



A Dataset of Routine Daily Activities in an Instrumented Home

Julien Cumin, Grégoire Lefebvre, Fano Ramparany, James L. Crowley

► **To cite this version:**

Julien Cumin, Grégoire Lefebvre, Fano Ramparany, James L. Crowley. A Dataset of Routine Daily Activities in an Instrumented Home. UCAmI 2017 - 11th International Conference on Ubiquitous Computing and Ambient Intelligence, Nov 2017, Philadelphie, United States. pp.413-425, 10.1007/978-3-319-67585-5_43 . hal-01639673

HAL Id: hal-01639673

<https://hal.archives-ouvertes.fr/hal-01639673>

Submitted on 20 Nov 2017

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

A Dataset of Routine Daily Activities in an Instrumented Home

Julien Cumin^{1,2}, Grégoire Lefebvre¹, Fano Ramparany¹, and
James L. Crowley²

¹ Orange Labs, France

{julien1.cumin, gregoire.lefebvre, fano.ramparany}@orange.com

² Univ. Grenoble Alpes, Inria, CNRS, Grenoble INP, LIG, F-38000 Grenoble, France
james.crowley@inria.fr

Abstract. We present a new dataset, called Orange4Home, of activities of daily living of one inhabitant in a smart home environment. We collected data from 236 heterogeneous sensors in a fully integrated instrumented apartment. Data collection spanned 4 consecutive weeks of working days for a total of around 180 hours of recording. 20 classes of varied activities were labeled *in situ*. We report the methodology adopted to establish a representative, challenging dataset, as well as present the apartment and sensors used to collect this data.

Keywords: Dataset · Activities of daily living · Smart home

1 Introduction

Development, evaluation and comparison of machine learning approaches for activity recognition in smart home environments all require realistic data labeled with ground truth. Acquiring labeled data of activities in a smart home is challenging and costly for a number of reasons. For one thing, few instrumented smart homes exist. Furthermore, providing accurate ground truth is costly and tedious. Giving a meaning to data collected by many heterogeneous sensors in an instrumented environment is difficult. For example, finding whether an occupant is currently cooking or not (which might impact the service that the system wants to provide to the occupant), using only low-level sensors such as door openings, luminosity, ambient noise, etc. is challenging. In addition, the range of possible activities is highly variable, and there are no guarantees that activities of an individual in a specific home are representative of the entire population of inhabitants or homes. Thus, establishing a representative dataset of activities in smart homes requires both technical work on instrumenting a home and an important effort to accurately label activities.

In this paper, we present a new labeled dataset, named *Orange4Home*³, of activities of daily living (ADL) of an occupant in an instrumented smart home environment. The Amiqua4Home smart home environment is a fully equipped

³ <http://amiqual4home.inria.fr/orange4home>

furnished 87 m² apartment, that has been instrumented with 236 data sources, which capture information about use of electrical equipment, water consumption, operation of doors, etc. The facility was constructed to serve as a resource for research in smart home services. We used this facility to acquire data about human daily activities, over a period of 20 working days. Labels for 20 classes of activity were noted *in situ* by the occupant, representative of the performed activities. This data is designed to be used for supervised offline activity recognition [4], supervised online activity recognition [6], unsupervised activity discovery [2], activity prediction [5], as well as other applications.

We report in Sect. 2 related datasets of activities in the home. We present in Sect 3 the methodology adopted to establish a realistic dataset of labeled activities, and in Sect. 4 the technical aspects of the experiment. We conclude in Sect. 5 on future uses of this dataset.

2 Related Work

The dataset presented in this work complements the already significant number of existing datasets of ADL, such as the Opportunity dataset, the Transfer Learning dataset, and other datasets of ADL.

