

# Simultaneous Learning from Human Pose and Object Cues for Real-Time Activity Recognition

Brian Reily<sup>1</sup>, Qingzhao Zhu<sup>1</sup>, Christopher Reardon<sup>2</sup>, and Hao Zhang<sup>1</sup>

**Abstract**—Real-time human activity recognition plays an essential role in real-world human-centered robotics applications, such as assisted living and human-robot collaboration. Although previous methods based on skeletal data to encode human poses showed promising results on real-time activity recognition, they lacked the capability to consider the context provided by objects within the scene and in use by the humans, which can provide a further discriminant between human activity categories. In this paper, we propose a novel approach to real-time human activity recognition, through simultaneously learning from observations of both human poses and objects involved in the human activity. We formulate human activity recognition as a joint optimization problem under a unified mathematical framework, which uses a regression-like loss function to integrate human pose and object cues and defines structured sparsity-inducing norms to identify discriminative body joints and object attributes. To evaluate our method, we perform extensive experiments on two benchmark datasets and a physical robot in a home assistance setting. Experimental results have shown that our method outperforms previous methods and obtains real-time performance for human activity recognition with a processing speed of  $10^4$  Hz.

## I. INTRODUCTION

Real-time human activity recognition is an essential capability of robots in human-centered robotics applications, such as assisted living, service robotics, human-robot teaming, and human-robot interaction [1]–[4]. It allows intelligent robots to understand human behaviors in a timely manner in order to effectively assist and interact with humans. Human activity recognition by robots in the real world is a difficult problem, complicated by both variations in human appearances and poses, and by challenges such as illumination changes or occlusions. Given these challenges, it is important for a robot to extract as much relevant information as possible from sensory observations. For example, as illustrated in Figure 1, this information can consist of humans themselves, such as human poses encoded by the human skeleton representation, and the context from the objects in the environment and objects that the human is interacting with, which provide additional cues to recognize activities. Moreover, in most real-world robotics applications, and especially in time-critical scenarios, activity recognition must occur in real-time so that a robot can promptly interact with and assist humans.

Due to the importance of activity recognition, many methods have been introduced over the past few decades [5]–[8].

\*This work was partially supported by NSF IIS-1942056, ARL DCIST CRA W911NF-17-2-0181, USAFA FA7000-18-2-0016, and ARO W911NF-17-1-0447.

<sup>1</sup>Brian Reily, Qingzhao Zhu and Hao Zhang are with Human-Centered Robotics Laboratory at the Colorado School of Mines, Golden, CO, 80401, USA. Email: {breily, zhuqingzhao, hzhang}@mines.edu.

<sup>2</sup>Christopher Reardon is with U.S. Army Research Laboratory, Adelphi, MD, 20783, USA. Email: christopher.m.reardon3.civ@mail.mil.



Fig. 1. A motivating example of integrating observations of human poses and of the objects involved in the task to understand human activities. The objects within the scene and in use by the humans provide additional cues beside human poses to recognize activities.

Especially, techniques based on skeletal data from structured-light cameras [9] have attracted increasing attention, due to skeletal data’s real-time performance and invariance to viewing distances and angles. For example, the methods can be implemented based on hand-crafted skeletal features [10], a concatenation of multiple types of features [11]–[13], and learning skeleton-based representations, e.g., by sparse optimization [14], [15] or deep learning [16], [17]. These methods generally have the limitation of not learning from the context of objects that the human is interacting with. Although several methods used object information [18], [19], they require explicit knowledge of the objects, such as object affordances that are typically manually defined to describe how an object can be interacted with. Moreover, previous methods cannot estimate the importance of the objects in recognizing human activities.

In this paper, we propose a principled method for real-time human activity recognition based on learning simultaneously from observations of humans and objects. Our approach formulates activity recognition as a regression-like optimization problem, and applies structured norms as regularization terms to promote sparsity and identify discriminative skeletal joints and object attributes. This formulation is inspired by the fact that many activities rely solely on a subset of joints (e.g., a waving activity uses only joints in the arm, not in the legs), or can be recognized based on context of objects in the scene (e.g., reading a book and typing on a laptop at a table appear similar if only the human pose is considered). By learning the importance of both to human activities, the proposed method is able to identify and integrate more relevant information to improve activity recognition accuracy. Because classification is integrated in our regression-like convex objective function

(i.e., no separate classifier is needed), our approach is capable of operating in real-time, which makes it suitable for robotics applications with real-time needs.

This paper has two major contributions:

- 1) We propose a novel principled method that formulates human activity recognition as simultaneously learning from human and object observations based on a unified regularized optimization framework. The method identifies both discriminative skeletal joints and discriminative object attributes, and integrates classification with sparse representation learning in order to enable a high processing speed for real-time recognition.
- 2) We implement a new iterative optimization algorithm to solve the formulated regularized optimization problem that has dependent model parameters, which holds a theoretical guarantee to find the optimal solution.

## II. RELATED WORK

Activity recognition has been shown to be a critical ability for robots to work with people in real-world human-centered robotics applications [20]–[25]. This section provides a brief review of existing methods of skeleton-based representations and object-assisted activity recognition.

### A. Skeleton-Based Representations for Activity Recognition

Among diverse human representations, skeleton-based representations attracted extensive attentions since the availability of structured-light or color-depth cameras. Skeleton-based representations can be based on joint positions, displacement, orientation, and a combination of multiple joint features [5].

