# Learning multimodal representations for sample-efficient recognition of human actions*

Miguel Vasco[1]     Francisco S. Melo[1]     David Martins de Matos [1]     Ana Paiva[1]     Tetsunari Inamura[2]

*Abstract*— **Humans interact in rich and diverse ways with the environment. However, the representation of such behavior by artificial agents is often limited. In this work we present *motion concepts*, a novel multimodal representation of human actions in a household environment. A motion concept encompasses a probabilistic description of the kinematics of the action along with its contextual background, namely the location and the objects held during the performance. Furthermore, we present Online Motion Concept Learning (OMCL), a new algorithm which learns novel motion concepts from action demonstrations and recognizes previously learned motion concepts. The algorithm is evaluated on a virtual-reality household environment with the presence of a human avatar. OMCL outperforms standard motion recognition algorithms on an one-shot recognition task, attesting to its potential for sample-efficient recognition of human actions.**

## I. INTRODUCTION

Humans are able to interact with their environment in rich and diverse ways. Such richness and variety make it impossible to program an artificial agent that is able to recognize all possible actions performed by a human user. One common approach is to program agents to learn to recognize new human actions from demonstrations. However, it is unrealistic to assume that such learning will depend on large amounts of data, as required by many current learning algorithms. Instead, the agent should be able to learn and recognize novel actions from just a few demonstrations provided by the human.

To attain such efficient learning, the learning process should take into account the multimodal information provided by the human to create a rich representation of the novel action. However, the conventional methodology of learning human action representations considering only motion pattern data neglects the rich contextual background of the demonstration. This negligence results in a limited representation of the human action, hindering its recognition and introducing difficulties in the distinction between actions with similar motion patterns.

In this work, we address the problem of learning and recognizing human actions, from few demonstrations provided by a human in a household environment. We propose a novel representation for multiple demonstrations of a given action, named *motion concept*. The motion concept encompasses a probabilistic motion primitive description of the motion patterns observed, augmenting it with their contextual background information, namely the location of the action and the objects used during the demonstrations. Moreover, the motion concept takes into account information provided directly through interaction with a human and allows the agent to reason on the importance of each contextual feature for its recognition.

Furthermore, we present the Online Motion Concept Learning (OMCL) algorithm, responsible for the creation of new motion concepts through interaction with a human user. The algorithm is able to recognize motion concepts from a single training demonstration and continuously update motion concepts as more demonstrations are provided. We evaluate the algorithm's performance on an offline "one-shot" motion recognition task, showing the importance of contextual information for the recognition of motion concepts built from a single training demonstration. The obtained results attest to the potential of OMCL for sample-efficient recognition of human actions.

## II. RELATED WORK

The question of learning motion and action representations has been addressed in literature, in part due to the widespread availability of low-cost motion sensing devices [1]. Several representations have been proposed to model human action based on motion data. Xia *et al.* [2] propose a view-invariant action representation based on histograms of 3D joint position, in relation to a fixed coordinate system, obtained from Kinect depth maps. The temporal evolution of these representations are modelled according to discrete HMMs. The authors in [3] propose a novel feature for human action recognition based on the differences in the position of joints in the skeletal model of the human, employing a naive-bayes-nearest-neighbor (NBNN) classifier for recognition of the action classes. The authors in [4] propose an interpretable representation of an action based on the sequence of joints in the skeletal model of the human which, at each time instant, are considered to be the most informative of a given demonstration. The informative criteria are based on predefined measures, such as the mean and variance of joint angle trajectories. Vemulapalli *et al.* [5] model an action as a curve in the Lie group manifold. The curve of each action is generated based on a novel skeletal-based representation that explicitly models the geometric relationships between various body parts using rotations and translations in 3D space. However, all the presented representations of human

actions are built solely resorting to motion data. As such, they neglect the rich contextual background of an action, which is fundamental for the distinction of action classes with similar motion patterns.

Moreover, several deep-learning frameworks have been recently proposed for motion and action recognition. Du *et al.* [6] proposed a hierarchical recurrent neural-network (RNN) framework for skeleton based action recognition, in which the human skeleton is divided accordingly to the human's physical structure. Multiple bidirectional RNNs (BRRN) are trained for each segmented section of the skeleton model, and their output is fused hierarchically by the upper-layers of the framework. Simonyan *et al.* [7] propose a convolutional network architecture for action recognition in video, that incorporates spatial and temporal networks. While these architectures obtain impressive recognition results, their requirement of large amounts of training data make it unsuitable for the recognition of novel actions from few demonstrations.