Opportunity [7] is a dataset where activities of 4 different occupants are recorded using a high number of both environmental (i.e. fixed in the home) sensors and body-worn (i.e. placed on the occupant) sensors. This dataset provides 3 levels of labeling of activities, from atomic arm gestures to high level activities such as eating a sandwich. However, the *Opportunity* dataset was recorded in a single experimental environment of only 1 room, which lessens its representativeness of real, inhabited homes. Moreover, nearly half of the sensors used are body-worn sensors, which is not realistic for smart home systems aimed at the general population, which is what we are interested in in this work. Finally, the sequence of activities captured in this dataset are very limited in time (about half an hour).

The *Transfer Learning dataset* [9] is another dataset of activities of daily living, recorded in 3 different real homes inhabited by 3 different persons. Labeling was done *in situ* by the occupant for a period of 13/18/25 days respectively, depending on the home. This dataset unfortunately contains the labeling of only 8 classes of activities (not including an “*Other*” activity), which may be too limited to be representative of all significant activities of a home. Moreover, only 23/21/14 sensors respectively were present in each home, which is possibly too little to realistically represent the complexity of future smart home systems.

ARAS [1] is a dataset of activities of daily living recorded in two real houses for a full month. This dataset contains the labels of 27 different activity classes labeled *in situ* by the inhabitant to a good degree of accuracy. Each house was equipped with 20 binary sensors, which is unfortunately restrictive in terms of algorithmic evaluation: it is not possible to use this dataset to experiment on algorithms that deal with heterogeneous data, nor is it possible to study sensors redundancy due to the small number of sensors.

A dataset by Tapia et al. presented in [8] provides the recordings of 77/84 sensors respectively in 2 real homes, inhabited by 2 different persons, for a period of 14 consecutive days. 33 different activities are labeled, each being part of categories of activities. This dataset was recorded with the intent of evaluating recognition of activities useful to healthcare applications, such as care to elderly people, which makes its representativeness of systems aimed at the general population debatable. Moreover, the authors of this work report significant difficulties in *in situ* labeling of activities, which was too coercive for both occupants, leading to imprecise or missing labels which had to be fixed by hand after the experiments.

The MavPad 2005 dataset [10] is a dataset recorded in a student apartment instrumented with 76 sensors of various kinds (light, temperature, humidity, motion, doors, water leak, smoke and CO₂), for a duration of 7 weeks during which one occupant was present in the home. This dataset is unfortunately not labeled with activities, which limits its usefulness to applications such as sensor events prediction, while not being usable for activity-related problems.

In this work, we aim at recording a dataset of activities which combines the positive points of these state-of-the-art datasets into one: our dataset provides the recordings of a high number of heterogeneous sensors scattered seamlessly in a real home, labeled *in situ* with a significant number of representative classes of activities performed by one occupant for a duration of 20 days.

3 Methodology

3.1 Goals of the Experiment

Based on the positive and negative points of state-of-the-art datasets of activities in the home, presented in Sect. 2, we based our methodology for recording the Orange4Home dataset on the following goals that we intend to reach:

1. label accurately all 4 main context information of the home for the entirety of the experiment;
2. record realistic routines of daily living of the general public (i.e. not of a specific population such as elderly people);
3. record data in a realistic environment, which is as close to a real home as possible;
4. record data in a pervasively instrumented environment, where as many objects as possible are instrumented, with as many different types of sensors as possible;
5. record data on a sufficiently long time scale such that the dataset is usable to test activity prediction approaches.

To fulfill goal 1, we need to provide information of *identity*, *time-of-day*, *place* and *activity* [4]. The Orange4Home dataset contains all 4 information, as there is only one occupant (identity is thus unique), all events are timestamped, and both place and activity are labeled *in situ* by the occupant (see Sect. 4.3).

The experiment spanned 4 consecutive weeks of working days, in order to fulfill goal 5.

We present in Sect. 4.1 the instrumented apartment used as the recording environment to fulfill goal 3, and in Sect. 4.2 the sensors and data types recorded during the experiment to fulfill goal 4.

To fulfill goal 2, we need to establish a home occupancy scenario, which will guide the choices of routines the occupants need to perform. We present this scenario in Sect. 3.2.