Relative *spatial displacement* between a pair of body joints is one of the most commonly applied skeleton-based features. For example, the normalized joint positions were employed to compute pairwise relative distances between joints as features to categorize activities using extreme learning machines [26]. Euclidian distances between joint pairs were computed as skeletal features to recognize activities [12]. Rahmani *et al.* [27] chose a reference joint, such as the torso center, and computed relative distances to other body joints. In addition, *orientation* of a segment between two joints in space or time is widely studied. Boubou and Suzuki [28] implemented the histogram of oriented velocity vector features that calculate joint velocities between frames and use distributions of joint orientations to classify human activities. Yang and Tian [29] designed a descriptor to include joint orientation differences between frames as features. *Joint positions* were also directly applied as input into long short-term memory networks and recurrent neural networks to recognize activities [30]. Recent methods *integrated multiple features*. Luo *et al.* [31] fused sparse-coding skeleton features with a representation named center-symmetric motion local ternary pattern, which extracts spatial and temporal gradients as features. A learning-based method is proposed in [14] to estimate the weights of feature modalities for activity recognition.

While most previous methods only consider skeletal information to build representations, our proposed method focuses

on integrating cues from both humans and objects for activity recognition, in a unified optimization framework.

### B. Object-Assisted Activity Recognition

A few methods were implemented to take into account of objects for activity recognition. Koppula *et al.* [19], [32] used human and object trajectories as a particle of an anticipatory temporal conditional random field for activity recognition. Yu *et al.* [13] combined human body joint and object positions as representations and used a boosting technique to identify human activities. Wei *et al.* [18] implemented a human-object interaction model that combines spatiotemporal human body joint displacements with object recognition and localization in sequence of frames. Hu *et al.* [33] implemented a heterogeneous method for joint features learning that concatenates spatial displacements of the skeleton data over a sequence of frames, as well as color and depth patterns and their gradients around each joint. Besides extracting joints and objects as independent features, probabilistic models based on human poses and object interactions were also designed for activity recognition [34].

Previous object-assisted activity recognition methods cannot estimate the importance of objects and skeletal features in recognizing activities. Several methods [32], [34] also require predefined knowledge about the objects such as object affordance. Our proposed method provides the capability of not only integrating human and object cues, but also estimate their importance when recognizing human activities.

## III. THE PROPOSED JOINT LEARNING APPROACH

*Notation.* In this paper, matrices are denoted using boldface uppercase letters, and vectors using boldface lowercase letters. For a matrix  $\mathbf{M} = \{m_{ij}\} \in \mathbb{R}^{p \times q}$ , we refer to its  $i$ -th row as  $\mathbf{m}^i$  and its  $j$ -th column as  $\mathbf{m}_j$ , and  $m_{ij}$  as the element in the  $i$ -th row and the  $j$ -th column. The Frobenius norm of a matrix  $\mathbf{M}$  is computed as  $\|\mathbf{M}\|_F = \sqrt{\sum_{i=1}^p \sum_{j=1}^q m_{ij}^2}$ . For dimensionality,  $d_T^j$  represents the dimensionality of the  $j$ -th body joint of the human and  $d_O^m$  indicates the dimensionality of the  $m$ -th attribute modality of an object.

### A. Problem Formulation

We assume that the input data instance set,  $\mathbf{X} = \{\mathbf{T}, \mathbf{O}\}$ , consists of paired observations of a human and objects.  $\mathbf{T} = [\mathbf{t}_1, \dots, \mathbf{t}_N] \in \mathbb{R}^{d_T \times N}$  denotes the matrix of observations of the human, where  $\mathbf{t}_i \in \mathbb{R}^{d_T}$  is the feature vector describing the human's  $J$  joints in the  $i$ -th data instance. Subsections in  $\mathbf{t}_i$  describe individual joints, with  $\mathbf{t}_i^j \in \mathbb{R}^{d_T^j}$  describing the  $j$ -th joint in the  $i$ -th data instance. Each body joint is described by its displacement from the center of the body (typically, the torso 'joint' in many skeletal representations). Observations of the objects are encoded in the matrix  $\mathbf{O} = [\mathbf{o}_1, \dots, \mathbf{o}_N] \in \mathbb{R}^{d_O \times N}$ , where  $\mathbf{o}_i \in \mathbb{R}^{d_O}$  is the feature vector representing all the objects in the  $i$ -th data instance. Each object feature vector is sub-divided to encode multiple objects, each with  $M$  attribute modalities, such that  $\mathbf{o}_i^{o_m} \in \mathbb{R}^{d_O^m}$  represents the features describing the  $m$ -th attribute modality of the  $o$ -th object in the  $i$ -th data instance.

Activity category labels for each training data instance are denoted in the category indicator matrix  $\mathbf{Y} = [\mathbf{y}^1; \dots; \mathbf{y}^N] \in \mathbb{R}^{N \times C}$ , where  $\mathbf{y}^i \in \mathbb{R}^C$  denotes the category indicator vector for the  $i$ -th data instance and  $C$  denotes the number of human activity categories. Specifically,  $y_{ic}$  indicates the probability that the  $i$ -th data instance  $\mathbf{x}_i = \{\mathbf{t}_i, \mathbf{o}_i\}$  belongs to the  $c$ -th activity category. In the training phase, these probabilities are either 0 (if the data instance does not belong to that category) or 1 (if the data instance belongs to that category).

We formulate human activity recognition based upon both skeletal observations and object observations as a regression-like optimization problem:

$$\min_{\mathbf{W}, \mathbf{U}} \|\mathbf{T}^\top \mathbf{W} + \mathbf{O}^\top \mathbf{U} - \mathbf{Y}\|_F^2 \quad (1)$$

where  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_C] \in \mathbb{R}^{d_T \times C}$  represents a weight matrix indicating the importance of  $\mathbf{T}$  to the activity category labels, and  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_C] \in \mathbb{R}^{d_O \times C}$  is a weight matrix doing the same for  $\mathbf{O}$ .  $\mathbf{w}_c \in \mathbb{R}^{d_T}$  represents weights of joints with respect to  $c$ -th category, with subsections  $\mathbf{w}_c^j \in \mathbb{R}^{d_T^j}$  representing the weights of the  $j$ -th joint to the  $c$ -th category. Similarly,  $\mathbf{u}_c \in \mathbb{R}^{d_O}$  represents weights of object attributes with respect to  $c$ -th category, with subsections  $\mathbf{u}_c^{o_m} \in \mathbb{R}^{d_O^m}$  representing weights of the  $m$ -th attribute of the  $o$ -th object to the  $c$ -th category.

