

CHARM: A Hierarchical Deep Learning Model for Classification of Complex Human Activities Using Motion Sensors

Eric Rosen¹ and Doruk Senkal¹

Abstract—In this paper, we report a hierarchical deep learning model for classification of complex human activities using motion sensors. In contrast to traditional Human Activity Recognition (HAR) models used for event-based activity recognition, such as step counting, fall detection, and gesture identification, this new deep learning model, which we refer to as CHARM (Complex Human Activity Recognition Model), is aimed for recognition of high-level human activities that are composed of multiple different low-level activities in a non-deterministic sequence, such as meal preparation, house chores, and daily routines. CHARM not only quantitatively outperforms state-of-the-art supervised learning approaches for high-level activity recognition in terms of average accuracy and F1 scores, but also automatically learns to recognize low-level activities, such as manipulation gestures and locomotion modes, without any explicit labels for such activities. This opens new avenues for Human-Machine Interaction (HMI) modalities using wearable sensors, where the user can choose to associate an automated task with a high-level activity, such as controlling home automation (e.g., robotic vacuum cleaners, lights, and thermostats) or presenting contextually relevant information at the right time (e.g., reminders, status updates, and weather/news reports). In addition, the ability to learn low-level user activities when trained using only high-level activity labels may pave the way to semi-supervised learning of HAR tasks that are inherently difficult to label.

I. INTRODUCTION

Human Activity Recognition (HAR), time-domain classification of human activities using sensor data, has been garnering increased interest with the advent of wearable sensing technologies. In today’s wearable devices, HAR enables important downstream applications, such as Human-Machine Interaction (HMI), gesture recognition for user interaction, as well as health and wellness tracking. While a wide variety of different sensor modalities can be used for HAR, such as vision or audio, this work focuses on a subclass of HAR that uses body-worn Inertial Measurement Units (IMUs) [1]. IMUs are ubiquitous in today’s mobile and wearable devices, as they provide a low-cost, low-power, self-contained, and privacy-focused sensing modality for both indoor and outdoor HAR.

Previous works on HAR using IMUs have largely focused on event-based activities, such as step counting, man-down detection, and gesture identification [2]. Event-based activities, also termed low-level activities, are characterized by a discrete event that occurs within a finite time window, typically on the order of a few seconds or less. State-of-the-art approaches use machine learning techniques with

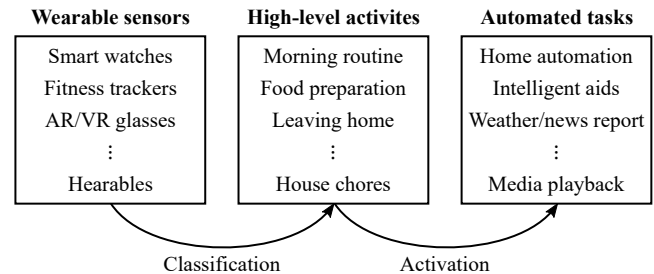


Fig. 1. Examples of how classification of high-level activities can be used for Human-Machine Interaction (HMI). User can choose to activate a user-defined automated task, whenever a pre-determined high-level activity is performed.

hand-crafted features or neural networks to train supervised classification models to infer whether the event has occurred within a specific time window. In this work, we investigate complex activities, also termed high-level activities, which contain multiple event-based activities in a non-deterministic sequence and occur over a long, highly variable duration (e.g., a person cleaning a room can manipulate and navigate the environment in highly varied ways and may take longer to complete the task for larger rooms). Related works in deep learning for computer vision have investigated recognizing high-level activities [3], but high-level HAR using wearable sensors, such as IMUs, have not been studied extensively [1]. Our experimental results suggest that traditional machine learning approaches have challenges in high-level HAR in terms of consistent precision and recall across all classes. In addition, collecting supervised labels for low-level activities to support high-level activity recognition is costly and typically requires domain expertise.

We postulate that a scalable high-level HAR solution would need to autonomously extract relevant low-level motion representations to support high-level activity recognition without any supervised low-level motion labels. We are motivated by the observation that high-level activities can be represented by a composition of low-level activities, similar to how a sentence is composed of individual words. For example, a high-level activity such as “cleaning a room” involves sequencing specific locomotion modes such as walking or standing and manipulation gestures such as picking up objects or opening drawers.

In this paper, we present CHARM (Complex Human Activity Recognition Model), a two-stage neural network architecture for classifying high-level activities from wearable sensor data that learns to represent low-level motion patterns without any low-level motion labels. The first stage uses

(Corresponding author: Doruk Senkal.)

¹Eric Rosen and Doruk Senkal are with Meta Platforms, Inc., 1 Hacker Way, Menlo Park, CA 94025 (e-mail: ericrosen@fb.com; dsenkal@fb.com)

a low-level neural encoder to compress short sequences of motion sensor data into a continuous feature representation. The second stage infers high-level activities from sequences of low-level encoder outputs strided across the raw data stream. CHARM exploits the natural structure of motion patterns during high-level activities, namely that they are complex compositions of low-level motion behaviors that may be sequenced in many ways. Our model architecture enables the low-level neural encoder to focus on shorter, localized motion patterns that characterize low-level motion patterns (over seconds) and the high-level encoder to focus on global patterns found across the long-horizon motion sensor stream (over minutes).

