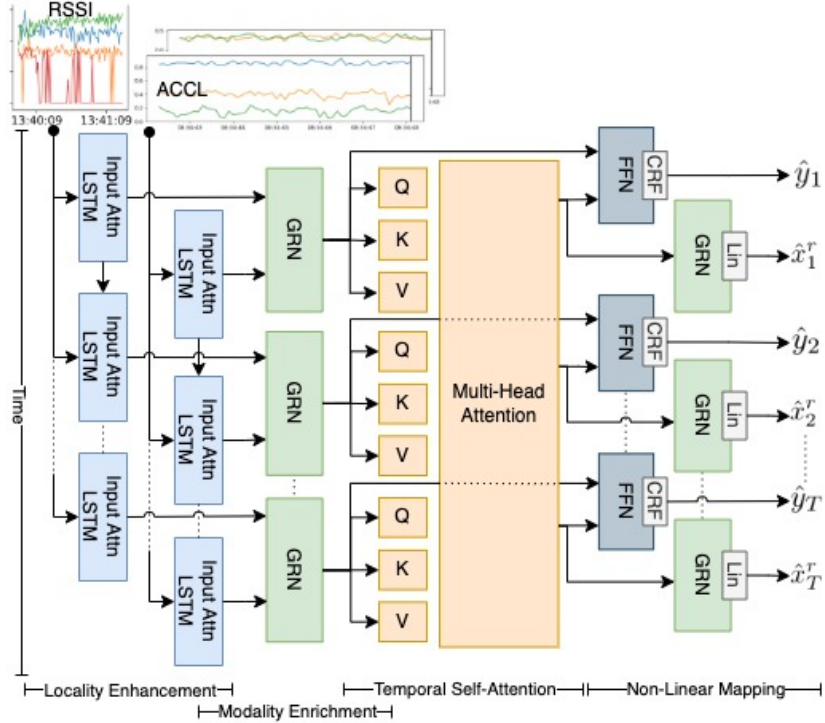


Fig. 1: DCMN architecture. DCMN processes time-varying RSSI and accelerometer features as two separate modalities. Four main neural network layers are presented to tackle challenges in merging two modalities.

3. **Modelling local and global temporal dynamics.** The true correlations between inputs both intra-modality (i.e. RSSI signal among access points) and inter-modality (i.e. RSSI signal against accelerometer fluctuation) are dynamics. These dynamics can effect one another within a local context (e.g. cyclical patterns) or across long-term relationships. Can we capture local and global relationship across different modalities?

The DCMN architecture, shown in Figure 1, addresses the aforementioned challenges through a fusion of four neural network layers which are described in the following sections.

3.1 Locality Enhancement and Internal Feature Selection with Attention LSTM

We consider the use of a long short-term memory (LSTM) network that naturally captures local patterns and has an appropriate inductive bias for the time ordering of the inputs. Given feature vectors $\mathbf{x}_t^u = [x_t^1, \dots, x_t^u]$ with $u \in \{r, a\}$

representing RSSI (i.e. access points) and accelerometer (i.e. spatial direction) features, and $t \leq T$ representing time index, we aim to learn a summarised temporal embedding \mathbf{h}_t^u for each modality u at each time step t ¹. A LSTM enables this learning, a mapping from \mathbf{x}_t to \mathbf{h}_t , by updating hidden state vectors through time with

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

where $\mathbf{h}_t \in \mathbb{R}^d$ is the temporal embedding of the LSTM at time t , and d is the embedding dimension (common across DCMN).

LSTM with an input attention. Instead of using a standard LSTM, we adopt an attention mechanism to be the input of an LSTM. The network is inspired by [21]. An input attention LSTM can adaptively select the relevant feature for each modality at each time step. Given the k -th input time series $\mathbf{x}_k = (x_1^k, \dots, x_T^k) \in \mathbb{R}^T$, an input attention is constructed via a multilayer perceptron utilizing the previous hidden state \mathbf{h}_{t-1} of an LSTM with:

$$e_t^k = \mathbf{v}_e \tanh(\mathbf{W}_e \mathbf{h}_{t-1} + \mathbf{U}_e \mathbf{x}^k) \quad (2)$$

and

$$\alpha_t^k = \frac{\exp(e_t^k)}{\sum_i^u \exp(e_t^i)} \quad (3)$$

where $\mathbf{v}_e \in \mathbb{R}^T$, $\mathbf{W}_e \in \mathbb{R}^{T \times T}$, and, $\mathbf{U}_e \in \mathbb{R}^{T \times T}$ are parameters to learn, and $u \in \{r, a\}$ is the type of modality (RSSI, or accelerometer). Bias terms are omitted for clarity. α_t^k is the attention weight measuring the importance of the k -th feature at time t . With these attention weights, the feature series can be adaptively adjusted with:

$$\hat{\mathbf{x}}_t = (\alpha_t^1 x_t^1, \alpha_t^2 x_t^2, \dots, \alpha_t^u x_t^u) \quad (4)$$

and Equation 1 can be computed accordingly by replacing \mathbf{x}_t with $\hat{\mathbf{x}}_t$. For regularisation, Dropout is added after Equation 1.