### B. Learning Discriminative Joints and Object Attributes

When recognizing activities, specific body joints and object attributes are typically more discriminative. For example, joints in the arm are much more important when a human is retrieving an object from a shelf than joints in the leg would be. Similarly, attributes describing the object being retrieved would allow a robot to understand, for instance, whether the human is about to work on a laptop or read a book.

In order to identify discriminative joints, we introduce the *skeletal norm* on the weight matrix  $\mathbf{W}$ , defined as:

$$\|\mathbf{W}\|_S = \sum_{c=1}^C \sum_{j=1}^J \|\mathbf{w}_c^j\|_2 \quad (2)$$

This skeletal norm enforces the  $\ell_2$ -norm within a joint feature and the  $\ell_1$ -norm between joints in order to force sparsity and identify discriminative joints (Figure 2).

$$\min_{\mathbf{W}, \mathbf{U}} \|\mathbf{T}^\top \mathbf{W} + \mathbf{O}^\top \mathbf{U} - \mathbf{Y}\|_F^2 + \lambda_1 \|\mathbf{W}\|_S \quad (3)$$

Similarly, we also introduce a new *attribute norm* to learn the importance of various attribute modalities of the objects. Attribute modalities can describe the color histograms (e.g., red, green, and blue), shape (e.g., gradient features), texture, or relationships of the object to the human (e.g., distances). We define the attribute norm over the weight matrix  $\mathbf{U}$  as:

$$\|\mathbf{U}\|_A = \sum_{c=1}^C \sum_{o=1}^O \sum_{m=1}^M \|\mathbf{u}_c^{o_m}\|_2 \quad (4)$$

The  $\ell_2$ -norm is employed to enforce similar weights within an attribute modality, and the  $\ell_1$ -norm is used between these

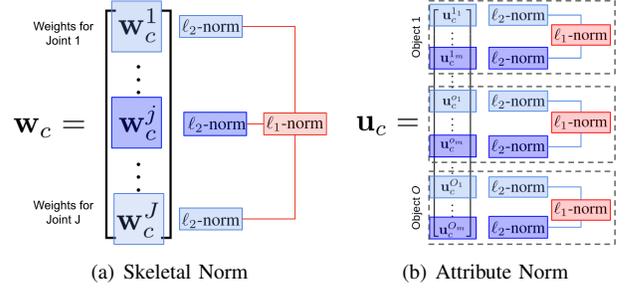


Fig. 2. Illustrations of the proposed regularization norms. Figure 2(a) shows the *skeletal norm*  $\|\mathbf{W}\|_S$ . The  $\ell_2$ -norm is used within joints and the  $\ell_1$ -norm is used between joints to enforce sparsity and the identification of discriminative joints. Figure 2(b) shows the *attribute norm*  $\|\mathbf{U}\|_A$ . The  $\ell_2$ -norm is used within attribute modalities and the  $\ell_1$ -norm is applied between modalities to enforce sparsity and identify discriminative attributes.

attribute modalities to enforce sparsity in order to identify discriminative object attributes (Figure 2).

Then, our final formulation formulates activity recognition as a regularized optimization problem:

$$\min_{\mathbf{W}, \mathbf{U}} \|\mathbf{T}^\top \mathbf{W} + \mathbf{O}^\top \mathbf{U} - \mathbf{Y}\|_F^2 + \lambda_1 \|\mathbf{W}\|_S + \lambda_2 \|\mathbf{U}\|_A \quad (5)$$

where  $\lambda_1$  and  $\lambda_2$  denote the hyperparameters to balance the importance of the loss function and regularization norms.

### C. Recognizing Human Activities

After we solve the regularized optimization problem in Eq. (5) using Algorithm 1, we obtain the optimal weight matrices  $\mathbf{W}^* = [\mathbf{w}_1^*, \dots, \mathbf{w}_C^*] \in \mathbb{R}^{d_T \times C}$  and  $\mathbf{U}^* = [\mathbf{u}_1^*, \dots, \mathbf{u}_C^*] \in \mathbb{R}^{d_O \times C}$ . Each column  $\mathbf{w}_c^*$  and  $\mathbf{u}_c^*$  denotes the importance of, respectively, an observation of a human  $\mathbf{t}$  and objects  $\mathbf{o}$  to recognize the  $c$ -th activity category. Given a new observation  $\mathbf{x} = \{\mathbf{t}, \mathbf{o}\}$ , the activity category  $y(\mathbf{t}, \mathbf{o})$  is classified by:

$$y(\mathbf{t}, \mathbf{o}) = \max_c \mathbf{t}^\top \mathbf{w}_c^* + \mathbf{o}^\top \mathbf{u}_c^* \quad (6)$$

Since the objective function in our formulation can be used to perform classification, no separate classifiers is needed.