We quantitatively test CHARM’s capability to accurately classify high-level activities from motion sensor data by using the publicly available OPPORTUNITY data set [4], [5], which includes four users performing four different high-level activities: “morning routine”, “coffee time”, “lunch”, and “cleanup”. To test if CHARM’s low-level neural encoder learns semantically meaningful low-level motion representations when trained end-to-end for high-level activity recognition, we use Principal Component Analysis (PCA) to visualize low-level feature representations. We qualitatively find that the low-level neural encoder automatically learns to characterize relevant low-level motion patterns for high-level activity recognition without any explicit labels. This opens the door for new opportunities for efficiently learning to detect low-level activities on wearables by using labeled data of daily high-level activities. To our knowledge, this is the first demonstration of a deep neural network architecture classifying high-level human activities and automatically extracting low-level motion patterns using motion sensor data. The ability to detect high-level activities in a scalable manner may enable many new down-stream applications for Human-Machine Interaction, where the user can choose to tie a user-defined automated task to a high-level activity, Fig. 1.

II. RELATED WORK

Human Activity Recognition (HAR) can be divided into two broad categories based on the sensing modality [1]: vision-based HAR and sensor-based HAR. Vision-based HAR focuses on using visual information such as images, whereas sensor-based HAR focus on other modalities like motion, audio, and biosensors located on the human body or objects in a smart home. Our review of related work focuses only on HAR using body-worn sensors, such as Inertial Measurement Units that consist of accelerometers and gyroscopes for sensing linear accelerations and angular velocities at the location they are worn.

In [6], a Convolutional Neural Network (CNN) was used for recognizing human activities with accelerometer data and experimentally demonstrated that the deep learning architectures can achieve high accuracy for classification of locomotion modalities. In [7], CNNs and other classical machine learning methods (Random Forests with hand-crafted features, naive bayes, and support vector machines) were

used to perform HAR on elementary activities with wrist-based accelerometers. It was found that with a large enough data set, the deep learning approach outperformed the other models. Choosing an appropriate feature representation of the raw motion sensor data has been shown to be important for successful classification as demonstrated by methods leveraging hand-crafted features based on the sensor data and using bag-of-words approaches [8], [9].

An advantage of deep learning approaches is that learned feature representations from raw motion sensor data can be reused for classifying different activities and can be more robust to new on-body sensor locations via transfer learning [10]. These works have demonstrated the effectiveness of deep learning models for HAR with motion sensors, but have only focused on event-based activities, such as locomotion and event-based activity recognition. Other works have investigated recurrent models such as Long Short Term Memory (LSTM) networks for HAR [11], which are powerful inference models because they can operate on arbitrary length sequences. However, [12] found that simple convolutional network architecture outperforms canonical recurrent models on a wide variety of sequence modeling tasks, which inspired our particular instantiation of CHARM in our experiments.

Detecting complex activities from motion sensor data has also been recognized as an important area of research, with a large focus on leveraging low-level motion patterns or body components to aid in classification. For example, in [13] a context-free grammar based representation was leveraged to decompose and classify complex activities by representing them as compositions of body-part and gesture layers. In [14], frequent patterns from low-level actions were mined to construct intermediate representations for complex activity recognition. These works all assume that the requisite low-level components have already been labeled, in addition to the high-level activities. CHARM instead automatically learns to represent the low-level motion patterns by only using labels of high-level, complex activities.

III. PROBLEM FORMULATION

We follow [1] in defining the HAR problem. We suppose that a user is performing an activity belonging to a predefined activity set A :

$$A = \{A_i\}_{i=1}^m, \quad (1)$$

where m denotes the number of activity classes. The user produces a stream of sensor data s that captures signals of the activity:

$$s = \{d_1, d_2, \dots, d_t, \dots, d_n\}, \quad (2)$$

where d_t represents the sensor data with dimension q ($d_t \in \mathbb{R}^q$) collected at time t .

A solution to the HAR problem is to define a model F that predicts the activity A based on the sensor stream s , while the true sequence of activities being performed (ground truth) is denoted as A^* .

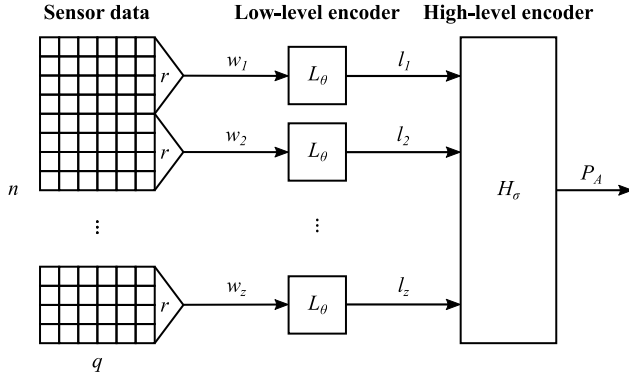


Fig. 2. A schematic of the proposed deep learning architecture for high-level human activity recognition using wearable sensors, consisting of a low-level encoder to learn representations of short time-scale, event-based manipulation and locomotion activities present in the sensor data, which then feeds into a high-level encoder to make predictions about what high-level activity is occurring.