3.2 Home Occupancy Scenario

In this experiment, we imagine that the home is a coworking apartment in which the subject of the experiment, Bob, comes to work alone every working day, from around 08:00 to 17:00. This apartment is a pervasive environment filled with sensors of various kinds, which transmit their data to a centralized system in charge of the home.

Bob is interested in having personalized services in this coworking environment, based on his activities. As such, he will label his routines for a duration of 20 days (i.e. 4 weeks of working days). Since this is Bob’s coworking apartment, his activities will not only be work-related, but also lunch and leisure related.

Bob does not live in the home outside of working hours, since this is a coworking environment. Therefore, data outside these hours is not provided.

3.3 Activities in Orange4Home

As mentioned in Sect 3.1, identity, time-of-day, places and activities are key information of context. Identity is both the identity in the literal sense of an occupant, but also their social role in the home. Time-of-day can be any semantic temporal information such as a part of the year, whether it’s a week day or not, etc. Places is a geographical location in the home which holds a specific function for the occupant, such as the bathroom or the kitchen. An activity is a set of tasks (or operations), a task being a set of actions [3].

Selecting the actual sets of labels used in the dataset for identity, time-of-day and place values is rather direct: names, timestamps and rooms of the home, respectively, are sufficient, since we can extract semantic labels from those values (e.g. all semantic labels for time-of-day can be obtained from a timestamp). It is less clear what set of labels should be used for activities. For example, in the *Opportunity* dataset [7], there are 3 *levels* of activity labeling, from complex activities like “*Sandwich time*” to atomic gestures like “*Interact with the bottle with the left arm*” (which match the activities, tasks and actions hierarchy).

In this dataset, we only label activities. Actions do not carry any functional meaning and are therefore not useful labels to characterize context in the home. Tasks, like activities, do carry functional meaning, but we believe that contextual information provided by tasks is too limited for general-purpose context-aware services (but could be useful for specific applications). Furthermore, tasks and

actions are too numerous and possibly too short-lived to be effectively and accurately labeled, be it from an occupant in a real-world setting, but also in an experimental setting of data collection.

The activities we aim to record are those that are quite frequent, and fairly recurrent in time. They must also fit the scenario established in Sect. 3.2. We present below the list of activities, grouped by places (see Fig. 3 and Fig. 4) in which they can occur:

- **Entrance** *Entering, Leaving*;
- **Kitchen** *Preparing, Cooking, Washing the dishes*;
- **Living Room** *Eating, Watching TV, Computing*;
- **Toilet** *Using the toilet*;
- **Staircase** *Going up, Going down*;
- **Walkway** (no activity specific to this place);
- **Bathroom** *Using the sink, Using the toilet, Showering*;
- **Office** *Computing, Watching TV*;
- **Bedroom** *Dressing, Reading, Napping*;
- **Common to all places** *Cleaning*.

We have established this list of activities based on the existing places and appliances available in the apartment, such that all realistic and common classes of activities that are usually performed in a coworking home are represented in the dataset. We have also tried to balance classes such that no place contains the majority of classes.

As we can see, some activities can occur in multiple different places (e.g. *Watching TV* in both the Living Room and the Office, *Using the toilet* in both the Toilet and the Bathroom, etc.). It is very common for occupants to be able to perform certain activities in multiple different places, which is not a possibility often captured in state-of-the-art datasets of ADL. Our experiment is set in a real home, in which forbidding those situations would not be realistic; allowing that some activity classes can happen in different places also adds complexity to the dataset, making it more challenging as a benchmark for ADL-related algorithmic problems.

There is no trivial link between activity and sensor data in the general case. Activities can vary in duration and in the way they are performed by the occupant. For example, the activity *Preparing* will be heavily dependent on what the occupant is intending to cook: differences in which cupboards are opened and in what order can for example appear. The occupant can also do mistakes or change their mind on what they want to do in this activity, adding more sensor events.

3.4 occupant’s Routine in Orange4Home

The experiment ran from January 30th 2017 to February 24th 2017, during working days: 08:00 to 17:00 on average, from Monday to Friday.