3.2 Modality Enrichment Layer with GRN

Unlike other approaches in indoor localisation which typically utilise one modality (RSSI), here we propose to add accelerometer data as an additional modality to enrich the RSSI data as a temporal embedding. To achieve our target, we use Gated Residual Network (GRN), introduced by [14], to integrate dual inputs into one integrated embedding. The GRN takes in a (primary) input $\mathbf{x} \in \mathbb{R}^d$ and another (secondary) input $\mathbf{y} \in \mathbb{R}^d$ and yields:

$$GRN(\mathbf{x}, \mathbf{y}) = LayerNorm(x + GLU(\Xi_1(\mathbf{x}, \mathbf{y}))), \quad (5)$$

with

$$\Xi_1(\mathbf{x}, \mathbf{y}) = \mathbf{W}_1 \Xi_2(\mathbf{x}, \mathbf{y}) + \mathbf{b}_1 \quad (6)$$

$$\Xi_2(\mathbf{x}, \mathbf{y}) = ELU(\mathbf{W}_2 \mathbf{x} + \mathbf{W}_3 \mathbf{y} + \mathbf{b}_2) \quad (7)$$

¹ We omit the modality symbol u for simplicity.

where ELU is the Exponential Linear Unit activation function [2], *GLU* is the Gating Linear Unit function [3], *LayerNorm* is standard layer normalization, Ξ_1, Ξ_2 are intermediate layers, and $\mathbf{W}_1, \mathbf{W}_2 \in \mathbb{R}^{d \times d}$, and $\mathbf{b}_1, \mathbf{b}_2 \in \mathbb{R}^d$ are weights and biases to learn. A GLU is used here to provide the flexibility to suppress the enrichment from the secondary input if it is not needed.

The temporal embedding for RSSI \mathbf{h}_t^r and accelerometer \mathbf{h}_t^a produced by Equation 1 are then fed to a *GRN* network:

$$\hat{\mathbf{h}}_t = GRN(\mathbf{h}_t^r, \mathbf{h}_t^a) \quad (8)$$

with the aim for the accelerometer data to enrich the RSSI temporal embedding with information regarding particular movements specific to each room. Note that this layer has all weights shared across each time step t . Dropout is also added before the GLU network.

3.3 Temporal Self-Attention Layer with Transformer Encoder

The DCMN employs a self-attention mechanism within the transformer [23] to pick-up global dependencies that may be challenging for RNN-based architectures to capture. This works in the opposite way when computing the temporal embedding \mathbf{h}_t by Attention LSTM, as it focuses on the local patterns surrounding each time step. Furthermore, the transformer allows asymmetric long-term learning by learning different query and key weight matrices, reflecting the information diffusion in the aggregate temporal embedding $\hat{\mathbf{h}}_t$.

Multihead Self-Attention. In general, attention mechanisms scale values $\mathbf{V} \in \mathbb{R}^{T \times d_v}$ based on relationships between keys $\mathbf{K} \in \mathbb{R}^{T \times d_{attn}}$ and queries $\mathbf{Q} \in \mathbb{R}^{T \times d_{attn}}$ as below:

$$Attention(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = Softmax(\mathbf{Q}\mathbf{K}^T / \sqrt{d_{attn}}) \mathbf{V}, \quad (9)$$

where *Softmax* is chosen to be the normalization function with a scaled dot-product attention as in [23].

To improve the learning capacity of the standard attention mechanism, multi-head attention is proposed in [23], employing different attention heads with different sets of $\mathbf{Q}, \mathbf{K}, \mathbf{V}$:

$$MultiHead(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = [H_1, \dots, H_m] W_H, \quad (10)$$

$$H_h = Attention(\mathbf{Q} \mathbf{W}_Q^h, \mathbf{K} \mathbf{W}_K^h, \mathbf{V} \mathbf{W}_V^h), \quad (11)$$

where $\mathbf{W}_K^h \in \mathbb{R}^{d \times d_{attn}}$, $\mathbf{W}_Q^h \in \mathbb{R}^{d \times d_{attn}}$, $\mathbf{W}_V^h \in \mathbb{R}^{d \times d_v}$ are head-specific weights for keys, queries and values, and $W_H \in \mathbb{R}^{(m \cdot d_v) \times d}$ linearly combines outputs concatenated from all heads H_h .