By learning the weight matrices for both human and object observations, our approach explicitly identifies discriminative human joints and object attributes. For example, consider the learned human observation weight matrix  $\mathbf{W}^*$ :

$$\mathbf{W}^* = \begin{bmatrix} \mathbf{w}_1^1 & \dots & \mathbf{w}_c^1 & \dots & \mathbf{w}_C^1 \\ \vdots & \ddots & \vdots & & \vdots \\ \mathbf{w}_1^j & \dots & \mathbf{w}_c^j & \dots & \mathbf{w}_C^j \\ \vdots & & \vdots & \ddots & \vdots \\ \mathbf{w}_1^J & \dots & \mathbf{w}_c^J & \dots & \mathbf{w}_C^J \end{bmatrix} \quad (7)$$

where  $\mathbf{w}_c^j$  represents the importance of the  $j$ -th joint to the  $c$ -th activity category. The sum-value of all elements within the sub-matrix  $\mathbf{w}_c^j$  indicates the relative importance of the  $j$ -th human joint when recognizing the  $c$ -th activity category. The sum-value of all elements in the row vector  $\mathbf{w}_c^j$  indicates the importance of the  $j$ -th human body joint to recognize all activity categories. Similarly,  $\mathbf{U}$  can also be used to analyze and identify which attributes of which objects are important when recognizing human activities.

**Algorithm 1:** An iterative algorithm to solve the formulated optimization problem in Eq. (5).

**Input :**  $\mathbf{T} = [\mathbf{t}_1, \dots, \mathbf{t}_N] \in \mathbb{R}^{d_T \times N}$ ,  
 $\mathbf{O} = [\mathbf{o}_1, \dots, \mathbf{o}_N] \in \mathbb{R}^{d_O \times N}$  and  
 $\mathbf{Y} = [\mathbf{y}^1; \dots; \mathbf{y}^N] \in \mathbb{R}^{N \times C}$ .  
**Output :**  $\mathbf{W}^* = \mathbf{W}(i) \in \mathbb{R}^{d_T \times C}$  and  
 $\mathbf{U}^* = \mathbf{U}(i) \in \mathbb{R}^{d_O \times C}$ .

- 1: Let  $i = 1$ . Initialize  $\mathbf{W}$  and  $\mathbf{U}$  randomly.
- 2: **repeat**
- 3:     Calculate  $\mathbf{D}_S^c(i+1)$  for  $c \in 1, \dots, C$ .
- 4:     Calculate  $\mathbf{D}_A^c(i+1)$  for  $c \in 1, \dots, C$ .
- 5:     Calculate  $\mathbf{w}_c(i+1)$  via Eq. (9) for each  $c \in 1, \dots, C$ .
- 6:     Calculate  $\mathbf{u}_c(i+1)$  via Eq. (11) for each  $c \in 1, \dots, C$ .
- 7:      $i = i + 1$ .
- 8: **until convergence;**
- 9: **return**  $\mathbf{W}^*$  and  $\mathbf{U}^*$

#### D. Optimization Algorithm

The formulated optimization problem in Eq. (5) is difficult to solve because the regularization norms  $\|\mathbf{W}\|_S$  and  $\|\mathbf{U}\|_A$  are not smooth and because we need to simultaneously find the optimal solutions for  $\mathbf{W}$  and  $\mathbf{U}$ , both of which the final solution depends on. To solve this, we propose a new iterative optimization solver as presented in Algorithm 1.

We calculate the derivative of Eq. (5) with respect to  $\mathbf{w}_c$  in order to solve  $\mathbf{W}$ :

$$\mathbf{T}\mathbf{T}^\top \mathbf{w}_c + \mathbf{T}\mathbf{O}^\top \mathbf{u}_c - \mathbf{T}\mathbf{y}_c + \lambda_1 \mathbf{D}_S^c \mathbf{w}_c = \mathbf{0} \quad (8)$$

$$\mathbf{w}_c = (\mathbf{T}\mathbf{T}^\top + \lambda_1 \mathbf{D}_S^c)^{-1} \mathbf{T} (\mathbf{y}_c - \mathbf{O}^\top \mathbf{u}_c) \quad (9)$$

where  $\mathbf{D}_S^c$  is a block diagonal matrix with  $\frac{1}{2\|\mathbf{w}_c^j\|_2} \mathbf{I}_{d_T^j}$  as the  $j$ -th block.

Similarly, we compute the derivative of Eq. (5) with respect to  $\mathbf{u}_c$  in order to solve  $\mathbf{U}$ :

$$\mathbf{O}\mathbf{O}^\top \mathbf{u}_c + \mathbf{O}\mathbf{T}^\top \mathbf{w}_c - \mathbf{O}\mathbf{y}_c + \lambda_2 \mathbf{D}_A^c \mathbf{u}_c = \mathbf{0} \quad (10)$$

$$\mathbf{u}_c = (\mathbf{O}\mathbf{O}^\top + \lambda_2 \mathbf{D}_A^c)^{-1} \mathbf{O} (\mathbf{y}_c - \mathbf{T}^\top \mathbf{w}_c) \quad (11)$$

where  $\mathbf{D}_A^c$  is a block diagonal matrix having  $O$  blocks. Each of these diagonal blocks is composed of  $M$  diagonal blocks, where the  $m$ -th diagonal block is  $\frac{1}{2\|\mathbf{u}_c^m\|_2} \mathbf{I}_{d_O^m}$ .

Because the solution to  $\mathbf{w}_c$  depends on both  $\mathbf{u}_c$  and  $\mathbf{D}_S^c$  (which is dependent on  $\mathbf{w}_c$ ), and the solution to  $\mathbf{u}_c$  depends on both  $\mathbf{w}_c$  and  $\mathbf{D}_A^c$  (which is dependent on  $\mathbf{u}_c$ ), an iterative optimization algorithm is necessary to address this problem. The proposed optimization solver is detailed in Algorithm 1, which alternately solving  $\mathbf{w}_c$  and  $\mathbf{u}_c$  until convergence. This proposed algorithm holds a theoretical guarantee to converge to the optimal solution:

**Theorem 1:** Algorithm 1 is guaranteed to converge to the optimal solution to the formulated regularized optimization problem in Eq. (5).