Multimodal approaches to the creation of action representations have also been explored in literature. The authors in [8] represent actions as an ensemble model and have proposed novel features for depth data which capture human motion and human-object interaction data from a demonstration. Using image data, Yao *et al.* [9] learn action representations composed of attributes, words that describe the properties of human actions, and action-parts, the objects and poselets that are related to the actions. The authors in [10] model an action by integrating multiple feature channels from several entities (such as objects, scenes and people), extracted from video sequences. The representations are obtained through a "multiple instance learning" (MIL) model, where a given action label is associated with a group of instances.
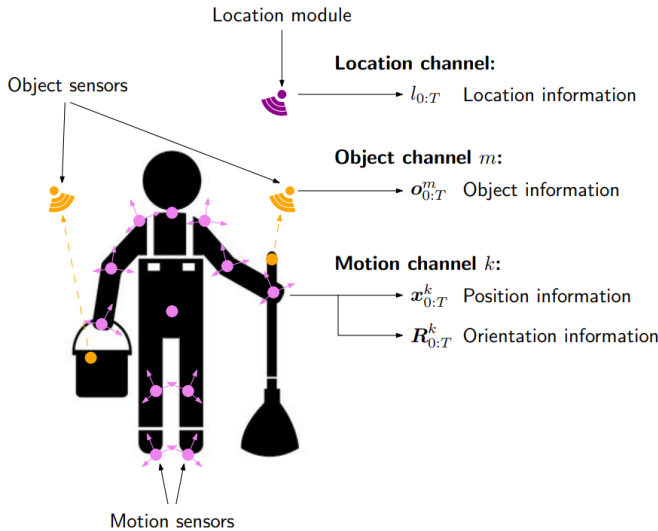


Fig. 1: Depiction of the setup used throughout the paper. Sensors correspond to input channels and provide streams of observations of length $T + 1$

The importance of the action representation and recog-

nition fields is attested by the extensive literature on the matter. Yet, the question regarding the creation of multimodal action representations from few demonstrations still requires addressing.

## III. CONCEPTUAL REPRESENTATION OF AN ACTION

### A. The setup

We consider the following setup for learning from demonstration. A human user demonstrates an action, which may involve interaction with objects in the environment. The environment comprises a number of locations of interest, and the human user may be in any of these locations at the time of demonstration.

We assume that the environment is engineered with a number of sensors, providing information regarding the *location* and *pose* of the human user as well as the *objects* that the user interacts with. In this work, we are not concerned with the actual sensing and admit that the system may include both "internal" sensors (e.g., data gloves to capture pose information, etc.) or external (e.g., cameras or optical trackers). In practical terms, the sensors act as input channels for the system, and as such we henceforth refer to sensors generally as a *channels*. For example, a sensor deployed to provide object information is referred simply as an *object channel*, and sensors deployed to track the human pose are referred as *motion channels* (see Fig. 1 for an illustration). Finally, the location of the user in the environment is provided by a dedicated sensing module, referred to as the *location channel*. A demonstration by a human user yields a number of data streams arising from the different input channels. In particular,

- Each motion channel $k, k = 1, ..., K$ provides two streams of length $T$, $x_{0:T}^k$ and $R_{0:T}^k$, where each $x_t^k$ indicates the position of a body element (joint, limb) at time step $t, t = 0, ..., T$, measured with respect to a common fixed world frame, and each $R_t^k$ is a rotation matrix representing the orientation of that same body element at time step $t$;

- Each object channel $m, m = 1, ..., M$, provides one stream of length $T$, $o_{0:T}^m$, where each individual observation $o_t^m$ corresponds to a binary vector indicating the objects (from a predefined finite set of objects $\mathcal{O}$) that the user is interacting with at time step $t, t = 0, ..., T$, according to channel $m$;

- Finally, the location module provides a stream of length $T$, $l_{0:T}$, where $l_t$ indicates the location of the user at time step $t$. We assume that the location of the user takes values in a finite set $\mathcal{L}$ of possible locations.

Learning a representation of an action will consist of taking the streams from the different input channels and compile them into a unique, compact representation that we refer to as a motion prototype, described in the continuation.

## B. Motion prototype

The central constituent in our proposed action representation is the *motion prototype*, providing a compact representation for a single demonstration of an action by a human user. In particular, motion prototypes capture in a probabilistic manner motion information (extracted from the motion channels) and object and location information.