We assume access to batches of supervised data sets D_b that consists of pairs of sensor data streams and ground truth high-level activity labels:

$$D = \{s, A^*\}. \quad (3)$$

The goal is to find a model that minimizes the discrepancy between the predicted classes and the ground truth sequence of activities. We can model this as an optimization problem by defining a loss function \mathcal{L} that captures the discrepancy between $F(s)$ and A^* as a real number:

$$\mathcal{L}(F(s), A^*) \in \mathbb{R}. \quad (4)$$

Given a supervised data set of sensor streams and ground truth activity labels, we search for an F that minimizes the loss function \mathcal{L} . We defer to [1], [2] for additional details of a conventional HAR pipeline.

IV. PROPOSED ALGORITHM

To solve the high-level HAR problem with motion sensor data, we propose a hierarchical neural network architecture with two stages, which we term CHARM: a low-level neural encoder L_θ and a high-level neural encoder H_σ , where θ are the parameters to the low-level neural encoder L and σ are the parameters to the high-level encoder H respectively. The CHARM network architecture uses the low-level encoder to repeatedly featurize sequential short time windows in the raw sensor data, and then uses the high-level neural encoder to classify the high-level activity based on the sequence of outputs from the low-level encoder, Fig. 2. By structuring the network to reuse the same low-level encoder across the input sensor data stream, the low-level encoder produces feature representations of the shorter motion patterns that are invariant to when they occur in the global sequence (similar to how a convolutional filter is invariant to translation in space via a limited kernel size). Note that the CHARM approach is agnostic to the particular deep neural network architecture of the low-level and high-level neural encoder.

TABLE I

THE HIGH-LEVEL AND LOW-LEVEL ACTIVITIES IN THE OPPORTUNITY DATA SET THAT WERE USED IN THE EXPERIMENTS. NOTE THAT CHARM IS TRAINED END-TO-END ONLY ON THE HIGH-LEVEL ACTIVITY LABELS.

High-level activities	Low-level activities	
	Locomotion	Manipulation
Morning routine	Standing	Stir, Sip
Coffee time	Walking	Cut, Spread, Bite
Lunch	Sitting	Open, Close
Cleanup	Lying	Lock, Unlock

More formally, from the data set D , for any given sensor stream sample $s = \{d_1, d_2, \dots, d_t, \dots, d_n\}$ with n samples (the input s is a tensor with shape $[n, q]$ and d_t is a vector with shape $[q]$ for the q motion stream channels), the first stage of the CHARM network applies a sliding window of size r across the input stream with stride of equal length, producing a window sequence $s_w = \{w_1, w_2, \dots, w_t, \dots, w_z\}$ (s_w is a tensor with shape $[z, r, q]$ where z is equal to $\frac{n}{r}$, and w_t is a tensor with shape $[r, q]$). We then feed the batch of windowed frames s_w into the low-level neural encoder $L_\theta(s_w)$ to create a q' -dimensional featurization of the batch of window frames $s_l = \{l_1, l_2, \dots, l_t, \dots, l_z\}$ (s_l is a tensor with shape $[z, q']$, and l_t is a vector with shape $[q']$). The end of the first stage produces s_l since it represents the featurization of all the windowed data frames. The second stage of CHARM takes the output of the first stage s_l , and feeds the input to the high-level neural encoder H_σ which produces a normalized m -dimensional vector P_A that represents the probability of each activity within the pre-determined activity set A being present in the raw sensor stream s .

When referring to the full CHARM approach, we use the notation E_ϕ , where $\phi = \{\theta, \sigma\}$ refers to the collection of the low-level and high-level neural encoder parameters. During inference time, for a given sensor data stream s , the model E_ϕ outputs probability distribution P_A for the pre-determined activity set A :

$$P_A = E_\phi(s). \quad (5)$$

To estimate the parameters ϕ , we perform stochastic gradient descent with batches of data D_b . We use the negative log-likelihood as our objective function.

Before training and validation, we apply standard pre-processing steps to the raw sensor data that are agnostic to the specific task and data set. We normalize the data by subtracting off the mean and dividing by the standard deviation for each of the sensor channels separately. If there are any sequences that have multiple simultaneous high-level labels present in a given data sequence, we do not use that sequence for training or validation.