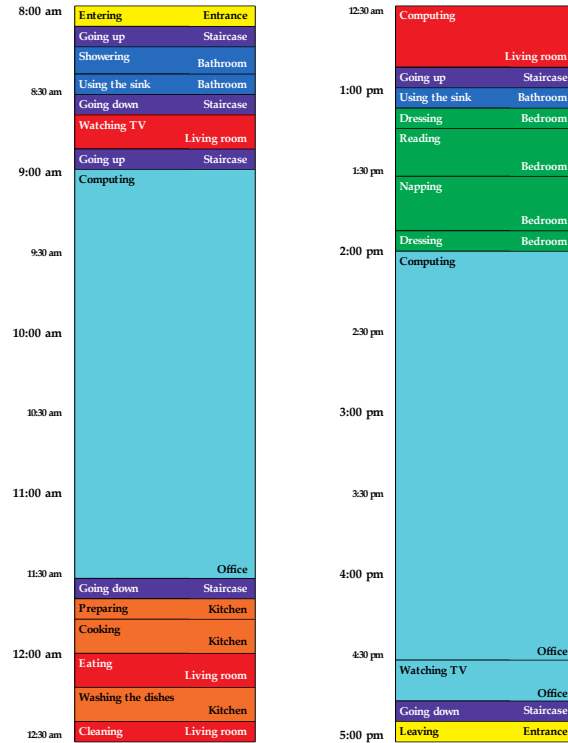


Fig. 1. Standard day routine. Activity is indicated on the left side, and the place in which they are performed on the right side.

Since the occupant is not normally working in this apartment, imposing a specific routine to the occupant is a necessity. Indeed, with very limited familiarity with the home, and with only 20 working days of experiment, an occupant might have the tendency to limit their activities to those they do in their own home and workplace, or with appliances that they are familiar with (although a mock experiment for a single day was performed prior to the real experiment so that the occupant gained more familiarity with the home). Moreover, imposing a routine on the user allows us to make sure that all activities are sufficiently represented in the dataset (so that it is for example usable with supervised machine learning techniques) and that sequences of activities are sufficiently frequent so that the dataset can be used to experiment activity prediction approaches.

The planning has been established around a standard day routine, presented on Fig. 1. This standard day routine has been carefully established with the following goals in mind:

- this routine must contain all activity classes (except for *Using the toilet* which is unpredictable);
- this routine must span an entire working day;
- this routine must fit our coworking scenario presented in Sect. 3.2.



Fig. 2. Routine for February 21st, 2017. Activity is indicated on the left side, and the place in which they are performed on the right side.

The standard day routine chosen therefore includes the occupant showering, using the sink and watching TV news in the morning, cooking and eating their food in the apartment, spending leisure time after lunch on their computer, reading and napping, and spending the rest of their time during the morning and the afternoon working on their computer in the Office.

This routine is strictly followed during the first two weeks (with leniency on the times at which activities start and end, and with slight differences on Fridays), and minor to major changes (such as interversion, omission, shortening, etc. of activities) to this daily routine are applied to establish the planning for the last two weeks. For example, we present on Fig. 2 the daily routine followed on Tuesday of the fourth week (February 21st, 2017). We can see fairly significant deviations from the standard daily routine this day: the occupant leaves the home at lunch time, the occupant does not perform any of the leisure activities during the lunch break, but instead goes back to work upon reentering the home (after the activity *Using the sink*).