Formally, we represent a motion prototype as a tuple $P = (\tau, \rho, \lambda)$, where $(\rho, \lambda)$ summarize the associated context information - namely object and location information - and $\tau$ summarizes the motion observed in the demonstration. Specifically,

- $\rho = \{\rho_m, m = 1, ..., M\}$, where $M$ is the total number of object channels. For every object $o \in \mathcal{O}$,

$$\rho_m(o) = \mathbb{P}\left[o_{t,o}^m = 1, t = 0, ..., T\right]$$

where $o_{t,o}^m$ is a random variable indication whether object $o$ was observed in object channel $m$ at time step $t$. In other words, in our proposed representation we assume that the observation of an object $o \in \mathcal{O}$ in channel $m$ at any moment during the human demonstration can be described probabilistically as a Bernoulli random variable with parameter $\rho^m(o)$.

- For every location $l \in \mathcal{L}$,

$$\lambda(l) = \mathbb{P}\left[l_t = l, t = 0, ..., T\right]$$

where $l_t$ is a random variable indicating the location of the human demonstrator at time step $t$. In other words, in our representation we assume that the location of the human user during the demonstration can be described probabilistically as a categorical distribution with parameters $\lambda(l), l \in \mathcal{L}$.

Finally, we have that $\tau = \{\tau_k, k = 1, ..., K\}$, where $K$ is the number of motion channels and each $\tau_k$ is a sequence of *motion primitives* $\{\phi_n, n = 1, ..., N\}$. The concept of motion primitive has been widely explored both to describe animal motion and to represent robot motion [11], [12]. For our purposes, a motion primitive $\phi_n$ is a probability distribution over the space of trajectories. In other words, given an arbitrary trajectory $(x_{0:T}, R_{0:T})$,

$$\phi_n(x_{0:T}, R_{0:T}) = \mathbb{P}\left[x_{0:T}^n = x_{0:T}, R_{0:T}^n = R_{0:T}\right]$$

For the purpose of learning and recognition, it is convenient to treat an action not as comprising a single trajectory $(x_{0:T}, R_{0:T})$ but, instead, as a sequence of smaller trajectories,

$$\left\{(x_{0:t_1}, R_{0:t_1}), (x_{t_1:t_2}, R_{t_1:t_2}), ..., \left(x_{t_{N-1}:t_N}, R_{t_{N-1}:t_N}\right)\right\}$$

which are then encoded as a sequence of motion primitives $\{\phi_n, n = 1, ..., N\}$, each $\phi_n$ providing a compact description of $(x_{t_{n-1}:t_n}, R_{t_{n-1}:t_n})$. Each motion primitive $\phi_n$ is selected among a library $\Phi$ of available motion primitives to maximize the
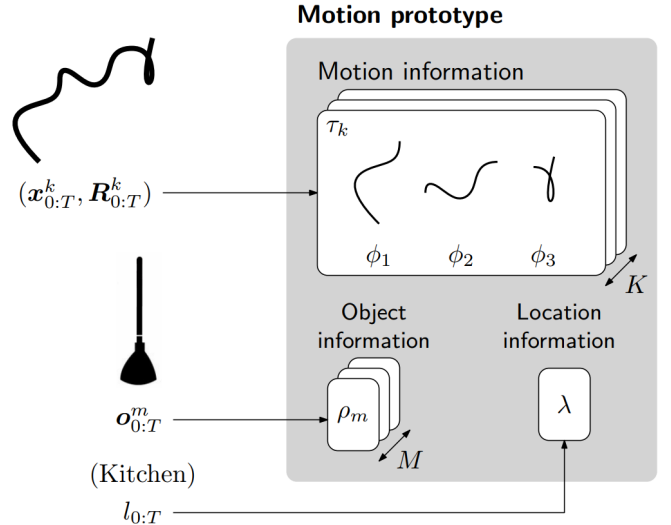


Fig. 2: Summary representation of a motion prototype.

likelihood of the observed trajectory, i.e.,

$$\phi_n = \underset{\phi \in \Phi}{\arg\max} \, \phi\left(x_{t_{n-1}:t_n}, R_{t_{n-1}:t_n}\right) \tag{1}$$

Summarizing, a motor prototype compactly encodes a demonstration of an action in the form of a tuple $(\tau, \rho, \lambda)$, where $\tau$ is a collection of trajectories (one for each motion channel), each represented as a sequence of motor primitives; $\rho$ is a collection of probability distributions (one for each object channel), describing how the human interacted with each object in the environment; and $\lambda$ is a probability distribution describing where the human was located during the demonstration (see Fig. 2).

## C. Motion concept

It is possible for a single action to be performed in multiple different ways. A motion prototype, while providing a convenient representation for a single demonstration (and corresponding context), is insufficient to capture the diversity that a broader notion of "action" entails.

We introduce *motion concept* as a higher-level representation of an action. A motion concept seeks to accommodate the different ways by which an action may be performed while, at the same time, encode distinctive aspects that are central in recognizing such action. For example, to distinguish actions such as "waving goodbye" and "washing a window", it is important to note that the latter involves interaction with an object (such a sponge) while the first does not.