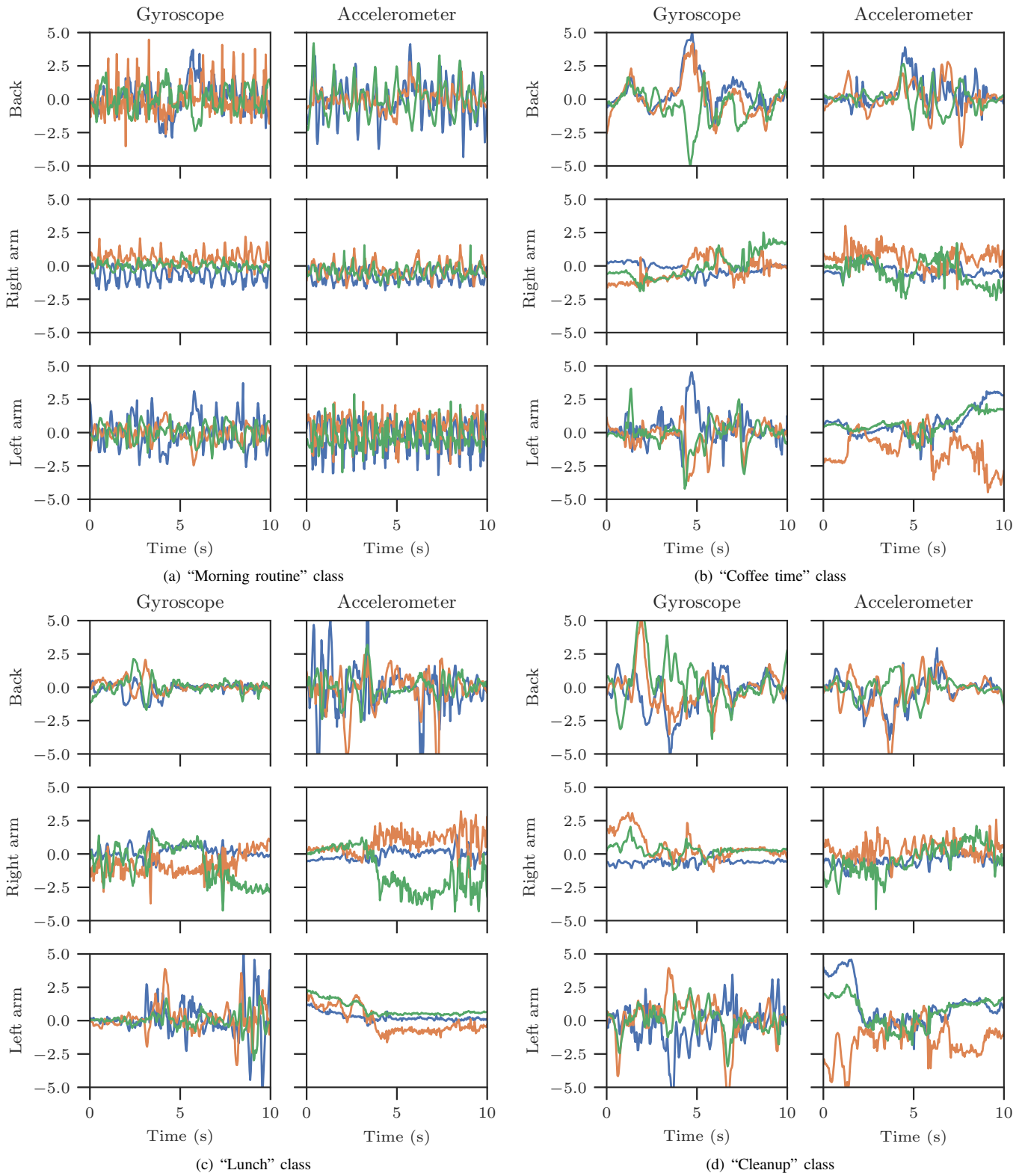


Fig. 3. Raw motion sensor data from the four different high-level activity classes. Each visualization shows 10 s of data from the three sensor locations (Left arm, Right arm, and Back) for the triaxial gyroscope (rad/s) and accelerometer (gee) output. X, Y, and Z sensor axes are colored blue, orange, and green respectively.

V. EXPERIMENTAL RESULTS

A. Data Set

For the experiments we used the OPPORTUNITY data set, which is a publicly available data set for evaluating HAR models [4], [5]. The data set contains readings of motion sensor data from four users performing daily living activities in a room simulating an apartment flat, Table I. To collect the motion sensor data, participants wore a jacket outfitted with triaxial IMUs (30 Hz data sampling rate) located at lower left arm, lower right arm, and upper back positions, Fig. 3. Participants were instructed to perform 5 daily activities:

Morning routine: Users began lying on a canvas chair. At the user’s own pace, they stood up and went out the door of the apartment for a walk. After a leisurely walk, they then re-entered the flat and closed the door behind them.

Coffee time: The users prepared and drank a cup of coffee. Users grabbed a cup located on a rack, put coffee mix and milk into a machine, stirred the ingredients together, and drank the coffee at their leisure.

Lunch: Users prepared a sandwich and ate it. Sandwich preparation involved cutting two slices of bread with a knife and spreading cheese on the slices with a knife. The users then put salami on the slices, and the sandwiches were then placed onto a plate for eating.

Cleanup: After eating and drinking, users cleaned up their apartment. Users put food back to the original locations, put any dishes into the dishwasher, and wiped the tables with a towel.

Relaxing: Users walked around the apartment at their leisure. The users turned off appliances in the room, and then laid down back onto their chair to rest.

Although not explicitly instructed to, during the execution of these daily activities, the users performed atomic manipulation and locomotion actions, e.g., reaching for drawers, walking etc. To model the temporal decomposition of these daily-activities, the OPPORTUNITY data set contains three relevant label tracks aligned with the sensor data: a daily activity track, a locomotion activity track, and a manipulation activity track. In this paper, we treat the daily activity labels as examples of high-level activities, while treating manipulation and locomotion labels as examples of low-level activities.

For our experiments, we leveraged the triaxial gyroscope and accelerometer data collected at the lower-left arm, lower-right arm, and upper back positions. This is because they are the most representative locations for existing wearable technology, e.g., wrist-based and head-mounted devices. Note that because of this, for all of our experiments, each data point in the sensor sequence had 18 channels.