This way, the dataset contains a routine of sufficiently recurring activities in the first 2 weeks, while also containing minor to major variations on that routine in the 2 last weeks; this allows the dataset to be used as a challenging

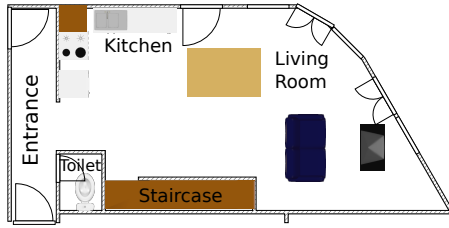


Fig. 3. Ground floor of the apartment

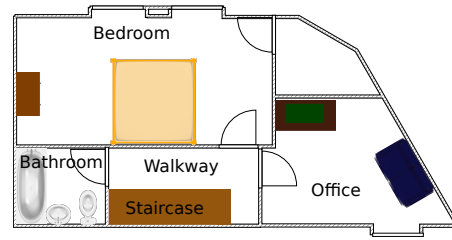


Fig. 4. First floor of the apartment

benchmark, not only for activity recognition problems but also for problems of activity prediction in time (over minutes, hours, days or even weeks time scales), with any type of algorithmic approach including supervised models. The variations in the last two weeks prevents activity prediction from being a trivial task with this dataset, while maintaining realism and coherency in the way routines evolve from one week to another.

4 Data Collection

4.1 Overview of the Apartment

The *Amiqual4Home* project⁴ is an experimental platform which comprises prototyping workshops, living labs and various mobile tools to be used in ambient intelligence research projects. Among those living labs, there is an instrumented apartment⁵: an 87 m² two-story home fully instrumented with sensors and actuators of various kinds throughout its rooms. The purpose of this apartment is to provide an environment in which to experiment on smart home systems in a real setting of an instrumented home. In particular, it can be used to record sensors' data during long periods of inhabited times, as was done in this work.

Figure 3 and Fig. 4 presents the layout of the apartment, annotated with the names of places of the apartment.

4.2 Sensors and Data Collection

The apartment has been furnished with the explicit goal of being an experimental apartment; as such, it is instrumented with many fully integrated sensors that are either not visible or not a hindrance to the occupant, as opposed to sensors installed in a standard apartment. The extensive list of sensors provides many different kinds of data, such as doors or cupboards opening, ambient noise, temperature, CO₂ levels, presence, switches being pressed, electrical information about appliances, hot and cold water consumption, luminosity, heaters information, weather information, etc. Moreover, all sensors are associated to the room

⁴ <https://amiqual4home.inria.fr/>

⁵ <http://amiqual4home.inria.fr/tools/smart-home>

Place	Data type				Total
	Binary	Integer	Real number	Categorical	
Entrance	3	1	2	3	9
Kitchen	13	21	18	0	52
Living room	16	6	8	7	37
Toilet	3	1	1	0	5
Staircase	3	0	0	0	3
Walkway	9	0	1	0	10
Bathroom	9	6	8	3	26
Office	9	3	3	5	20
Bedroom	17	4	6	7	34
Global	1	13	20	6	40
Total	83	55	67	31	236

Table 1. Number of sensors per place and per type of data in Orange4Home.

they are installed in (or to the entire apartment, if it is a global information about the home). A total of 236 different data sources are present in the dataset.

OpenHAB⁶ is used to collect data from all those sensors during the experiment. A MySQL⁷ database is used to persist the data collected by openHAB. Videos from cameras placed on the ceiling of most rooms were also recorded to serve as ground truth; they were used to correct the few labeling mistakes made during the experiment, but are not provided in this dataset.

We present in Table 1 the number of sensors in each place of Orange4Home and for each of the 4 main types of data produced by sensors: binary (door opening, presence, switches, etc.), integers (total cold water consumption, appliance power, humidity, etc.), real numbers (luminosity, voltage, CO₂ levels, noise levels, etc.) and categorical data (weather, heater modes, AC modes, wind direction, etc.). We can see that there are large variations in the number of sensors per place, and on the types of sensors per place, which is expected because each place has a different use and thus different objects that need to be instrumented. We can also observe that all 4 types of data are well-represented, making Orange4Home a dataset of truly heterogeneous sensors.