Following the feature and modality enhancement layers, we next apply multi-head self-attention. All enriched temporal embeddings $\hat{\mathbf{h}}_t$ are first grouped into

a single matrix – i.e. $\mathbf{h} = [\hat{\mathbf{h}}_1, \dots, \hat{\mathbf{h}}_T]$ – and multi-head attention is applied at each time step $t, t \leq T$:

$$\tilde{\mathbf{h}} = \text{MultiHead}(\mathbf{h}, \mathbf{h}, \mathbf{h}), \quad (12)$$

to yield $\tilde{\mathbf{h}} = [\tilde{\mathbf{h}}_1, \dots, \tilde{\mathbf{h}}_T]$. We choose m attention heads such that $d_V = d_{attn} = d/m$ with d as the embedding dimension.

3.4 Non-linear Mapping Layer through CRF and GRN

We apply additional non-linear processing to the outputs $\tilde{\mathbf{h}}$ of the self-attention layer. Similar to the transformer encoder in [23]. This makes use of a multilayer perceptron (MLP) in combination with skip connections as follows:

$$\check{\mathbf{h}}_t = \tanh\left(\Xi_3(\hat{\mathbf{h}}_t + \tilde{\mathbf{h}}_t) + \text{MLP}\left(\Xi_3(\hat{\mathbf{h}}_t + \tilde{\mathbf{h}}_t)\right)\right), \quad (13)$$

$$\Xi_3(\mathbf{x}) = \text{LayerNorm}(\mathbf{x}) \quad (14)$$

$$\text{MLP}(\mathbf{x}) = \mathbf{W}_1 \delta(\mathbf{W}_2 \mathbf{x} + \mathbf{b}_2) + \mathbf{b}_1 \quad (15)$$

where $\delta(\cdot)$ is the Mish activation function [18], and $\mathbf{W}_1 \in \mathbb{R}^{d \times (4 \cdot d)}$, $\mathbf{W}_2 \in \mathbb{R}^{(4 \cdot d) \times d}$, and $\mathbf{b}_1 \in \mathbb{R}^d$, and $\mathbf{b}_2 \in \mathbb{R}^{4 \cdot d}$. The *MLP*, transforming the size of the temporal embedding $(\hat{\mathbf{h}}_t + \tilde{\mathbf{h}}_t)$ to 4 times of its original size d with Mish activation function, is done to refine the aggregate temporal embeddings since the self-attention does not impose additional non-linearity. It also contains two residual connections to learn the identity function if needed: one for the self-attention and another for the MLP. The MLP layer has all weights shared across each time step t . Dropout is applied after the attention layer.

Final Prediction. Finally, we apply two different layers to produce two different outputs. The room-level predictions is produced via a single conditional random field (CRF) layer in combination with a linear layer applied to the refined temporal embeddings to produce the final predictions as

$$\hat{y}_t = \text{CRF}\left(\Xi_4(\check{\mathbf{h}}_t)\right) \quad (16)$$

$$\Xi_4(\mathbf{x}) = \mathbf{W}_p \mathbf{x} + \mathbf{b}_p \quad (17)$$

where $\mathbf{W}_p \in \mathbb{R}^{d \times n}$, and $\mathbf{b}_p \in \mathbb{R}^n$ are weight and bias to learn. Even though transformer can take into account neighbour information before generating its own temporal embedding at time step t , its decision is independent; it does not take into account the actual decision made by other temporal embeddings t . We use a CRF layer to cover just that to maximize the probability of the temporal embeddings of the entire time steps, so it can better model cases where temporal embeddings closest to one another must be compatible (i.e. minimizing the possibility for impossible room transitions). When finding the best sequence of room location \hat{y}_t , the Viterbi Algorithm is used as a standard for CRF layer.

The RSSI value reconstruction is produced via a GRN network applied to the refined temporal embeddings as

$$\hat{x}_t^r = \mathbf{W}_r \text{GRN}(\tilde{\mathbf{h}}_t) + \mathbf{b}_r \quad (18)$$

where $\mathbf{W}_r \in \mathbb{R}^{d \times r}$, and $\mathbf{b}_r \in \mathbb{R}^r$ are weight and bias to learn. The RSSI reconstruction is used for backloss regularisation which is explained in detail in the next section.