*Proof:* See supplementary materials<sup>1</sup>. ■

The time complexity of Algorithm 1 is dominated by Steps (5) and (6), because Steps (3) and (4) are trivial, executing

in linear time of  $\mathcal{O}(Cd_T)$  and  $\mathcal{O}(Cd_O)$ , respectively. Steps (5) and (6) can be solved as a system of linear equations instead of performing the matrix inverse, and are respectively  $\mathcal{O}(d_T^2)$  and  $\mathcal{O}(d_O^2)$ .

## IV. EXPERIMENTAL RESULTS

We assess our approach's performance on two benchmark activity recognition datasets and using a physical robot as a case study. In the experiments, we used one type of skeletal features and three different object attribute modalities. The skeletal feature used is the displacement of each body joint from the central torso joint. The object attributes used are color, shape, and object-joint distance. The color attribute is implemented using red, green, and blue histograms. The shape attribute is implemented using the Histogram of Oriented Gradients (HOG) [35] features. Finally, the object-joint distance attribute is implemented through calculating the distance in 3D space from the object to each skeletal joint to model the interaction between objects and the human. We also evaluate our case study using a multinomial probability distribution as the object attributes, showing our method's ability to identify the most important objects with respect to specific human activities.



Fig. 3. Example images from the CAD-60 dataset.

#### A. Results on Cornell Activity Dataset

The Cornell Activity Dataset (CAD-60) has been widely used as a standard benchmark dataset for activity recognition [36] in robotics. The dataset consists of activities performed by 4 humans, with each activity execution recorded as color images, depth images, and annotated skeleton joint positions for 15 joints. Our experiments used six activities that involve the use of objects, including *writing on whiteboard*, *cooking (stirring)*, *cooking (chopping)*, *working on computer*, *rinsing mouth with water*, and *talking on the phone*, as illustrated in Figure 3. All objects existing in a scene were utilized, even if they did not relate to the human activity. For example, the whiteboard was still included in our experiments when recognizing the cooking (stirring) human activity scenes.

The quantitative experimental results are listed in Table I. It is observed that the proposed approach is able to identify 98.11% of these activity executions correctly. This table also compares our results with other state-of-the-art methods for human activity recognition, which indicates that our simultaneous learning from both human observations and object observations provides superior performance. This table also shows that limiting our approach to only one of our proposed sparsity inducing norms degrades our performance. This drop off is significant when we only deploy the *skeletal norm* in the formulation, showing that object attributes provide useful discriminative information.

<sup>1</sup>hcr.mines.edu/publication/HAR\_Supp.pdf

TABLE I

ACCURACY OBTAINED BY OUR APPROACH ON THE CAD-60 DATASET AND COMPARISONS TO PREVIOUS APPROACHES.

| Approach                                   | Accuracy      |
|--|---------------|
| Feature and Body Part Learning [14]        | 83.93%        |
| Joint Heterogeneous Features Learning [33] | 84.10%        |
| Spatiotemporal Interest Point [37]         | 87.50%        |
| Feature-Level Fusion [38]                  | 87.50%        |
| Pose Kinetic Energy [39]                   | 91.90%        |
| Sparse Coding Dictionary Learning [40]     | 94.12%        |
| Kinect + Pose machine [41]                 | 95.58%        |
| Our Approach (only <i>skeletal norm</i> )  | 86.86%        |
| Our Approach (only <i>attribute norm</i> ) | 96.18%        |
| <b>Our Approach</b>                        | <b>98.11%</b> |

### B. Results on MSR Activity 3D Dataset

We further evaluate our approach based on the MSR Daily Activity 3D Dataset [42], a commonly used public dataset for benchmarking human activity recognition approaches. This dataset consists of daily human activities performed by ten human subjects, which is recorded with color images, depth images, and annotated joint positions for 20 skeletal joints. Examples of color and depth images from activity categories involving the use of objects are illustrated in Figure 4.



Fig. 4. Example images from the MSR Daily Activity 3D dataset.

The quantitative experimental results are listed in Table II. It is observed that our approach achieves an activity recognition accuracy of 97.71%. Table II also compares our results with existing state-of-the-art approaches that have also been tested on this dataset, showing again that our simultaneous learning from both human body joints and object attributes is effective to recognize human activities. Moreover, the results indicate the importance of the proposed *attribute norm*, as omitting this norm causes our approach to lose the ability of identifying most important object attributes, thus decreasing the recognition accuracy.

### C. Results in Home Assistance Case Studies

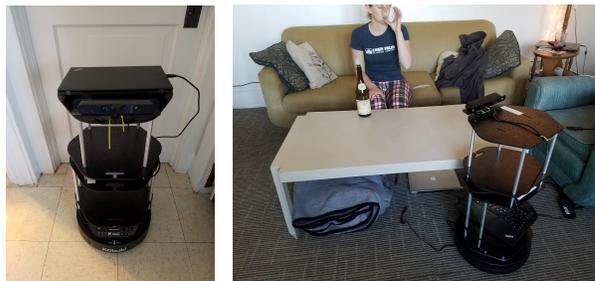
In addition to evaluating our method on the publicly available benchmark datasets, we also implemented our approach on a physical robot in order to demonstrate its performance in case study. We deployed our approach on a Turtlebot robot participating in a simulated home assistance scenario (Figure 5). The robot has a color-depth sensor onboard to extract 3D skeleton data and a lightweight netbook for processing.