Two major challenges in using the OPPORTUNITY data set for this study were the semantics of the null class and class imbalance: in the OPPORTUNITY data set sequences labeled as “null” had no interpretable meaning to them as they occurred between labeled high-level activities and therefore contained motion patterns from multiple different high-level activities, making it a poor class to train and evaluate

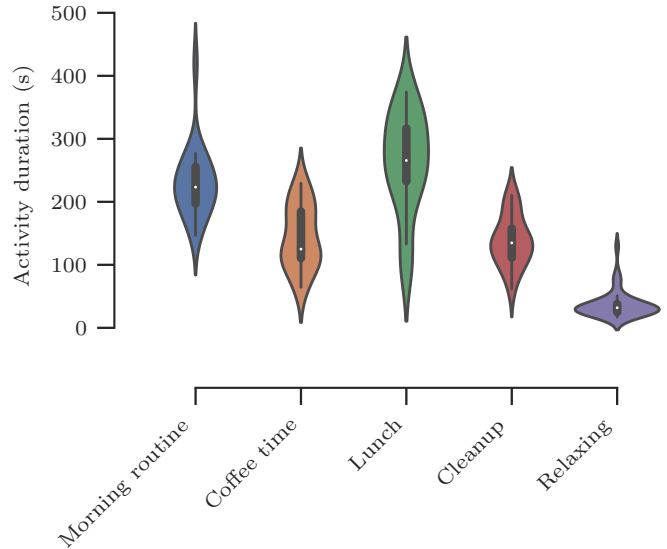


Fig. 4. Violin plots showing the duration of high-level activities in the OPPORTUNITY data set. Note that the “Relaxing” class is unique in how short in length the activities are (typically less than 100 s) compared to all the other activities (which range from 100 to 400 s).

on. For this reason, we do not include the “null” class in our experiments. In addition, for the high-level activity “relaxing”, we found there was a significant imbalance in the activity length compared to all the other high-level activities. Specifically, the longest instances of the “relaxing” activity were around 100 s, while the shortest length instances of all the other high-level activities were around 100 s and lasted up to several minutes, Fig. 4. We therefore did not use the “relaxing” class, and only focused on learning from the four high-level other activity labels (“Morning routine”, “Coffee time”, “Lunch”, and “Cleanup”). Note that we train our high-level activity recognition model only using these high-level labels, i.e., CHARM is never exposed to the low-level locomotion and manipulation labels included in the OPPORTUNITY data set.

B. Models

For our experiments, we evaluated 4 different machine learning approaches: Support Vector Machines (SVMs), Random Forests, a Multi-Layer Perceptron, and CHARM. For SVMs and Random Forests, we evaluated the models with both raw sensor data and with hand-crafted features common for event-based activities. For these approaches, we applied 5 standard hand-crafted features from [15] to each of channels in the raw sensor data (note that because there are 18 channels, the size of the resulting input vector dimension is 90). The 5 hand-crafted features we compute for each channels are: the mean, the variance, the difference between the maximum and minimum peaks in the signal, the minimum value, and the maximum value. We also performed hyperparameter sweeps for all models and reported the best results. The following subsections describe the hyperparameters for the models. For all the models, we train using data

TABLE II
PRECISION (P), RECALL (R) AND F1 SCORES FOR CLASSIFICATION OF HIGH-LEVEL ACTIVITIES USING SVM, RANDOM FOREST (RF), MULTI-LAYER PERCEPTRON (MLP), AND CHARM.

	Coffee time			Morning routine			Cleanup			Lunch			Class average		
	P	R	F1	P	R	F1	P	R	F1	P	R	F1	P	R	F1
SVM	1.00	0.91	0.95	0.82	1.00	0.90	1.00	0.59	0.74	0.88	1.00	0.94	0.93	0.88	0.88
RF	0.52	0.26	0.34	0.65	1.00	0.79	0.99	0.43	0.60	0.84	0.95	0.89	0.75	0.66	0.66
MLP	0.75	0.62	0.68	0.88	1.00	0.94	0.91	0.92	0.91	0.99	0.93	0.96	0.88	0.87	0.87
CHARM	0.91	0.93	0.92	0.93	1.00	0.96	0.97	0.88	0.93	0.96	0.97	0.96	0.94	0.95	0.94

files of 3 of the 4 users from the OPPORTUNITY data set, and validate using the last held-out user. All approaches had input sequences with 2560 data points (at 30 Hz, 85 s), each data point with 18 channels consisting of 3 sensor locations, each with 2 sensor modalities (gyroscope and accelerometer) and 3 sensor axes (X, Y, Z).

1) *CHARM*: The low-level encoder window size was 16 (about 500 ms), with a two-layer fully-connected network with a hidden-state size of 32 dimensions. The non-linear activation function between the layers was a Leaky ReLU with a negative slope of 0.01. The high-level encoder was a two-layer fully-connected network with an input sequence size of 160 data points, each with 32 dimensions from the low-level encoder. We applied a softmax to the final output from the high-level encoder to calculate a probability distribution over the high-level classes. We used a cross-entropy loss as our optimization criterion, and gave weights to each class inversely proportional to the number of labels of each class in the training data. We trained with a batch size of 1 for 10 epochs. We use the Adam optimizer with a learning rate of $5e^{-4}$. We performed dropout during training between all the layers with a probability of 5%. At test time, we selected the class with the highest probability.