3.5 Training with Backloss Regularisation

During the training process, the DCMN produces two kinds of outputs. Emission outputs (outputs produced by Equation 17 prior to prediction outputs) $\hat{\mathbf{e}} = [\Xi_4(\tilde{\mathbf{h}}_1), \dots, \Xi_4(\tilde{\mathbf{h}}_T)]$ are trained to generate the likelihood estimate of room predictions, while the backcasting outputs $\hat{\mathbf{x}}^r = [\hat{x}_1^r, \dots, \hat{x}_T^r]$ are used in an auto-encoding fashion to enhance the representation power of RSSI data. The final loss function can be formulated as a combination of both likelihood and backcasting losses:

$$\mathcal{L}(\hat{\mathbf{e}}, \mathbf{y}, \hat{\mathbf{x}}^r, \mathbf{x}^r) = \mathcal{L}_{NLL}(\hat{\mathbf{e}}, \mathbf{y}) + \sum_{i=1}^T \mathcal{L}_H(\hat{x}_i^r, x_i^r) \quad (19)$$

$$\mathcal{L}_{NLL}(\hat{\mathbf{e}}, \mathbf{y}) = \sum_{\hat{y}} \sum_{i=0}^T P(\Xi_4(\tilde{\mathbf{h}}_i) | \hat{y}_i) T(\hat{y}_i | \hat{y}_{i-1}) - \sum_{i=0}^T P(\Xi_4(\tilde{\mathbf{h}}_i) | y_i) T(y_i | y_{i-1}) \quad (20)$$

$$\mathcal{L}_H(\hat{y}, y) = \begin{cases} 0.5(\hat{y} - y)^2 & |\hat{y} - y| < \tau \\ \tau(|\hat{y} - y| - 0.5\tau), & \text{otherwise} \end{cases} \quad (21)$$

where $\mathcal{L}_{NLL}(\cdot)$ represents the negative log-likelihood and $\mathcal{L}_H(\cdot)$ denotes the backcasting loss, $\mathbf{y} = [y_1, \dots, y_T] \in \mathbb{R}^T$ is the actual room locations, $\mathbf{x}^r = [x_1^r, \dots, x_T^r] \in \mathbb{R}^{T \times r}$ is the actual RSSI signals across r access points, and τ is a hyper-parameter. $P(x | y)$ denotes the conditional probability, and $T(y_i | y_{i-1})$ denotes the transition matrix cost of having transitioned from y_{i-1} to y_i . Huber loss [8], which is widely used for regression problems, is used as the backcasting loss to alleviate the impact of the outliers in comparison with square loss.

4 Experiments

4.1 Dataset

This dataset was collected using privacy preserving RGB-D cameras, and wearable sensors in a residential smart home. For the data collection, 10 Access Points (APs) were installed throughout the home, which measure the RSSI (Received Signal Strength Indication) [12]. Participants wore wristband devices,

one on each hand, equipped with a tri-axial accelerometer. The devices wirelessly transmit data using the Bluetooth Low Energy (BLE) standard to the 10 APs. The outputs of these wearable sensors are a continuous numerical stream of the accelerometer readings and RSSI values which were both sampled at 20 Hz. To measure the accuracy of our proposed network, cameras are installed in the ground floor of the house as ground truth. These cover the kitchen, hallway, dining room, and living room. Due to privacy requirements, the RGB and depth data were discarded after extracting the room location information, and the data were only collected during the day for 2-3 hours daily. As a result, we perform indoor localisation in the following rooms: kitchen, living room, dining room, hallway, stairs, and porch. Note, our approach generalises to all rooms over the entire period, however we limit the time frame and rooms to only those we have a ground truth for.

Our dataset includes RSSI and accelerometer data corresponding to 10 heterosexual pairs living freely in a smart home for five days. Each pair consists of one person with PD and one person as the HC. From the 20 participants, 4 females and 6 males have PD. The average age of the participants is 63 and the average time since PD diagnosis for the person with PD is 8.9 years. The duration of data recorded by the RGB-D cameras for PD and HC is 53.8 and 52.8 hours, respectively (106.6 hours in total), which provides a relatively balanced label set for our room-level classification.

Data pre-processing. Two wearable sensor values are grouped together based on their modalities, i.e. twenty RSSI values corresponding to 10 APs for each wearable sensor, and six spatial directions corresponding to three spatial directions (x, y, z) for each wearable, at each time. This data is resampled to 1 Hz with a 10-second time window, which makes an input of size (10 x 20) for RSSI data and an input of size (10 x 6) for accelerometer data. Imputation for missing values, specifically for RSSI data, is applied by replacing the missing values with the lowest acceptable value (i.e. -120dB). Missing values exist in RSSI data whenever the wearable is out of range of an AP. Finally, all time-series measurements by the modalities are normalized to be within the range of zero and one before they are processed by the model.

4.2 Experimental Setup

Baseline. We compare DCMN with the following baselines for indoor localisation:

- Random Forest (RF) is the most basic baseline for our indoor localisation, where all time series features of RSSI and accelerometer are flattened and merged into one long feature vector for room prediction.
- DARNN [21] represents the attention LSTM that correlates multiple inputs and time steps using a dual-attention mechanism. For a representative comparison, each modality is represented by one DARNN network, where a simple MLP layer is used on top of them to merge the two networks into one output.