Formally, a motion concept consists of a tuple $\mathcal{M} = (\mathcal{P}, \eta, k_\rho, k_\lambda)$ used to represent some action $a$, where

- $\mathcal{P} = \{P_1, ..., P_\ell\}$, where each $P_i$ is a motion prototype describing one possible way by which the action $a$ can be performed;

- $\eta$ is a *designation* (a "name") provided by the user to refer to the action $a$ - for example, it may consist of
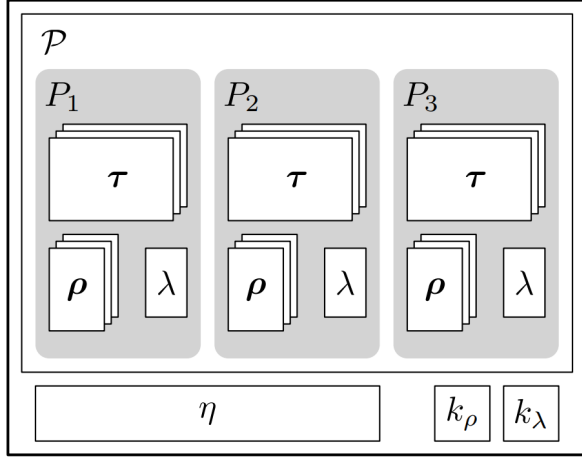
Fig. 3: Schematic representation of a motion concept.
.

a label for $a$ or an utterance that corresponds to the spoken designation of $a$;

- $k_p$ and $k_\lambda$ are two constants used to weight the importance of object information and location information in recognizing the action $a$.

We schematically depict a motion concept in Fig. 3.

## IV. LEARNING MOTION CONCEPTS

A motion concept provides a multimodal representation of a given action demonstrated by a human. However, in order for an agent to learn motion concepts from demonstration a novel algorithm is required. For that purpose, we now introduce the OMCL (Online Motion Concept Learning) algorithm, designed to construct motion concepts from the demonstration data provided by a human user.

The OMCL algorithm proceeds can roughly be understood as working on two different abstraction levels. At a lower level, OMCL takes the data from a single demonstration and constructs a motion prototype from such data. At a higher level, OMCL combines information from multiple demonstrations to build a motion concept that potentially contains multiple motion prototypes. We now detail OMCL at each of these two abstraction levels.

### A. Learning a motion prototype from data

Given a demonstration $(x_{0:T}, R_{0:T})$, OMCL starts by segmenting the single trajectory into multiple sub-trajectories

$$\left\{ (x_{0:t_1}, R_{0:t_1}), ..., (x_{t_{N-1}:t_N}, R_{t_{N-1}:t_N}) \right\} \quad (2)$$

which can be achieved using any segmentation method from the literature - OMCL is agnostic to the particular segmentation method used. OMCL uses online kernel density estimation[1] to construct new motion primitives from sub-trajectory data. The list of segmented sub-trajectories (Eq. 2) is used to update the previous set $\Phi$ of available motion

[1]In our implementation, we use the XOKDE++ algorithm from [13]

primitives, after which each sub-trajectory $\left(x_{t_{n-1}:t_n}, R_{t_{n-1}:t_n}\right)$ is evaluated against the updated $\Phi$ and a primitive $\phi_n$ is selected accordingly to Eq. 1. The resulting sequence of motion primitives, $\{\phi_1, ..., \phi_N\}$, corresponds to the trajectory representation $\tau$, as described in the previous section. Such procedure is repeated for all motion channels.

As for $\rho$, we use standard maximum likelihood estimation to compute the parameters $\rho(o), o \in \mathcal{O}$, from the data $o_{0:T}$, repeating this procedure for all object channels. Finally, we also use maximum likelihood estimation to compute the parameters $\lambda(l), l \in \mathcal{L}$, from the data $l_{0:T}$.

### B. Building a motion concept

From a provided demonstration of action $a$, the OMCL algorithm learns a motion prototype $P_a = (\tau_a, \rho_a, \lambda_a)$. If $P_a$ is the first motion prototype of action $a$ provided, a new motion concept $\mathcal{M}_a$ is built through the following procedure:

- The motion prototype is added to the empty list of prototypes $\mathcal{P}$,

$$\mathcal{P} = \{P_a\}$$

- The importance weights $k_\rho, k_\lambda$ are initialized to predetermined values,

$$k_\rho = k_{\rho,0}, \quad k_\lambda = k_{\lambda,0}$$

- The human is queried for the designation $\eta$ of the action.