2) *Multi-Layer Perceptron (MLP)*: The MLP had 4 fully-connected layers, each with a hidden size of 16. The non-linear activation function between the layers was a Leaky ReLU with a negative slope of 0.01. We applied a softmax to the output from the final layer to calculate a probability distribution over the high-level classes. We used a cross-entropy loss as our optimization criterion, and gave weights to each class inversely proportional to the number of labels of each class in the training data. We trained with a batch size of 1 for 10 epochs. We used the Adam optimizer with a learning rate of $5e^{-4}$ and performed dropout during training between all the layers with a probability of 5%.

3) *Support Vector Machines (SVMs)*: We found that Support Vector Machine performed better on the raw data rather than using the hand-crafted features. We used a radial-basis function as the kernel for the SVM, with a scaled gamma kernel coefficient.

4) *Random Forests (RF)*: For the Random Forests, we apply the hand-crafted features to the input sequence. We use 100 trees in the forest for estimation, and we use the Gini criterion for measuring quality of splits in the tree. Splits were made in the tree until the leaves only contained points of the same class.

C. Classification Performance

For each model, we report the precision, recall and F1 scores for all the high-level activity classes in Table II. For the precision, recall, and F1 score for each class, we bold the numbers of the model that performed best in that category. Overall, CHARM has the best performance across most of the metrics, with the multi-layer perceptron having the second-best performance, while the SVM and RF performed the worst, which was expected since the neural network approaches are able to learn nonlinear feature representations of the input data. The overall average accuracy of CHARM is also the highest of all the models (0.95). Almost all of CHARM’s precision, recall, and F1 scores were above 0.90 for all the classes, unlike the other models which do not perform as well and struggle to get a F1 score across all classes above 0.75. This suggests that CHARM is best suited for the task of high-level activity recognition because it does not over-fit to any of motion patterns in the individual classes but also learns to effectively generalize from the training data to the validation set.

D. Low-Level Semantic Representations

CHARM was inspired by the notion that a high-level activity consists of potentially repeated discrete event-based low-level motions that may be sequenced in many different ways to accomplish the task. Because we do not assume access to any labeled examples of low-level motions at training time, we are training both stages of the model end-to-end based only on high-level activity labels.

Convolutional filters in computer vision have been known to learn semantically meaningful hierarchical feature representations of concepts in images [16] when trained end-to-end for a downstream task. For example, when trained to detect objects and animals, visualizations of the feature representations in earlier layers show kernels that represent edge detectors, and later layers learn to detect shapes based on those edges. This inspired us to investigate whether CHARM’s low-level encoder also learns representations of semantically meaningful concepts in the motion patterns present in the high-level activity sequences.

To this end, we first trained the low-level and high-level encoder end-to-end on the high-level activity recognition tasks described in Section V-A. Then, we used the labeled data set for locomotion and manipulation activity tracks to generate labeled sequences with the same size as the low-level encoder (500 ms) for all of the low-level classes.

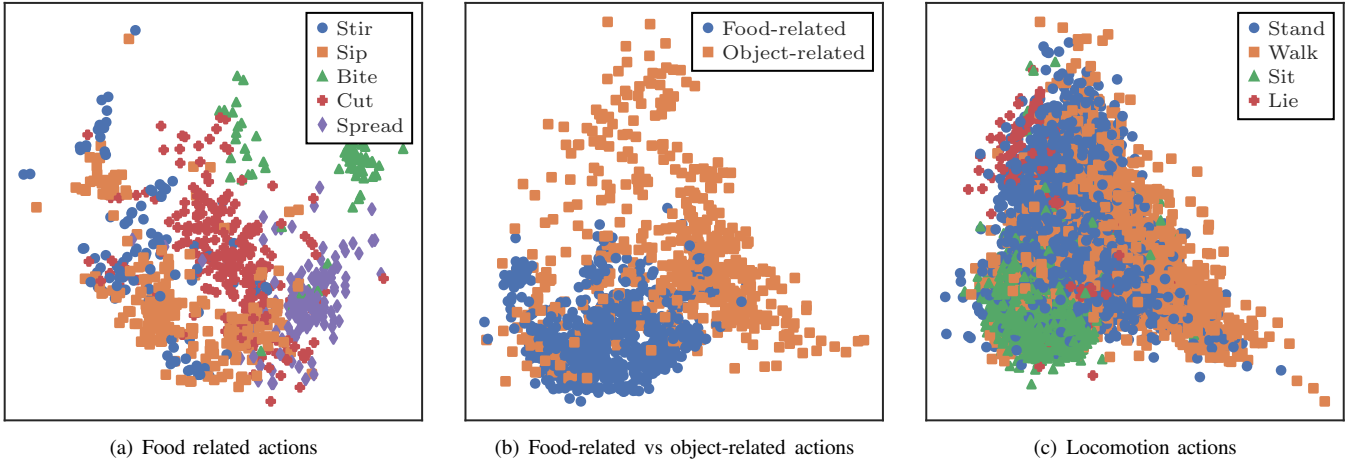


Fig. 5. 2-D visualization of low-level neural encoder representations for different low-level manipulation and locomotion activities. Even though the low-level encoder was not trained with any of these labels, the low-level neural encoder still learns to cluster related motion patterns. (a) Manipulation activities: stirring, sipping, biting, cutting, and spreading. (b) Food-related activities (stirring, sipping, biting, cutting, and spreading) vs object-related activities (unlocking, locking, opening, closing). (c) Locomotion modes: standing, walking, sitting, and lying down.

