

This is the peer reviewed version of the following article:

A General Pipeline for Online Gesture Recognition in Human-Robot Interaction / Villani, Valeria; Secchi, Cristian; Lippi, Marco; Sabattini, Lorenzo. - In: IEEE TRANSACTIONS ON HUMAN-MACHINE SYSTEMS. - ISSN 2168-2291. - 53:2(2023), pp. 315-324. [10.1109/THMS.2022.3227309]

*Terms of use:*

The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

25/04/2024 18:16

(Article begins on next page)

# A General Pipeline for Online Gesture Recognition in Human-Robot Interaction

Valeria Villani, Cristian Secchi, Marco Lippi and Lorenzo Sabattini

**Abstract**—Recent advances in robotics have allowed the introduction of robots assisting and working together with human subjects. To promote their use and diffusion, intuitive and user-friendly interaction means should be adopted. In particular, gestures have become an established way to interact with robots since they allow to command them in an intuitive manner. In this paper, we focus on the problem of gesture recognition in human-robot interaction (HRI). While this problem has been largely studied in the literature, it poses specific constraints when applied to HRI. We propose a framework consisting in a pipeline devised to take into account these specific constraints. We implement the proposed pipeline considering, as an example, an evaluation use case. To this end, we consider standard machine learning algorithms for the classification stage and evaluate their performance considering different performance metrics for a thorough assessment.

## I. INTRODUCTION

Over the last years, progress in the design and development of robotic systems has led to advanced solutions that are entering our daily lives in several application fields, such as social assistance, surveillance, tour guidance, rehabilitation, and search and rescue. From the point of view of human-robot interaction (HRI), in addition to the requirements for safe interaction, this trend has further implications on the modalities used to communicate with robots. First, it is important that communication from and to the robot is quick, smooth, co-located and requires low attentional demand. In other words, for efficient HRI, communicating to the robot and understanding its messages should not be a bottleneck and should not increase the overall complexity of teaming up with a robotic agent, as it is usually in human-human communication. Moreover, given the fact that most robotic assistants are intended to be used by people without expertise in robotics, it is also important that intuitive communication is enabled, so that users do not have to learn commands that are specific for the robot. In this regard, gestures represent a valid candidate for intuitive communication with robots.

Gesture-based interaction with machines or robots has been proposed by a large body of the literature, in both everyday and industrial applications [1]–[10]. The problem of gesture-based interaction can be split in two specific problems: recognizing gestures and mapping gestures into commands to the interaction system. The focus of our paper is on the first specific problem applied to the domain of robotics, that is gesture recognition in the context of HRI. While gesture recognition

has long been an active area of research in computer vision and machine learning [11], the use of gestures for interaction with robots poses specific constraints. These constraints define the prerequisites for successful gesture-based HRI. To improve the use of gestures for HRI it is, then, beneficial that approaches to gesture recognition are compliant by design with such constraints. The second specific problem refers to designing how gestures can be translated into commands to the robot, possibly in an intuitive manner [2]. While these two specific problems are independent, they jointly address the general problem of gesture-based interaction with robots.

Building upon these lines, in this paper we propose a pipeline for online gesture recognition in HRI. The approach has been designed to be general and can be applied to any set of gestures for interacting with any robot. To the best of our knowledge, this is the first attempt in this direction, since existing approaches have been designed to address specific use case scenarios. Conversely, we design the proposed pipeline starting from the understanding of how gesture-based interaction, in general, can be applied to HRI; no additional constraints or requirements derived by specific use cases are included to design the pipeline. As a result, the proposed pipeline is application agnostic.

To achieve this, in Section III-A we first discuss the specific constraints of gesture recognition for HRI, which have to be taken into account to use gesture-based interaction in daily life. Then, having these constraints in mind, in Section III-B we propose an algorithmic pipeline to implement gesture recognition in HRI. To show and discuss the application of the proposed pipeline, we consider the experimental scenario presented in [2], [12], [13]. In particular, we assess the capability of the pipeline to generalize across multiple subjects having different level of acquaintance with the use of gestures. Indeed, prospectively, the use of gestures for interacting with robots should be designed in such an intuitive way that it is available to any user, and does not require specific prior experience with gestures.