In this scenario, five activities were defined, with Figure 6 showing example color and depth images for each activity. These activities are *drinking wine*, *storing food*, *storing dishes*, *pouring wine*, and *eating*. Each activity was performed 20 times. In order to test our approach in learning

TABLE II

ACCURACY OBTAINED BY OUR APPROACH ON THE MSR DAILY ACTIVITY 3D DATASET AND COMPARISONS TO PREVIOUS APPROACHES.

| Approach                                   | Accuracy      |
|--|---------------|
| Sparse Coding Dictionary Learning [40]     | 68.75%        |
| BIPOD [10]                                 | 79.70%        |
| Spatiotemporal Interest Point [37]         | 80.00%        |
| Key-Pose-Motifs [43]                       | 83.47%        |
| Kinect + Pose machine [41]                 | 84.37%        |
| Feature-Level Fusion [38]                  | 88.80%        |
| 3D joint+CS-Mltp (concatenate) [31]        | 92.50%        |
| DL-GSGC [44]                               | 95.00%        |
| Joint Heterogeneous Features Learning [33] | 95.00%        |
| $\tau$ -test [45]                          | 95.63%        |
| Our Approach (only <i>skeletal norm</i> )  | 82.00%        |
| Our Approach (only <i>attribute norm</i> ) | 95.71%        |
| <b>Our Approach</b>                        | <b>97.71%</b> |



(a) TurtleBot Platform

(b) Scenario Setup

Fig. 5. The experiment setup used in the case studies. Figure 5(a) shows the Turtlebot platform equipped with an ASUS Xtion Pro camera and laptop for computation. Figure 5(b) shows the robot observing human activities.

simultaneously from observations of the human and the objects, these activities were defined to involve similar objects and human poses. For example, both drinking wine and pouring wine involve a glass and a bottle; however, drinking wine is performed while sitting down and pouring wine is performed while standing up. Similarly, both eating and drinking wine are activities performed by a sitting human, but they involve different objects (respectively, a bowl and a spoon versus a wine glass and a bottle).

The quantitative experimental results are presented in Table III. We can observe that the proposed approach achieves an overall accuracy of 98.33% in the case studies. Comparison with baseline real-time approaches is also listed in Table III, which shows that our approach is superior to two standard real-time machine learning methods as baselines. With only the *skeletal norm* or the *attribute norm* as the regularization, our approach achieves good accuracy but less than with the complete formulation that uses both norms.



Fig. 6. Color and depth images of the activities included in our case studies in a simulated home assistance scenario. From left to right, *drinking wine*, *storing food*, *storing dishes*, *pouring wine*, and *eating*.

TABLE III

ACCURACY OBTAINED BY OUR APPROACH IN THE CASE STUDIES AND COMPARISON TO BASELINE REAL-TIME APPROACHES.

| Approach                                   | Accuracy      |
|--|---------------|
| Support Vector Machine                     | 51.67%        |
| Decision Forest                            | 91.67%        |
| Our Approach (only <i>skeletal norm</i> )  | 95.00%        |
| Our Approach (only <i>attribute norm</i> ) | 96.67%        |
| Our Approach                               | <b>98.33%</b> |

We also tested our method with a different set of attributes in order to assess its ability to identify discriminative objects. In this setup, we defined five attribute modalities, where each modality is the probability that an object category appeared in the view of the robot. The 5 object categories used are the *wine bottle*, *glass*, *fridge*, *bowl*, and *spoon*. The probabilities that an object appeared in a scene were obtained from the YOLO object detection system [46], which uses a pre-trained neural network to identify common household objects. For example, for the activity of drinking wine, the probability of a bottle or glass appearing would be close to 1, and close to 0 for the remaining objects.

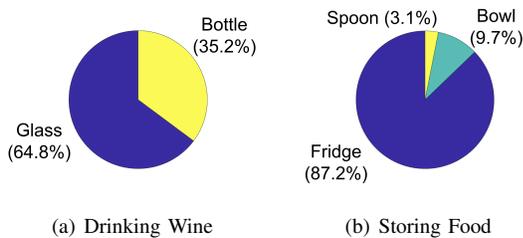


Fig. 7. Illustration of the distribution of weights in two columns of the weight matrix  $\mathbf{U}$ . Figure 7(a) demonstrates the weights for the drinking wine activity, where the glass and bottle are important. Figure 7(b) illustrates the weights for the human activity of storing food, where the fridge is the most relevant object.

Using this setup, our approach is able to recognize 96.67% of home activities correctly. Additionally, this setup allowed our approach to identify discriminative objects, as each column of the  $\mathbf{U}$  matrix contained only five values, each relating one object to that column’s associated human activity. Figure 7 displays two columns from the  $\mathbf{U}$  weight matrix. In Figure 7(a), we observe that the bottle and glass are the only objects receiving weights, as these are very indicative of the drinking wine activity. Similarly in Figure 7(b), we see that the fridge receives nearly 90% of the total column weight, identifying it as being very indicative to recognize the storing food activity. Similarly, Figure 8 displays this relationship for two of the rows of  $\mathbf{U}$ . Figure 8(a) demonstrates that for the bowl object, the most related activity is eating. Figure 8(b) shows this for the fridge object, which shows that the most relevant activity is storing food, the only activity category in which this object appears.

#### D. Discussion

1) *Real-Time Performance*: One of the major advantages of our approach is its ability to run at real-time speeds. For each dataset we evaluated, we analyzed the runtime of our activity recognition approach, summarized in Table IV. Due

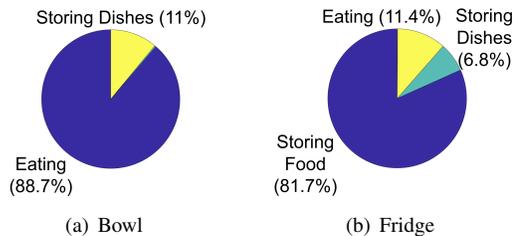


Fig. 8. Illustration of the distribution of weights in two rows of the weight matrix  $\mathbf{U}$ . Figure 8(a) illustrates the weights for the bowl object, where the eating is the most relevant activity. Figure 8(b) depicts the weights for the fridge object, where storing food is the most relevant activity.

to the efficiency of our proposed convex problem formulation that integrates classification within the loss function, our approach obtains recognition processing speeds in excess of  $2 \times 10^4$  Hz, when executing on an Intel i5 processor with 4Gb memory. Our approach provides a suitable solution for accurate and real-time recognition of human activities.