- DTML [25] represents the state-of-the-art model for multimodal and multivariate time series with a transformer encoder to learn asymmetric correlations across modalities.
- DeIT [22] represents a state-of-the-art pure transformer encoder for visual learning tasks which is combined with a distillation technique that learns from a teacher network to improve performance. We choose DTML as a teacher network as DTML is the closest model with natural multimodal capabilities similar to our DCMN.
- TENER [24] is a modified transformer encoder with direction and distance-aware attention in combination with conditional random field (CRF) to further enforce dependencies across temporal aspects.

For both DeIT and TENER, at each time step t , RSSI \mathbf{x}_t^r and accelerometer \mathbf{x}_t^a features are combined via a linear layer before they are processed by the networks. Nvidia Quadro RTX 6000 GPU was used for these experiments.

Hyperparameters. We hyperparameter tune DCMN as follows: the embedding dimension d in $\{128, 256\}$, the number of epochs in $\{200, 300\}$, and the learning rate in $\{0.01, 0.0001\}$. We set the dropout rate to 0.15. We use the RAdam optimizer [15] in combination with Look-Ahead algorithm [26] for the training with early stopping using the validation performance. Similar settings are used for other neural network baseline models. In our RF models, we perform a cross-validated parameter search for the number of trees ($\{200, 250\}$) and the minimum number of samples in a leaf node ($\{1, 5\}$). The Gini impurity is used to measure an optimal split.

4.3 Experimental Results

Given our particular interest in developing a system to better understand PD in home environments, we specifically design experiments in order to be better measure the performance. For example, we will consider if we can detect any significant difference in the performance of the systems when trained on a person with Parkinson’s versus trained on someone without. This may provide useful insight into the deployment of systems in the future, such as whether there is any benefit to a person with Parkinson’s collecting the training data over a HC.

Apart from training our models on all HC subjects (ALL-HC), we also perform two different kinds of cross-validation: 1) We leave one PD subject out as training data (LOO-PD), 2) we leave one HC subject out as training data (LOO-HC). For all of our experiments, we test our trained models on all PD subjects (excluding the one used as training data for LOO-PD). For prediction accuracy, we report our classification results by precision, accuracy, and F1-score, all averaged and standard deviated across the test folds.

Prediction Accuracy. Table 1 compares the accuracy of our DCMN and baselines for room-level classification in PD datasets. DCMN outperforms all baselines with consistent improvements in many cross validation types. The improvement is more significant on the LOO-PD validation, where the training

Table 1: Room-level classification accuracy of our DCMN and other baselines. Standard deviation is shown under (.), the best performer is bold, while the second best is italicized.

Training Data	Models	Precision	Accuracy	F1-Score
ALL-HC	RF	96.10	94.60	95.30
	DARNN	95.50	93.70	94.50
	DTML	95.20	93.50	94.30
	DeIT	93.80	93.80	93.80
	TENER	95.30	93.40	94.20
	DCMN	<i>95.60</i>	<i>93.90</i>	<i>94.70</i>
LOO-HC	RF	89.44 (7.19)	89.99 (4.36)	89.39 (6.03)
	DARNN	89.09 (6.59)	89.42 (4.52)	88.33 (6.65)
	DTML	89.95 (6.46)	<i>90.55 (3.35)</i>	<i>90.01 (5.14)</i>
	DeIT	88.14 (6.41)	88.38 (4.38)	86.99 (6.60)
	TENER	90.66 (2.62)	89.65 (3.59)	89.02 (5.20)
	DCMN	<i>90.08 (6.60)</i>	90.84 (3.23)	90.28 (5.15)
LOO-PD	RF	88.07 (8.41)	88.25 (6.02)	86.89 (8.61)
	DARNN	86.08 (7.99)	85.67 (6.21)	84.79 (4.73)
	DTML	87.16 (8.16)	87.41 (5.98)	86.51 (7.53)
	DeIT	83.66 (9.17)	83.50 (6.62)	81.13 (8.90)
	TENER	<i>89.29 (2.97)</i>	<i>88.31 (3.74)</i>	<i>87.79 (4.73)</i>
	DCMN	89.81 (2.74)	88.98 (3.53)	88.78 (3.67)

data for the accelerometer is more prone to variation depending on the severity of the disease.

The high accuracy of DTML, especially for the LOO-HC validation, is because of its ability to model the temporal dynamics of each modality and the ability to capture asynchronous relation across modalities. However, the accuracy suffers in the LOO-PD validation as the accelerometer data (and modality) may be erratic due to PD and should be excluded at times from contributing to room classification. DCMN achieves the same objective by correlating different modalities via GRN in combination with transformer encoder. However, for the LOO-PD validation, the DCMN performs very well due to its ability to suppress a noisy modality, i.e. noisy accelerometer data. Any model, that is able to suppress the accelerometer information as the DCMN does, might have performed well on the LOO-PD validation.