The novel motion concept $\mathcal{M}_a$ is added to the current list of built motion concepts $\Sigma = \{\mathcal{M}_1, ..., \mathcal{M}_A\}$, where $A$ is the total number of action classes previously demonstrated. If the motion prototype $P_a$ concerns an action $a$ previously demonstrated, it is then used to update the respective motion concept $\mathcal{M}_a$, following:

- The motion prototype is added to the list of prototypes $\mathcal{P}$,

$$\mathcal{P} = \{P_1, ..., P_a\}$$

- The contextual information of the motion prototype $(\rho_a, \lambda_a)$ is used to update the values of the importance weights $k_\rho, k_\lambda$. If, for the majority of motion prototypes $P_i = (\tau_i, \rho_i, \lambda_i) \in \mathcal{P}$:

$$\arg\max_{o \in \mathcal{O}} \rho_{i,m}(o) = \arg\max_{o \in \mathcal{O}} \rho_{a,m}(o), \forall m \in M$$

we increase the value of $k_\rho$ by a percentage $\alpha_k$ of its value. Otherwise, we decrease it by the same percentage. The same procedure is applied for the update of the $k_\lambda$ weight.

## V. RECOGNIZING MOTION CONCEPTS

Beyond creating and updating motion concepts, the OMCL algorithm is also responsible for the recognition of previously observed actions and for the assessment of the novelty of previously unobserved action classes. Given a motion prototype from an unknown action $P_* = (\tau_*, \rho_*, \lambda_*)$, OMCL compares $P_*$ with the prototypes contained in every motion

concept in the current set $\Sigma$. The cost of assigning $P_*$ to the motion concept $\mathcal{M}_i$ is given by:

$$\mathcal{C}(\mathcal{M}_i, P_*) = \mathcal{C}_\tau(\mathcal{P}, \tau_*) + k_\rho \mathcal{C}_\rho(\mathcal{P}, \rho_*) + k_\lambda \mathcal{C}_\lambda(\mathcal{P}, \lambda_*) \quad (3)$$

where,

- $\mathcal{C}_\tau(\mathcal{P}, \tau_*)$ is the average cost of comparing the sequence of motion primitives in $\tau_*$ to the sequence of motion primitives $\tau$ of every motion prototype $P \in \mathcal{P}$ of $\mathcal{M}_i$. To compute this cost, we use dynamic time warping [14] between the sequences of each motion channel, with 0-1 loss;

- $\mathcal{C}_\rho(\mathcal{P}, \rho_*)$ is the average distance between $\rho_*$ and the collection of object probability distributions $\rho$ of every motion prototype $P \in \mathcal{P}$ of $\mathcal{M}_i$. The algorithm is agnostic to the type of metric used to compute the distance between distributions;

- $\mathcal{C}_\lambda(\mathcal{P}, \lambda_*)$ is the average distance between $\lambda_*$ and the location probability distributions $\lambda$ of every motion prototype $P \in \mathcal{P}$ of $\mathcal{M}_i$;

- $k_\rho, k_\lambda$ are the object and location information weights of $\mathcal{M}_i$.

For recognition purposes, the cost is computed for all motion concepts available in $\Sigma$ and $P_*$ is assigned to the motion concept $\mathcal{M}_R \in \Sigma$ with the lowest assignment cost $\mathcal{C}_R$.

However, the question of the possible class-novelty of the demonstration still requires addressing. Therefore, subsequently, OMCL determines if $P_*$ belongs to the assigned motion concept or if it belongs to a novel action class. The decision takes into account the average wrong cost $\mathcal{C}_W$ of assigning $P_*$ to the remaining motion concepts $\mathcal{M} \in \Sigma \setminus \mathcal{M}_R$. In the case,

$$|\mathcal{C}_R - \mathcal{C}_W| \geq \Delta_\mathcal{C} \times \mathcal{C}_R \quad (4)$$

the provisional assignment of $P_*$ to $\mathcal{M}_R$ is confirmed, where $\Delta_\mathcal{C}$ is a predefined constant. Otherwise, the assignment can be considered quasi-random, and $P_*$ will be used to build a new motion concept (as detailed in the previous section).

## VI. EVALUATION

### A. Experimental Setup

The evaluation of the OMCL algorithm was performed in a virtual-reality (VR) environment. The user interacts with the VR environment using a Oculus Rift headset and hand motion controllers, as shown in Fig. 4. Thus, in this setup, the number of motion channels $K$ is equal to the number of object channels $M$, with $K = M = 3$.