We present on Fig. 5 and Fig. 6 the number of presence detections and the power consumption of the computer in the Office during the first 5 days of experiment. We can observe that both sensors provide valuable information for activity recognition: both sensors capture correctly that something is happening in the Office on mornings and afternoons, while nothing happens during the lunch break (there are data during the lunch break on Friday because the lunch break on Friday is shorter, and thus does not cover an entire hour interval). We expect that data from all sensors is sufficiently rich to capture all activity classes performed during this experiment.

⁶ <https://www.openhab.org>

⁷ <https://www.mysql.com>

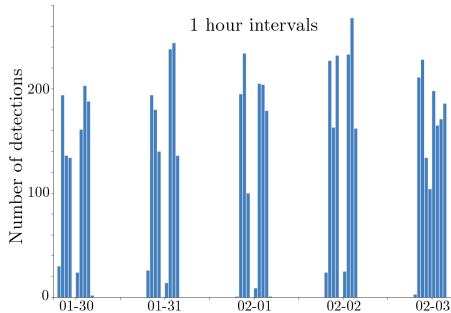


Fig. 5. Presence in the Office during the first week (each box is a 1 hour interval).

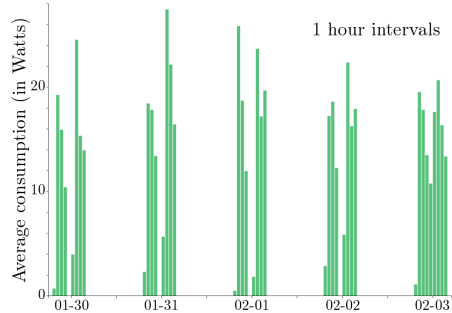


Fig. 6. Power consumption of the computer in the Office during the first week (each box is a 1 hour interval).

There are no body-worn sensors in this dataset. We indeed believe that body-worn sensors are not a realistic data source to use in a smart home system intended for the general public, where data collection should be as seamless as possible (as opposed, for example, to healthcare or wellness systems, where trading seamlessness for more health-related data sources is desirable). Smartphones could be an example of realistic body-worn data sources; however they are not necessarily always carried or even owned by an occupant, and are thus not included in this dataset.

4.3 *In Situ* Labeling

Labeling of activities was performed in real time by the occupant, using an Android application on a smartphone that the occupant carried throughout the home. This application sent events of virtual sensors (one for labels and one for comments) to OpenHAB’s APIs through WiFi.

This application allows the occupant to select the room they are in, and then selected the activity they will perform (the set of activities being restricted to those that can happen in the selected place). The occupant can then press the “*START*” button before beginning their activity, and the “*STOP*” button (which appears in a modal window) once the activity ended. The application also allows the occupant to send comments through the “*ERROR*” button. This was used to comment any unexpected event (such as the occupant pressing “*STOP*” later than the actual end of the activity), in order to greatly simplify the task of fixing labeling issues (which there were very few of) after the experiment.

We present in Table 2 the number of instances of each class of activity that was labeled during the 20 working days of data collection. We can observe that some activity classes are way more numerous than others (e.g. *Computing* in Office compared to *Cleaning* in Bathroom), which is also a reality in a real home setting. No activities other than those presented in Table 2 were performed (the occupant could select an activity named *Other* in the labeling application, but

Place	Activity	Number of instances	Total
Entrance	Entering	21	42
	Leaving	21	
Kitchen	Preparing	19	61
	Cooking	19	
	Washing the dishes	19	
	Cleaning	4	
Living room	Eating	19	71
	Watching TV	18	
	Computing	15	
	Cleaning	19	
Toilet	Using the toilet	8	8
Staircase	Going up	57	114
	Going down	57	
Bathroom	Using the sink	38	70
	Using the toilet	9	
	Showering	19	
	Cleaning	4	
Office	Computing	46	64
	Watching TV	14	
	Cleaning	4	
Bedroom	Dressing	30	63
	Reading	15	
	Napping	15	
	Cleaning	3	
Total			493

Table 2. Number of instances of each class of activity.

this was never needed). 493 instances of activities were performed in total during roughly 180 hours of experiment, for a total of 21 MB of data (in a MySQL database dump format).