Room-to-Room Transition Accuracy. We also compare the performance of our proposed architecture in terms of in-home mobility as shown in Table 2. We measure in-home mobility by ‘Daily Transitions’ and room-to-room transition duration. ‘Daily Transitions’ shows the average number of transitions daily

Table 2: Room-to-room transition accuracy of our DCMN and other baselines.

Data	Models	Daily Transition	Kitchen-Live	Kitchen-Dine	Dine-Live
Ground Truth		14.87 (10.59)	11.02 (17.45)	11.12 (11.77)	7.04 (5.59)
ALL-HC	RF	19.90 (41.4)	11.80 (12.26)	11.48 (11.36)	9.77 (8.29)
	DARNN	20.70 (27.20)	9.81 (8.29)	13.12 (12.46)	9.86 (8.44)
	DTML	25.62 (46.57)	<i>10.33 (9.21)</i>	13.52 (11.92)	12.32 (10.91)
	DeIT	<i>18.96 (26.49)</i>	10.20 (9.49)	7.86 (5.45)	<i>8.79 (6.52)</i>
	TENER	21.26 (25.49)	9.88 (9.51)	<i>11.53 (11.06)</i>	10.45 (7.39)
	DCMN	17.81 (23.25)	10.43 (9.50)	11.62 (9.62)	8.67 (11.13)
LOO-HC	RF	30.57 (39.54)	13.12 (14.96)	10.78 (8.62)	10.59 (12.45)
	DARNN	<i>25.87 (25.20)</i>	11.03 (11.78)	12.39 (10.78)	10.90 (9.37)
	DTML	42.96 (50.41)	11.25 (12.48)	9.80 (7.66)	9.68 (9.29)
	DeIT	26.07 (29.83)	10.42 (12.59)	8.86 (6.70)	7.82 (6.51)
	TENER	50.58 (69.71)	11.24 (12.19)	8.55 (6.54)	10.72 (11.17)
	DCMN	22.72 (25.40)	<i>11.15 (12.11)</i>	<i>10.06 (9.12)</i>	<i>9.47 (13.07)</i>
LOO-PD	RF	32.89 (47.70)	<i>10.93 (11.34)</i>	10.82 (8.75)	9.28 (9.38)
	DARNN	38.81 (48.67)	10.48 (11.07)	8.30 (6.82)	11.66 (14.38)
	DTML	51.19 (61.17)	11.37 (13.26)	10.05 (9.26)	<i>9.45 (9.78)</i>
	DeIT	29.80 (37.75)	10.65 (14.17)	9.77 (10.36)	9.81 (10.36)
	TENER	50.58 (69.71)	11.32 (14.54)	9.96 (8.05)	9.93 (9.96)
	DCMN	<i>30.71 (37.03)</i>	10.98 (14.05)	<i>10.09 (8.43)</i>	9.91 (11.00)

across different PD subjects. Given the layout of the home, we can assume that the hallway is a hub, and ‘Room-to-room Transition’² shows the transition duration (in seconds) between two rooms connected by the hallway. We choose the transition between (1) kitchen and living room, (2) kitchen and dining room, and (3) dining room and living room since these transitions are more common across PD subjects. ‘Daily Transition’ represents how active PD subjects are in their day-to-day activities, while ‘Room-to-room Transition’ may provide insight into how severe their disease is by the way they navigate their home environment.

DCMN performs the best on average for both ‘Daily Transition’ and ‘Room-to-room Transition’. For the ‘Daily Transition’, DeIT has the second best average offset of 10.07 transitions to the ground truth, while DCMN has the average offset of 8.88 transitions. For the ‘Room-to-room Transition’, DCMN has the average offset of 1.13 seconds to the ground truth, while the second best (RF) has an offset of 1.39 seconds.

Ablation Study. We compare the performance of DCMN and its variants where each of the main components is removed in Table 3:

² The transition is undirected

Table 3: An ablation study of DCMN on ALL-HC, LOO-HC, and LOO-PD. Standard deviation is shown under (.), the worst performer is bold, and the second worst is italicized.