We then fed all of the sequences into CHARM’s low-level encoder and fed the high-dimensional output through a Principal Component Analysis (PCA) to reduce the feature representation to 2 dimensions. This dimensionality reduction technique enables us to visualize the learned representations of the low-level encoder using a 2-D graph, where we can color the points based on the known label of the manipulation or locomotion sequence. We hypothesized that the low-level neural encoder would learn to cluster low-level motion patterns of the same class, even though the low-level encoder was never exposed to these labels. In Fig. 5, we present visualizations for different manipulation and locomotion classes.

In Fig. 5(a), we plot the learned representations of 5 manipulation behaviors: stirring, sipping, biting, cutting, and spreading. We qualitatively see that each of these manipulation behaviors have distinct clusters, which indicate that CHARM has learned to make distinct representations of these low-level motion patterns. Although “biting” and “sipping” are motion patterns that are very similar to each other in the raw motion sensor data (since both involve the user bringing their hand near their face), CHARM learns to separate these representations from each other and instead embed them closer to manipulation behaviors that are closely related to each other for classification of high-level activities. In addition, we notice that the “biting” activity is close to the “spreading” and “cutting” activities, whereas the “sipping” activity is close to the “stirring” activity. In practice, the “biting”, “spreading” and “cutting” activities almost always occur in the “Lunch” high-level activity, whereas “sipping” and “stirring” activities are almost exclusively performed during coffee preparation. We therefore hypothesize that the distance between embedded motion patterns inside CHARM’s low-level encoder are related to how correlated those motion patterns are to the high-level activities, and not how similar the motion patterns are to each other in raw

sensor data. In other words, CHARM’s low-level encoder embeds motion patterns close to each other based on how relevant they are for classification of high-level activities, and not just based on how similar the underlying motion sensor data is.

To determine if this behavior was consistent across other manipulation behaviors, we investigated two classes of low-level activities: “object-related” activities and “food-related” activities, Fig. 5(b). Object-related activities were typically performed in “cleanup” and “morning routine” high-level activities, and included manipulation gestures like “unlocking”, “locking”, “opening”, and “closing”. On the other hand, food-related activities were mostly used in the “lunch” and “coffee preparation” classes, and involved manipulation behaviors like “sipping”, “cutting”, “stirring”, “bitting”, and “spreading”. We plot the learned representations of these 2 classes in Fig. 5(b), and qualitatively see that there are two distinct clusters of these classes. This further supports our hypothesis that CHARM’s low-level encoder has learned to make distinct feature representations of low-level manipulation gestures, and that the distance between the embedded representations is proportional to how correlated the gestures are to the same high-level activities. In other words, low-level gestures that are prevalent in the same high-level activities are closer in CHARM’s learned embedding space.

We also wanted to investigate whether CHARM learned distinct feature representations of locomotion motion patterns as well. In Fig. 5(c), we plot the learned representations of 4 locomotion behaviors: “standing”, “walking”, “sitting”, and “lying-down”. Similar to the manipulation gestures, we see that each of these locomotion modes has a distinct cluster, although there is much higher-overlap, especially between the “sitting” and “standing” class. However, this is expected because we only have access to motion sensor data, which makes distinguishing standing and sitting extremely challenging. In addition, “lying-down” and “standing” have

extremely high overlap since users who were lying-down during the experiments eventually had to stand up in order to do any of the high-level activities, and “standing” and “walking” have high overlap since users often walked and stood at locations when doing any of the high-level activities. Overall, since almost all of the high-level activities involved almost all of the locomotion modes except for “lying-down”, which was almost exclusively done in the “Morning routine” class because users started lying down, most of the locomotion modes are closer in embedding space than the manipulation behaviors, but there still exist distinct clusters of these manipulation behaviors.

Overall, we postulate that these qualitative assessments validate that CHARM’s low-level encoder is able to learn semantically meaningful representations of low-level motion patterns, even when it is only trained end-to-end on high-level activity labels. The important caveat to this is that, although CHARM learns to generally distinguish individual low-level motion patterns from each other, the distance between clusters is more closely related to how correlated low-level motion pattern are to similar high-level activities, and not the underlying low-level motion patterns themselves. This is important because it implies that the choice of high-level activity label and how common certain low-level motion patterns are in each high-level activity class have a large impact on what kinds of low-level motion patterns can be automatically distinguished. This ultimately makes sense because the model is trained end-to-end on high-level activity labels and feature representations are optimized so as to maximize performance on detecting high-level activities, so low-level motion patterns that are present in the same high-level activity will become closer in the embedding space because they are useful signals for determining what high-level activity is occurring. In other words, when determining, for example, whether a user is preparing lunch, it is sufficient to just check if they are either biting, spreading, or cutting, since as long as one of these activities is occurring it is highly likely that the user is preparing lunch since none of these manipulation behaviors occur in the other high-level activities. Therefore, the model learns to project these motion patterns to similar locations in the embedding space.