The rest of the paper is organized as follows. Section II reports the state-of-the-art on gesture recognition in HRI. In Section III we discuss how HRI constrains the use of gesture-based interaction and, building upon such discussion, we present the proposed pipeline. Then, Section IV focuses on the classification stage, which is the core of the proposed approach. In Section V the experimental setting considered to validate the proposed architecture is presented. Then, in Section VI we report the results of implementing the proposed pipeline to the considered use case. Finally, Section VII follows with some concluding remarks.

Authors are with Department of Sciences and Methods for Engineering, University of Modena and Reggio Emilia, Reggio Emilia, Italy. Email: {vvillani, mlippi, csecchi, lsabattini}@unimore.it

## II. GESTURE RECOGNITION FOR HRI

A large part of existing approaches to gesture recognition in the domain of HRI rely on the use of vision systems [14], [15]. In [5] a stereo camera is mounted on the head of a personal service robot for elderly assistance. Whole-body gestures, such as walking or raising a hand, are then recognized from estimated 3-D human body components. Burke and Lasenby in [6] have proposed the use of a Kinect sensor to detect pantomimic hand gestures that control an unmanned aerial vehicle. Xu *et al.* in [8] have focused on the problem of background subtraction when using RGB-D cameras for hand gesture detection in home-like dynamic environment. In [7] vision-based gesture recognition is applied by Kim *et al.* to address the problem of HRI at a long distance, approximately five meters from the camera. The authors in [9] have considered to classify the arm trajectories, seen as sequences of motor primitives. To this end, RGB video sequences are used, with the subject standing in front of the robot in a constrained setup. The same requirement, that is the subject standing in front of the robot, is set by Ciciirelli *et al.* in [10], despite the use of multiple Kinect cameras to monitor the surrounding environment. Chandarana *et al.* in [16] used the infrared-based Leap Motion controller to detect hand gestures for teleoperating unmanned aerial vehicles (UAVs). UAV trajectories are built combining gestures that define portions of flight paths. Generally speaking, vision-based approaches require proper lighting conditions and camera angles, and the user has to be in the field of view of the camera. Hence, these characteristics pose a limitation to the use of such approaches in real world HRI applications.

To overcome the limitations linked to the use of vision systems and given the recent advances in pervasive computing, wearable devices have been used to detect and recognize gestures [17]. Indeed, unobtrusive, wireless and inexpensive body worn sensors, such as accelerometers and gyroscopes, possibly integrated in inertial measurement units (IMUs), are available on everyday mobile and ubiquitous systems, such as smartwatches, wrist bands and mobile phones. They provide information about body movements and, hence, can be used to track user activity [18], [19]. Additionally, they have been used for gesture recognition for human-machine interaction in several works [3], [20], [21]. For what concerns gesture-based interaction with robots, the authors in [4] have proposed the use of five IMUs and an ultra wideband positioning system to capture the human upper body shape and the relative position between the human and the robot. Villani *et al.* in [2], [12] have used inertial data recorded with a smartwatch to control both wheeled and aerial robots. Gestures are used to provide high-level commands, such as take off, land or stop, whereas robot velocity is determined by mapping user's wrist movements. A similar setting was considered by Carfi *et al.* in [22]. However, their framework was not designed for HRI. Indeed, gestures were executed in steady conditions, which means that the user was in a fixed predetermined pose between consecutive gestures. This condition limits the possibility to use this interaction means while the user is performing everyday activities. Additionally, inter-subject robustness of the classification approach was assessed to a limited extent since

TABLE I  
SPECIFIC CONSTRAINTS OF GESTURE RECOGNITION IN HRI AND  
PROPOSED APPROACHES TO COMPLY WITH THEM.