TABLE IV

EXPERIMENTAL RESULTS ON REAL-TIME PERFORMANCE.

| Metric                | CAD-60               | MSR                  | Home Asst.           |
|-----------------------|----------------------|----------------------|----------------------|
| Processing Speed (Hz) | $6.1 \times 10^4$    | $2.5 \times 10^4$    | $2.2 \times 10^4$    |
| Time Per Frame (sec)  | $1.6 \times 10^{-5}$ | $3.9 \times 10^{-5}$ | $4.5 \times 10^{-5}$ |

2) *Hyperparameter Selection*: In our problem formulation in Eq. (5), the hyperparameters  $\lambda_1$  and  $\lambda_2$  control the importance of the *skeletal norm* and *attribute norm*, respectively, and balance these two norms with the loss function. As our presented results have demonstrated, the method’s accuracy decreases with either of these hyperparameters assigned to 0. As each norm captures different information (i.e., weights of skeletal features or object attributes), both are necessary for the proposed approach to achieve its state-of-the-art accuracy. Also, we observe that as the values of these hyperparameters become too large, performance decreases as the loss function relating observations to activity labels becomes less important. Our presented results use the hyperparameter values of  $\lambda_1 = 0.1$  and  $\lambda_2 = 0.1$ .

#### V. CONCLUSION

In this paper, we introduce a principled method for activity recognition through simultaneous learning from observations of human poses and object attributes. The proposed approach is capable of identifying discriminative joints and object attributes, providing an interpretable understanding of their importance to human activities and the relationships between joints and objects. We formulate activity recognition as an optimization problem that uses a regression-like loss function to integrate teammate and object cues to perform activity recognition, and utilizes sparsity-inducing norms to estimate feature importance. We introduce an iterative algorithm guaranteed to find the optimal solution. We assess our proposed approach on two benchmark activity recognition datasets, and on an actual robot to show a case study. Experimental results have shown that our approach achieves state-of-the-art recognition accuracy and obtains real-time performance.