Training Data	Models	Precision	Accuracy	F1-Score
ALL-HC	DCMN - LSTM	95.50	93.60	94.40
	DCMN - GRN	95.40	93.70	94.50
	DCMN - Transformer	93.7	93.60	93.6
	DCMN - CRF	95.50	93.80	94.60
	DCMN - ACCL	95.40	93.50	<i>94.30</i>
LOO-HC	DCMN - LSTM	89.03 (6.70)	88.38 (4.38)	88.75 (5.62)
	DCMN - GRN	88.71 (6.78)	89.33 (3.55)	88.75 (5.34)
	DCMN - Transformer	89.15 (6.43)	89.62 (3.25)	89.08 (5.01)
	DCMN - CRF	89.63 (6.92)	90.02 (3.59)	89.51 (5.46)
	DCMN - ACCL	<i>88.74 (6.64)</i>	<i>88.55 (5.57)</i>	87.88 (6.67)
LOO-PD	DCMN - LSTM	89.53 (2.42)	88.78 (3.21)	88.66 (3.29)
	DCMN - GRN	86.56 (5.00)	87.80 (5.77)	87.56 (4.16)
	DCMN - Transformer	88.31 (4.70)	88.84 (6.97)	88.58 (4.75)
	DCMN - CRF	<i>87.07 (5.05)</i>	89.11 (5.34)	88.72 (4.20)
	DCMN - ACCL	89.42 (2.78)	<i>88.57 (3.88)</i>	<i>88.20 (4.40)</i>

- DCMN - LSTM: DCMN where the Input Attention LSTM is replaced by a positional encoding and a linear layer.
- DCMN - GRN: DCMN where the GRN is replaced by a simple linear layer to combine two modalities.
- DCMN - Transformer: DCMN where the transformer encoder is removed.
- DCMN - CRF: DCMN where the last layer CRF is replaced by a simple linear layer for room-level prediction.
- DCMN - ACCL: DCMN where the accelerometer modality is removed.

Each component improves the prediction accuracy, and DCMN having all the components produces the best accuracy. We observe that adding extra modalities does increase the performance of DCMN for indoor localisation. Without accelerometer data, the DCMN suffers a performance drop in many categories specifically in LOO-HC validation. GRN also plays an important role to merge modalities and suppress unnecessary modality if needed. Table 3 shows that without GRN ability to suppress noisy accelerometer data, the DCMN suffers a high performance drop specifically in LOO-PD validation. Removing the CRF layer has the least impact to the overall performance of DCMN. However, we believe that it helps DCMN in in-home mobility assessment as it constraints DCMN from predicting impossible transition between rooms.

5 Conclusion

We proposed DCMN, a dual modality deep learning model that jointly learns temporal representations of different rooms for indoor localisation. We evaluated our proposed model on data collected of people with and without Parkinson’s disease (PD) living freely in a smart home. During data collection, subjects were wearing wrist-worn wearable accelerometers that provided RSSI and accelerometer data. The novelty of our method is based on using an IoT platform to collect, and effectively utilise, data from multiple sensors to better measure in-home mobility within the context of PD. The use of the two data modalities in our approach provides enriched room signatures to discriminate adjacent rooms that tend to be similar from the perspective of the typically used RSSI signal. Due to a large variance in PD severity, however, some modalities (e.g. accelerometer data) may suffer from inconsistent representations across different subjects. As a result, combining the accelerometer representation into RSSI may degrade the performance of a model for indoor localisation. We demonstrated that our proposed model is able to cope with this problem by outperforming existing methods in predicting where PD subjects are within the home.

Furthermore, our data was collected in naturalistic settings to capture subjects’ natural in-home mobility as opposed to short measurements in clinical environments. We proposed the use of transformer models to learn global temporal relationships among RSSI and accelerometer data and, indirectly, capture a smooth room-to-room transition. We also add a CRF to further enforce correct room-to-room transitions and minimize involuntarily jumps between room predictions due to common RSSI (and accelerometer) fluctuations. Using both transformer and CRF models in tandem, our model managed to outperform other baseline multimodal models in predicting the number of daily room transitions. Our model also shows its capability in approaching the correct room-to-room transition duration, outperforming others with an average offset of 1.13 seconds to ground truth.

References

1. Arora, S., Venkataraman, V., Zhan, A., Donohue, S., Biglan, K., Dorsey, E., Little, M.: Detecting and monitoring the symptoms of parkinson’s disease using smartphones: A pilot study. *Parkinsonism and Related Disorders* **21**(6), 650–653 (2015)
2. Clevert, D.A., Unterthiner, T., Hochreiter, S.: Fast and accurate deep network learning by exponential linear units (elus) (2015)
3. Dauphin, Y.N., Fan, A., Auli, M., Grangier, D.: Language modeling with gated convolutional networks. In: Precup, D., Teh, Y.W. (eds.) *Proceedings of the 34th International Conference on Machine Learning*. *Proceedings of Machine Learning Research*, vol. 70, pp. 933–941. PMLR (06–11 Aug 2017)
4. Fisher, J.M., Hammerla, N.Y., Ploetz, T., Andras, P., Rochester, L., Walker, R.W.: Unsupervised home monitoring of parkinson’s disease motor symptoms using body-worn accelerometers. *Parkinsonism & Related Disorders* **33**, 44–50 (2016)