4.4 Postprocesses

Electrical consumption data for the connected plugs of the computer and the TV, in both the living room and the office (5 plugs in total), have been simulated for the first day (30th of January) of the experiment, using averaged data from other days of the experiment. Those plugs were indeed non-functional this day. Data from those plugs were also filtered to remove sporadic outliers most likely induced by intrinsic problems with the plugs.

Events are persisted at most every second for a sensor; therefore, if two events happened during the same second (such as pressing and releasing a switch), the first one was persisted one second earlier than the second event.

Apart from the previous two points, no other postprocesses were applied to the data. In particular, we report no missing sensor values for the entirety of the experiment, thanks to the constant and extensive work done in maintaining the experimental apartment for the sake of the Amigual4Home project.

5 Conclusions

We presented in this paper a dataset of labeled activities of daily-living of an occupant in a fully-integrated, instrumented smart home. This dataset spans 20 working days of realistic routines and contains 236 heterogeneous data sources capturing an extensive amount of events happening in the home. Experimenting in a home that the occupant is not used to forced us to carefully craft routines of activities (the process of which we have reported here) such that the occupant feels comfortable and confident enough to act as if they were in their own home.

We believe this dataset can be used as a realistic benchmark for various different kinds of algorithmic problems we are facing in smart homes: activity recognition, activity segmentation, etc. In particular, the length of the experiment as well as the established routine and its variability makes this dataset a good benchmark for activity prediction in smart homes for different time scales.

Acknowledgments. We thank Nicolas Bonnefond and Stan Borkowski for their technical and organizational help. This work benefited from the support of the French State through the *Agence Nationale de la Recherche* under the Future Investments program referenced ANR-11-EQPX-0002.

References

1. Alemdar, H., Ertan, H., Incel, O.D., Ersoy, C.: Aras human activity datasets in multiple homes with multiple residents. In: 7th International Conference on Pervasive Computing Technologies for Healthcare. pp. 232–235 (2013)
2. Cook, D.J., Krishnan, N.C., Rashidi, P.: Activity discovery and activity recognition: A new partnership. *IEEE transactions on cybernetics* 43(3), 820–828 (2013)
3. Crowley, J.L., Coutaz, J., Rey, G., Reignier, P.: Perceptual components for context aware computing. In: International conference on ubiquitous computing. pp. 117–134. Springer (2002)
4. Cumin, J., Lefebvre, G., Ramparany, F., Crowley, J.L.: Human activity recognition using place-based decision fusion in smart homes. In: International and Interdisciplinary Conference on Modeling and Using Context. pp. 137–150. Springer (2017)
5. Li, K., Fu, Y.: Prediction of human activity by discovering temporal sequence patterns. *IEEE transactions on pattern analysis and machine intelligence* 36(8), 1644–1657 (2014)
6. Ordóñez, F.J., Roggen, D.: Deep convolutional and lstm recurrent neural networks for multimodal wearable activity recognition. *Sensors* 16(1), 115 (2016)
7. Roggen, D., Calatroni, A., Rossi, M., Holleczeck, T., Förster, K., Tröster, G., Lukowicz, P., Bannach, D., Pirkl, G., Ferscha, A., et al.: Collecting complex activity datasets in highly rich networked sensor environments. In: Networked Sensing Systems (INSS), 2010 Seventh International Conference on. pp. 233–240. IEEE (2010)
8. Tapia, E.M., Intille, S.S., Larson, K.: Activity recognition in the home using simple and ubiquitous sensors. In: *Pervasive*. vol. 4, pp. 158–175. Springer (2004)
9. Van Kasteren, T., Englebienne, G., Kröse, B.J.: Transferring knowledge of activity recognition across sensor networks. In: *Pervasive*. vol. 10, pp. 283–300. Springer (2010)
10. Youngblood, G.M., Cook, D.J.: Data mining for hierarchical model creation. *IEEE Transactions on Systems, Man, and Cybernetics* 37(4), 561–572 (2007)