## REFERENCES

- [1] G.-J. M. Kruijff, M. Janiček, S. Keshavdas, B. Larochelle, H. Zender, N. J. Smets, T. Mioch, M. A. Neerinx, J. V. Diggelen, and F. Colas, "Experience in system design for human-robot teaming in urban search and rescue," in *Field and Service Robotics*, pp. 111–125, Springer, 2014.
- [2] T. Fong, "Human-robot teaming: Communication, coordination, and collaboration," tech. rep., NASA, 2017.
- [3] R. Schulz, P. Kratzer, and M. Toussaint, "Preferred interaction styles for human-robot collaboration vary over tasks with different action types," *Frontiers in Neurorobotics*, vol. 12, p. 36, 2018.
- [4] H. Zhang, C. Reardon, C. Zhang, and L. E. Parker, "Adaptive human-centered representation for activity recognition of multiple individuals from 3d point cloud sequences," in *IEEE International Conference on Robotics and Automation*, 2015.
- [5] F. Han, B. Reily, W. Hoff, and H. Zhang, "Space-time representation of people based on 3d skeletal data: A review," *Computer Vision and Image Understanding*, vol. 158, pp. 85–105, 2017.
- [6] M. Vrigkas, C. Nikou, and I. A. Kakadiaris, "A review of human activity recognition methods," *Frontiers in Robotics and AI*, vol. 2, p. 28, 2015.
- [7] J. K. Aggarwal and M. S. Ryoo, "Human activity analysis: A review," *ACM Computing Surveys*, vol. 43, no. 3, p. 16, 2011.
- [8] J. K. Aggarwal and L. Xia, "Human activity recognition from 3d data: A review," *Pattern Recognition Letters*, vol. 48, pp. 70–80, 2014.
- [9] J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, and A. Blake, "Real-time human pose recognition in parts from single depth images," in *Conference on Computer Vision and Pattern Recognition*, 2011.
- [10] B. Reily, F. Han, L. E. Parker, and H. Zhang, "Skeleton-based bio-inspired human activity prediction for real-time human-robot interaction," *Autonomous Robots*, vol. 42, no. 6, pp. 1281–1298, 2018.
- [11] S. Z. Masood, C. Ellis, A. Nagaraja, M. F. Tappen, J. J. LaViola, and R. Sukthankar, "Measuring and reducing observational latency when recognizing actions," in *International Conference on Computer Vision Workshops*, 2011.
- [12] J. Wang, Z. Liu, Y. Wu, and J. Yuan, "Learning actionlet ensemble for 3d human action recognition," *Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 5, pp. 914–927, 2014.
- [13] G. Yu, Z. Liu, and J. Yuan, "Discriminative orderlet mining for real-time recognition of human-object interaction," in *Asian Conference on Computer Vision*, pp. 50–65, Springer, 2014.
- [14] F. Han, X. Yang, C. Reardon, Y. Zhang, and H. Zhang, "Simultaneous feature and body-part learning for real-time robot awareness of human behaviors," in *International Conference on Robotics and Automation*, 2017.
- [15] S. R. Fanello, I. Gori, G. Metta, and F. Odone, "Keep it simple and sparse: Real-time action recognition," *The Journal of Machine Learning Research*, vol. 14, no. 1, pp. 2617–2640, 2013.
- [16] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. Van Gool, "Temporal segment networks for action recognition in videos," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.
- [17] 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.
- [18] P. Wei, Y. Zhao, N. Zheng, and S.-C. Zhu, "Modeling 4d human-object interactions for event and object recognition," in *International Conference on Computer Vision*, 2013.
- [19] H. S. Koppula, R. Gupta, and A. Saxena, "Learning human activities and object affordances from rgb-d videos," *International Journal of Robotics Research*, vol. 32, no. 8, pp. 951–970, 2013.
- [20] S. Rossi, F. Ferland, and A. Tapus, "User profiling and behavioral adaptation for hri: A survey," *Pattern Recognition Letters*, vol. 99, pp. 3–12, 2017.
- [21] P. Stone, G. A. Kaminka, S. Kraus, and J. S. Rosenschein, "Ad hoc autonomous agent teams: Collaboration without pre-coordination.," in *AAAI Conference on Artificial Intelligence*, 2010.
- [22] S. V. Albrecht and P. Stone, "Autonomous agents modelling other agents: A comprehensive survey and open problems," *Artificial Intelligence*, vol. 258, pp. 66–95, 2018.
- [23] J. Heard, R. Heald, C. E. Harriott, and J. A. Adams, "A diagnostic human workload assessment algorithm for human-robot teams," in *International Conference on Human-Robot Interaction*, 2018.
- [24] S. al Mahi, M. Atkins, and C. Crick, "Learning to assess the cognitive capacity of human partners," in *International Conference on Human-Robot Interaction*, 2017.
- [25] H. Zhang and L. E. Parker, "Bio-inspired predictive orientation decomposition of skeleton trajectories for real-time human activity prediction," in *IEEE International Conference on Robotics and Automation*, 2015.
- [26] X. Chen and M. Koskela, "Online rgb-d gesture recognition with extreme learning machines," in *International conference on multimodal interaction*, 2013.
- [27] H. Rahmani, A. Mahmood, D. Q. Huynh, and A. Mian, "Real time action recognition using histograms of depth gradients and random decision forests," in *Winter Conference on applications of Computer Vision*, 2014.
- [28] S. Boubou and E. Suzuki, "Classifying actions based on histogram of oriented velocity vectors," *Journal of Intelligent Information Systems*, vol. 44, no. 1, pp. 49–65, 2015.
- [29] X. Yang and Y. Tian, "Effective 3d action recognition using eigen-joints," *Journal of Visual Communication and Image Representation*, vol. 25, no. 1, pp. 2–11, 2014.
- [30] W. Zhu, C. Lan, J. Xing, W. Zeng, Y. Li, L. Shen, and X. Xie, "Co-occurrence feature learning for skeleton based action recognition using regularized deep lstm networks," in *AAAI Conference on Artificial Intelligence*, 2016.
- [31] J. Luo, W. Wang, and H. Qi, "Spatio-temporal feature extraction and representation for rgb-d human action recognition," *Pattern Recognition Letters*, vol. 50, pp. 139–148, 2014.
- [32] H. S. Koppula and A. Saxena, "Anticipating human activities using object affordances for reactive robotic response," *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 1, pp. 14–29, 2016.
- [33] J.-F. Hu, W.-S. Zheng, J. Lai, and J. Zhang, "Jointly learning heterogeneous features for rgb-d activity recognition," in *Conference on Computer Vision and Pattern Recognition*, 2015.
- [34] J.-F. Hu, W.-S. Zheng, J. Lai, S. Gong, and T. Xiang, "Exemplar-based recognition of human-object interactions," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 26, no. 4, pp. 647–660, 2016.
- [35] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *Conference on Computer Vision and Pattern Recognition*, 2005.
- [36] J. Sung, C. Ponce, B. Selman, and A. Saxena, "Unstructured human activity detection from rgb-d images," in *International Conference on Robotics and Automation*, 2012.
- [37] Y. Zhu, W. Chen, and G. Guo, "Evaluating spatiotemporal interest point features for depth-based action recognition," *Image and Vision Computing*, vol. 32, no. 8, pp. 453–464, 2014.
- [38] Y. Zhu, W. Chen, and G. Guo, "Fusing multiple features for depth-based action recognition," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 6, no. 2, p. 18, 2015.
- [39] J. Shan and S. Akella, "3d human action segmentation and recognition using pose kinetic energy," in *International Workshop on Advanced Robotics and its Social Impacts*, 2014.
- [40] J. Qi, Z. Wang, X. Lin, and C. Li, "Learning complex spatio-temporal configurations of body joints for online activity recognition," *IEEE Transactions on Human-Machine Systems*, vol. 48, no. 6, pp. 637–647, 2018.
- [41] S. Das, M. Koperski, F. Bremond, and G. Francesca, "Action recognition based on a mixture of rgb and depth based skeleton," in *International Conference on Advanced Video and Signal Based Surveillance*, 2017.
- [42] W. Li, Z. Zhang, and Z. Liu, "Action recognition based on a bag of 3d points," in *Computer Vision and Pattern Recognition Workshops*, 2010.
- [43] C. Wang, Y. Wang, and A. L. Yuille, "Mining 3d key-pose-motifs for action recognition," in *Conference on Computer Vision and Pattern Recognition*, 2016.
- [44] J. Luo, W. Wang, and H. Qi, "Group sparsity and geometry constrained dictionary learning for action recognition from depth maps," in *International Conference on Computer Vision*, 2013.
- [45] C. Lu, J. Jia, and C.-K. Tang, "Range-sample depth feature for action recognition," in *Conference on Computer Vision and Pattern Recognition*, 2014.
- [46] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," *CoRR*, vol. abs/1804.02767, 2018.