5. Fraiwan, L., Khnouf, R., Mashagbeh, A.R.: Parkinson's disease hand tremor detection system for mobile application. *Journal of Medical Engineering & Technology* **40**(3), 127–134 (2016)
6. Greenland, J.C., Williams-Gray, C.H., Barker, R.A.: The clinical heterogeneity of parkinson's disease and its therapeutic implications. *European Journal of Neuroscience* **49**(3), 328–338 (2019)
7. Heidarvincheh, F., McConville, R., Morgan, C., McNaney, R., Masullo, A., Mirmehdi, M., Whone, A.L., Craddock, I.: Multimodal classification of parkinson's disease in home environments with resiliency to missing modalities. *Sensors* **21**(12) (2021)
8. Huber, P.J.: *Robust Estimation of a Location Parameter*, pp. 492–518. Springer New York, New York, NY (1992)
9. Ibrahim, M., Torki, M., ElNainay, M.: Cnn based indoor localization using rss time-series. In: 2018 IEEE Symposium on Computers and Communications (ISCC). pp. 01044–01049 (2018). <https://doi.org/10.1109/ISCC.2018.8538530>
10. Ibrahim, M., Torki, M., ElNainay, M.: Cnn based indoor localization using rss time-series. In: 2018 IEEE Symposium on Computers and Communications (ISCC). pp. 01044–01049 (2018)
11. Jankovic, J.: Parkinson's disease: clinical features and diagnosis. *Journal of Neurology, Neurosurgery & Psychiatry* **79**(4), 368–376 (2008)
12. Kozłowski, M., Byrne, D., Santos-Rodríguez, R., Piechocki, R.: Data fusion for robust indoor localisation in digital health. In: 2018 IEEE Wireless Communications and Networking Conference Workshops (WCNCW). pp. 302–307 (2018)
13. Li, M.H., Mestre, T.A., Fox, S.H., Taati, B.: Vision-based assessment of parkinsonism and levodopa-induced dyskinesia with pose estimation. *Journal of neuro-engineering and rehabilitation* **15**(1), 1–13 (2018)
14. Lim, B., Arik, S.O., Loeff, N., Pfister, T.: Temporal fusion transformers for interpretable multi-horizon time series forecasting (2019)
15. Liu, L., Jiang, H., He, P., Chen, W., Liu, X., Gao, J., Han, J.: On the variance of the adaptive learning rate and beyond (2019)
16. Masullo, A., Burghardt, T., Damen, D., Hannuna, S., Ponce-López, V., Mirmehdi, M.: Calorinet: From silhouettes to calorie estimation in private environments (2018)
17. Masullo, A., Burghardt, T., Damen, D., Perrett, T., Mirmehdi, M.: Person re-id by fusion of video silhouettes and wearable signals for home monitoring applications. *Sensors* **20**(9) (2020)
18. Misra, D.: Mish: A self regularized non-monotonic activation function (2019)
19. Pandey, A., Sequeira, R., Kumar, S.: Joint localization and radio map generation using transformer networks with limited rss samples. In: 2021 IEEE International Conference on Communications Workshops (ICC Workshops). pp. 1–6 (2021)
20. Pfister, F.M., Um, T.T., Pichler, D.C., Goschenhofer, J., Abedinpour, K., Lang, M., Endo, S., Ceballos-Baumann, A.O., Hirche, S., Bischl, B., et al.: High-resolution motor state detection in parkinson's disease using convolutional neural networks. *Scientific reports* **10**(1), 1–11 (2020)
21. Qin, Y., Song, D., Chen, H., Cheng, W., Jiang, G., Cottrell, G.: A dual-stage attention-based recurrent neural network for time series prediction (2017)
22. Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., Jegou, H.: Training data-efficient image transformers and distillation through attention. In: Meila, M., Zhang, T. (eds.) *Proceedings of the 38th International Conference on Machine Learning*. vol. 139, pp. 10347–10357. PMLR (18–24 Jul 2021)
23. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L.u., Polosukhin, I.: Attention is all you need. In: Guyon, I., Luxburg, U.V., Bengio,

- S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (eds.) *Advances in Neural Information Processing Systems*. vol. 30. Curran Associates, Inc. (2017)
24. Yan, H., Deng, B., Li, X., Qiu, X.: Tener: Adapting transformer encoder for named entity recognition (2019)
 25. Yoo, J., Soun, Y., Park, Y.c., Kang, U.: Accurate multivariate stock movement prediction via data-axis transformer with multi-level contexts. p. 2037–2045. *KDD '21*, Association for Computing Machinery, New York, NY, USA (2021)
 26. Zhang, M.R., Lucas, J., Hinton, G., Ba, J.: Lookahead optimizer: k steps forward, 1 step back (2019)
 27. Zhang, W., Liu, K., Zhang, W., Zhang, Y., Gu, J.: Deep neural networks for wireless localization in indoor and outdoor environments. *Neurocomputing* **194**, 279–287 (2016)