Specific constraints in HRI	Proposed solution
C1 – Wearable and lean infrastructure	Smartwatch or other wrist-worn device embedding an IMU
C2 – Online gesture recognition	Choice of efficient classification stage
C3 – Online application in dynamic conditions (e.g., moving user)	Sliding windows and specific decision rules
C4 – Robustness to inter-subject variability	Different user's confidence with the gestures under analysis
C5 – Minimization of false positives	Model selection with specific performance metrics

it was tested with gestures performed by the same subjects involved in the training phase.

Furthermore, wearable sensors based on surface electromyography (sEMG) are being used for gesture recognition [23]–[25]. Quite often, IMU and EMG data are combined together to improve gesture recognition. Indeed, while inertial data provide information about hand position, EMG sensors allow to fully understand complex finger or hand gestures. This is the case, for example, of the work by Jiang *et al.* in [26], where sEMG and IMU sensing fusion allows to recognize several air and surface gestures with two distinct force levels. Georgi *et al.* in [27] have proposed the simultaneous usage of IMU and EMG sensors for gesture-based interfaces and Hidden Markov Models (HMMs) are used as classifiers to discriminate between the defined gesture classes.

From an algorithmic point of view, most approaches to gesture recognition resort to machine learning techniques to deal with high-dimensional, multimodal streams of data that are characterized by a large variability. Different machine learning algorithms applied to gesture recognition have been compared in [24], [28], [29]. Comparisons by Trigueiros *et al.* in [28] and Wahid *et al.* [24] focused on hand gestures, as most of the approaches proposed in the literature (e.g., [20], [21], [23], [27]), whereas realistic daily life activities are considered by Sagha *et al.* in [29].

The above mentioned approaches consider the use of gestures for HRI in specific case studies and cannot be easily scaled to other applications. Our aim in this paper is to address the problem of gesture-based HRI from a general perspective, identifying the specific constraints of this application domain and proposing a pipeline that can be applied to any case study, with some fine tuning.

## III. PROPOSED ARCHITECTURE

We hereby describe the proposed pipeline for online gesture recognition in HRI. It addresses several constraints specific for gesture-based interaction with robots. A diagram representing the proposed architecture is depicted in Figure 1.

### A. Specific constraints of gesture recognition in HRI

As introduced in Section I, when dealing with gesture recognition in HRI, it is important that using gestures does not limit the user interacting with the robot, insofar gestures