Furthermore, we designed a virtual household environment (as seen in Fig. 5), composed of 4 different sections: Kitchen (K), Living-Room (LR), Dining-Room (DR) and Bathroom (BR). In other words $\mathcal{L} = \{K, LR, DR, BR\}$. Each section contains objects specific of that section (e.g. "Tooth-brush" is contained in the "Bathroom" area) as well as a number of common objects that can be found in multiple sections



Fig. 4: Participant demonstrating an example of the "Vacuum-clean" action, along with the corresponding VR avatar.
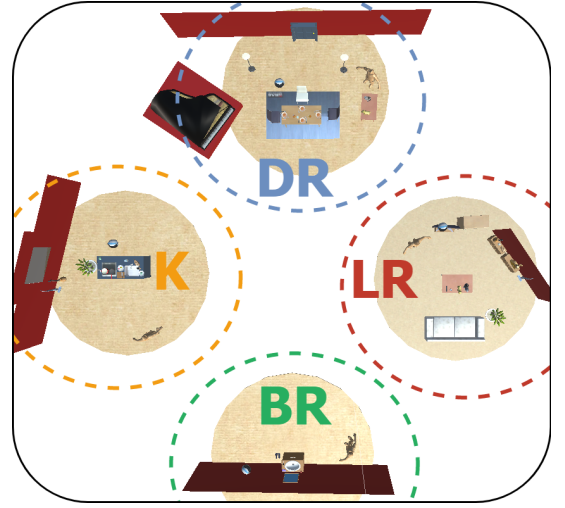


Fig. 5: Virtual household environment with the area of every discrete location discriminated: "DR"-Dining Room, "K"-Kitchen, "BR"-Bathroom, "LR"-Living Room.

of the environment (e.g. "Cup" can be found in "Kitchen", "Living-Room", "Dining-Room").

### B. One-shot recognition (OSR) task

In the one-shot recognition (OSR) task we evaluate the performance of OMCL in the recognition of actions when provided a single training demonstration of each class to create the associated motion concepts. We asked 10 participants to perform, on the virtual household environment, two demonstrations of (randomly-ordered) 22 action classes, after a tutorial period of adaptation to the VR setting. Each action was recorded for 6 seconds, storing the motion data, from the VR headset and motion controllers, and the contextual data (object and location information) of the performance. We provided to the participants no information regarding which objects to use or where to perform the action. The complete list of actions selected for this task is presented

TABLE I: List of action classes performed for the OSR task, along with the most common objects used in the performances and their most common locations in the household environment (following the nomenclature of Fig. 5).

| Motion Class | Location | Objects |
|---|---|---|
| Bow | All | None |
| Comb hair | BR | Hairbrush |
| Cut | K | Knife, Apple, Banana, Pear |
| Drink | All | Mug, Glass, Bottle |
| Eat at Table | DR | Knife, Fork, Chopsticks |
| Fry | K | Frying Pan |
| High-Five | All | Hand |
| Hug | All | Body |
| Knock on door | LR | Door |
| Pet | DR, LR | Cat, Dog |
| Play Guitar | LR | Guitar |
| Play Piano | DR | Piano |
| Shake Hands | All | Hand |
| Stir Pot | K | Spoon, Pot |
| Sweep | K, LR | Broom |
| Throw | All | All |
| Vacuum clean | K, LR | Vacuum-cleaner |
| Wash Hands | BR | Soap |
| Wash Plates | K | Sponge, Dish |
| Wash Window | K | Sponge |
| Wave | All | None |
| Wring Sponge | K | Sponge |

in Table I, along with the most common objects used in the performances and their most common locations in the household environment.

The action classes were chosen due to their simplicity, as complex manipulation of objects in a virtual environment is difficult, and the fact that participants could performed them stationary, to minimize the discomfort of locomotion in virtual space. Moreover, we selected actions with very similar motion patterns but distinct object and location contexts (e.g. "Wash Hands"/"Wash Plates" and "Wave"/"Wash Window" actions) and actions with highly variant motion patterns, object and location contexts (e.g. "Throw" action).

We optimize the values of the $(k_{\lambda,0}, k_{\rho,0})$ parameters of OMCL by grid-search validation, training with one random sample of each action class and evaluating the remaining samples in the training partition of the dataset. The training procedure is repeated 10 times per tuple of parameter values and the optimized values are selected based on the total accuracy of the model. The $\Delta_{\mathcal{C}}$ parameter is optimized following the same grid-search procedure: fixing the values of $(k_{\lambda,0}, k_{\rho,0})$ obtained previously, we build a motion concept from a single randomly-selected training sample of each action class. Subsequently, we evaluate the number of times OMCL assesses the test samples (provided without explicit class labels) as examples of the correct, corresponding, motion concept. The final parameter values are presented in Table II.