VI. CONCLUSION

In this paper, we investigated feasibility of using wearable sensor data and deep learning algorithms for the classification of high-level activity recognition. Motivated by the observation that high-level activities can be represented by a composition of low-level activities, we developed a hierarchical deep learning architecture for HAR, which infers high-level activities from low-level encoder outputs.

Our results suggest that representations of semantically meaningful manipulation and locomotion behaviors can automatically be learned from end-to-end training of deep neural network models for high-level activity recognition using wearable sensor data to enable new applications in Human-Machine Interaction, such as intention recognition

and human-centered automation, as well as health and wellness tracking. New future research directions may include: using automatically learned low-level feature representations for improving event-based detection via transfer learning and application of CHARM architecture for classification of other high-level activities.

REFERENCES

- [1] J. Wang, Y. Chen, S. Hao, X. Peng, and L. Hu, “Deep learning for sensor-based activity recognition: a survey,” *Pattern Recognition Letters*, vol. 119, pp. 3–11, 2019.
- [2] A. Bulling, U. Blanke, and B. Schiele, “A tutorial on human activity recognition using body-worn inertial sensors,” *ACM Computing Surveys (CSUR)*, vol. 46, no. 3, pp. 1–33, 2014.
- [3] Q. Yang, “Activity recognition: linking low-level sensors to high-level intelligence,” in *Proc. International Joint Conference on Artificial Intelligence (IJCAI’09)*, Pasadena, California, USA, July 2009, pp. 20–25.
- [4] D. Roggen, A. Calatroni, M. Rossi, T. Holleczeck, K. Förster, G. Tröster, P. Lukowicz, D. Bannach, G. Pirkel, A. Ferscha, *et al.*, “Collecting complex activity datasets in highly rich networked sensor environments,” in *Proc. IEEE International Conference on Networked Sensing Systems (INSS’10)*, Kassel, Germany, June 2010, pp. 233–240.
- [5] R. Chavarriaga, H. Sagha, A. Calatroni, S. T. Digumarti, G. Tröster, J. d. R. Millán, and D. Roggen, “The opportunity challenge: A benchmark database for on-body sensor-based activity recognition,” *Pattern Recognition Letters*, vol. 34, no. 15, pp. 2033–2042, 2013.
- [6] Y. Chen and Y. Xue, “A deep learning approach to human activity recognition based on single accelerometer,” in *Proc. IEEE International Conference on Systems, Man, and Cybernetics*, Hong Kong, China, Oct. 2015, pp. 1488–1492.
- [7] H. Gjoreski, J. Bizjak, M. Gjoreski, and M. Gams, “Comparing deep and classical machine learning methods for human activity recognition using wrist accelerometer,” in *Proc. International Joint Conference on Artificial Intelligence (IJCAI’16)*, New York, NY, USA, July 2016, p. 970.
- [8] M. Kheirhahan, S. Mehta, M. Nath, A. A. Wanigatunga, D. B. Corbett, T. M. Manini, and S. Ranka, “A bag-of-words approach for assessing activities of daily living using wrist accelerometer data,” in *Proc. IEEE International Conference on Bioinformatics and Biomedicine (BIBM’17)*, Kansas City, MO, USA, Nov. 2017, pp. 678–685.
- [9] M. F. Aslan, A. Durdu, and K. Sabanci, “Human action recognition with bag of visual words using different machine learning methods and hyperparameter optimization,” *Neural Computing and Applications*, vol. 32, no. 12, pp. 8585–8597, 2020.
- [10] F. J. O. Morales and D. Roggen, “Deep convolutional feature transfer across mobile activity recognition domains, sensor modalities and locations,” in *Proc. ACM International Symposium on Wearable Computers (ISWC’16)*, Heidelberg, Germany, Sept. 2016, pp. 92–99.
- [11] M. Edel and E. Köppe, “Binarized-BLSTM-RNN based human activity recognition,” in *Proc. IEEE International Conference on Indoor Positioning and Indoor Navigation (IPIN’16)*, Alcalá de Henares, Madrid, Spain, Oct. 2016.
- [12] S. Bai, J. Z. Kolter, and V. Koltun, “An empirical evaluation of generic convolutional and recurrent networks for sequence modeling,” *arXiv:1803.01271*, 2018.
- [13] M. S. Ryoo and J. K. Aggarwal, “Recognition of composite human activities through context-free grammar based representation,” in *Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’06)*, New York, NY, USA, June 2006, pp. 1709–1718.
- [14] Y. Liu, L. Nie, L. Han, L. Zhang, and D. S. Rosenblum, “Action2activity: recognizing complex activities from sensor data,” in *Proc. International Joint Conference on Artificial Intelligence (IJCAI’15)*, Buenos Aires, Argentina, July 2015, pp. 1617–1623.
- [15] M. Zhang and A. A. Sawchuk, “A feature selection-based framework for human activity recognition using wearable multimodal sensors,” in *Proc. IEEE International Conference on Body Area Networks (BodyNets’11)*, Dallas, Texas, USA, May 2011, pp. 92–98.
- [16] Q.-s. Zhang and S.-c. Zhu, “Visual interpretability for deep learning: a survey,” *Frontiers of Information Technology & Electronic Engineering*, vol. 19, no. 1, pp. 27–39, 2018.