The performance of the OMCL algorithm is evaluated against a Gaussian Emission Hidden Markov Model (GMM-HMM), optimized through the same training procedure, yet resorting only to the motion data of the recorded actions. Using the total accuracy of the model as the selection criteria, the optimized number of hidden states in the model

TABLE II: Optimized parameter values of the OMCL algorithm.

| Parameter | Value |
|---|---|
| $k_{\lambda,0}$ | 0.005 |
| $k_{\rho,0}$ | 0.05 |
| $\Delta_{\mathcal{C}}$ | 0.9 |

TABLE III: Accuracy on the one-shot recognition task for the GMM-HMM, OMCL-N and OMCL algorithms.

| GMM-HMM (%) | OMCL-N (%) | OMCL (%) |
|---|---|---|
| $37.6 \pm 21.2$ | $68.8 \pm 19.7$ | $90.5 \pm 20.8$ |

is $h_{\text{HMM}} = 16$ and the optimized number of components in the GMMs is $k_{\text{GMM}} = 3$. Moreover, to fairly compare both algorithms, we include in the evaluation procedure a modified OMCL model (OMCL-N), in which we neglect the contribution of the contextual features (object and location information) to the recognition cost (Eq. 3). In other words, the motion concepts in OMCL-N are built solely considering motion data. In the OSR task, the recognition rates of the GMM-HMM algorithm, OMCL-N and OMCL algorithms in the test partition of the dataset are presented in Table III. Moreover, their confusion matrices are presented in Fig. 6.

In the OSR task, the OMCL-N algorithm significantly outperforms the GMM-HMM algorithm, with an accuracy of $68.8 \pm 19.7\%$ against $37.6 \pm 21.2\%$. This result validates the methodology of solving the recognition problem not through the direct comparison of low-level joint data, which is prone to noise and measurement errors, but through the comparison of previously learned motion primitives, generalized from the data. Yet, due to the diversity of motions patterns for the same action, contextual information still seems fundamental for the one-shot recognition, as the regular OMCL algorithm significant out-performs both methods (with $90.5 \pm 20.8\%$ accuracy rate).

The difference in performance between the algorithms can also be verified by the confusion matrices shown in Fig. 6. OMCL-N (Fig. 6b) presents a significantly more diagonal matrix compared to the matrix of GMM-HMM (Fig. 6a). Yet, the recognition of actions with similar motion patterns is still difficult, as both algorithms are not able to successfully distinguish between the "Wash Hands"/"Wash Dishes" actions (marked in green in Fig. 6) as well as between the "Wave"/"Wash-Window" actions (in blue in Fig. 6). The OMCL algorithm (Fig. 6c) shows a significant improvement in the recognition of the actions classes, indicated by the near-diagonal confusion matrix. Moreover, OMCL is able to distinguish between actions with similar motion patterns ("Wash Window"/"Wash Dishes", "Wave"/"Wash-Window") by taking into account the contextual information of the action (object and location data). However, OMCL is still unable to recognize the action "Throw" (marked in red in Fig. 6) due to the similarity of its motion pattern to the class "High-Five" and the variance of objects used and locations where it can be performed. Indeed, the consideration of
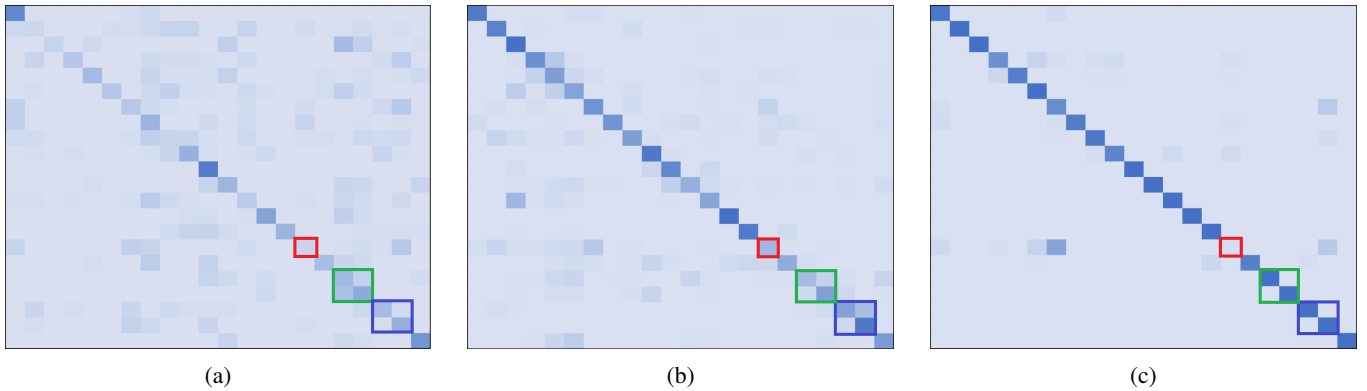
Fig. 6: Confusion Matrices of the GMM-HMM (6a), OMCL-N (6b) and OMCL (6c) algorithms on the one-shot recognition task. We highlight the accuracy on the "Throw" action class (red), the accuracy on the "Wash Hands" and "Wash Dishes" classes (green) and the accuracy on the "Wash Window" and "Wave" classes (blue).

contextual information in the recognition of the "Throw" action seems to worsen the accuracy performance of the algorithm in comparisson with the solely-motion-based version of OMCL. Yet, for the remaining action classes, the contextual information of the actions seem to play a fundamental role in the improvement of the recognition performance of OMCL.

## VII. CONCLUSION

In this paper we present motion concepts, a novel multimodal representation for human actions in a household environment, based on the kinematics of the demonstration, the objects interacted with during the action and the location where it was performed. Moreover, we present OMCL, a new algorithm for the creation and the recognition of motion concepts from demonstrations provided by human users.

We evaluated OMCL on a one-shot recognition task, which showed that the motion concept representation proposed is suitable to be used in action recognition from a single demonstration. Moreover we attest to the importance of contextual information of an action to recognize actions with similar motion patterns. We plan to further evaluate the algorithm on an online motion concept learning task.

The question of learning rich representations of actions in an environment is an ever-evolving subject. We plan to develop further work on the representation of actions performed by multiple agents and the extension of the motion concept representation for task learning. We believe that action representations are fundamental tools for attaining a profound understanding of human behavior in an environment and, ultimately, for the widespread use of artificial agents in household environments.

## REFERENCES

[1] R. Lun and W. Zhao, "A survey of applications and human motion recognition with microsoft kinect," *International Journal of Pattern Recognition and Artificial Intelligence*, vol. 29, no. 05, p. 1555008, 2015.

[2] L. Xia, C.-C. Chen, and J. K. Aggarwal, "View invariant human action recognition using histograms of 3d joints," in *Computer vision and pattern recognition workshops (CVPRW), 2012 IEEE computer society conference on*. IEEE, 2012, pp. 20–27.

[3] X. Yang and Y. L. Tian, "Eigenjoints-based action recognition using naive-bayes-nearest-neighbor," in *Computer vision and pattern recognition workshops (CVPRW), 2012 IEEE computer society conference on*. IEEE, 2012, pp. 14–19.

[4] F. Ofli, R. Chaudhry, G. Kurillo, R. Vidal, and R. Bajcsy, "Sequence of the most informative joints (smij): A new representation for human skeletal action recognition," *Journal of Visual Communication and Image Representation*, vol. 25, no. 1, pp. 24–38, 2014.

[5] R. Vemulapalli, F. Arrate, and R. Chellappa, "Human action recognition by representing 3d skeletons as points in a lie group," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 588–595.

[6] Y. Du, W. Wang, and L. Wang, "Hierarchical recurrent neural network for skeleton based action recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1110–1118.

[7] K. Simonyan and A. Zisserman, "Two-stream convolutional networks for action recognition in videos," in *Advances in neural information processing systems*, 2014, pp. 568–576.

[8] J. Wang, Z. Liu, Y. Wu, and J. Yuan, "Mining actionlet ensemble for action recognition with depth cameras," in *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. IEEE, 2012, pp. 1290–1297.

[9] B. Yao, X. Jiang, A. Khosla, A. L. Lin, L. Guibas, and L. Fei-Fei, "Human action recognition by learning bases of action attributes and parts," in *Computer Vision (ICCV), 2011 IEEE International Conference on*. IEEE, 2011, pp. 1331–1338.

[10] N. Ikizler-Cinbis and S. Sclaroff, "Object, scene and actions: Combining multiple features for human action recognition," in *European conference on computer vision*. Springer, 2010, pp. 494–507.

[11] T. Flash and B. Hochner, "Motor primitives in vertebrates and invertebrates," *Current opinion in neurobiology*, vol. 15, no. 6, pp. 660–666, 2005.

[12] J. Kober and J. Peters, "Learning motor primitives for robotics," in *Robotics and Automation, 2009. ICRA'09. IEEE International Conference on*. IEEE, 2009, pp. 2112–2118.

[13] J. Ferreira, D. M. de Matos, and R. Ribeiro, "Fast and extensible online multivariate kernel density estimation," *arXiv preprint arXiv:1606.02608*, 2016.

[14] M. Müller, "Dynamic time warping," *Information retrieval for music and motion*, pp. 69–84, 2007.