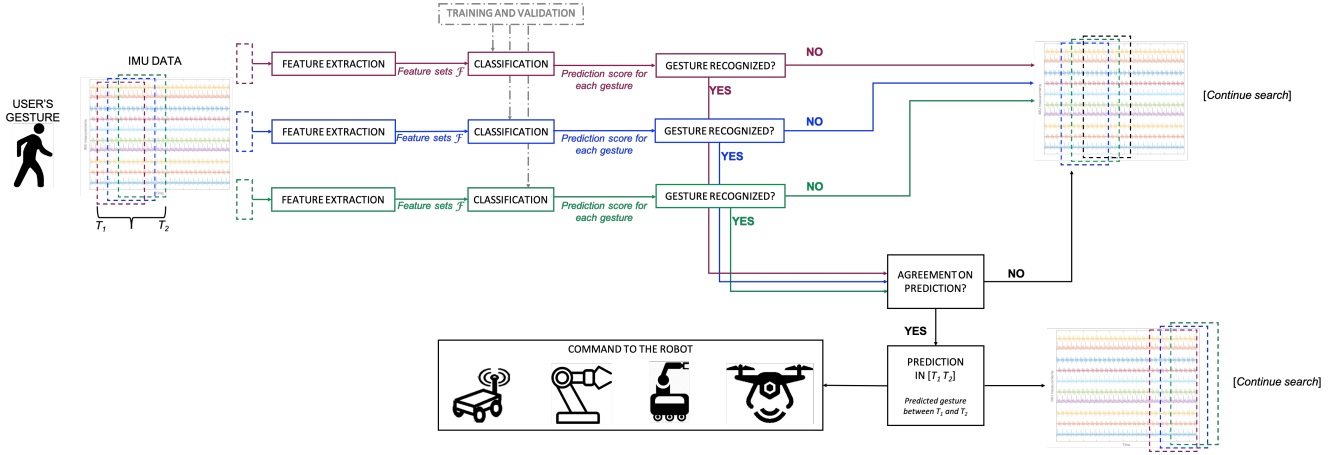


Fig. 1. Flowchart for the proposed pipeline for online gesture detection and classification in HRI. The proposed pipeline can be implemented in Robotic Operating System (ROS), which can deal with data acquisition, data processing, gesture recognition and communication with the robotic platform.

are not perceived as a slowdown for the interaction. As a result, the use of gestures in the context of HRI poses the following specific constraints for the problem of gesture recognition.

- C1 Gesture-based interaction should rely on a lean infrastructure that requires poor or no installation and does not limit the user's freedom to move around and/or with the robot.
- C2 Gesture recognition should be performed online and not introduce any perceivable delay between the execution of a gesture and its effect, meant as command to the robot. In other words, the algorithmic pipeline in charge of gesture recognition should be fast enough, requiring limited computational burden. Moreover, the architecture should guarantee immediate communication of a detected gesture to the robot.
- C3 It should be possible to detect gestures online and in dynamic conditions, while users are performing other ordinary activities and movements.
- C4 The system should be easily adapted to different users, thus requiring robust generalization capabilities.
- C5 While seeking to improve gesture recognition performance, it is particularly important to minimize the rate of false positives since they would initiate an unintended communication towards the robot.

It is noteworthy that some of these constraints are general for gesture-based interaction and apply also to other domains, as discussed, for example, in [14], [30]. Nevertheless, as regards HRI, they define the prerequisites for successful gesture-based HRI. These constraints are summarized in the left column of Table I, while the right column describes how the proposed system tackles them. More details are given in the following subsections.

### B. Overview of the proposed pipeline

To track the movements of the user, we consider inertial data recorded by a wrist-worn device, and thus focus on forearm gestures. This allows to comply with constraint C1 since wearable devices recording inertial data are not cumbersome to wear and are easily available on the market. Inertial data

are then analyzed, by considering sliding windows of fixed length, with one sample shift. Recorded data are continuously processed and the occurrence of a gesture is continuously verified, thus complying with C3. In particular, to search for gestures, for each sliding window a set of statistical features are computed for each of the measured signals. The extracted features represent the input for the classification stage. In this paper, we start considering different classical classification algorithms, such as  $K$ -nearest neighbors, support vector machines, random forests and neural networks. Among them, we then select the algorithm that best complies with the specific constraints for gesture-based HRI, in a validation use case. To select the most appropriate algorithm for the considered application domain, different performance metrics are considered, tailoring the severity of false positives in HRI, so as to comply with C5. Moreover, starting from the results achieved in our case study, the determination of the classification stage is driven by considerations related to the need for inter-subject robustness and computational burden, thus complying with C2 and C4. The right column of Table I summarizes how the proposed pipeline complies with the specific constraints for gesture-based HRI.

With reference to Figure 1, all the considered algorithms share the same general structure. They require previous training and validation with respect to the specific gestures of the considered use case. As output, each algorithm provides a score referred to the probability that one of the considered gestures occurred in the current sliding window. A threshold is set on such scores; the threshold can be tuned in order to balance between false positives and false negatives. If the confidence of the classifier is below such threshold, then the process moves to the following sliding window. Otherwise, in case the confidence of the classifier is above the threshold for at least one gesture, a decision is not taken immediately, but the system checks for the same condition also in the subsequent sliding windows. If the condition is met, then a decision is taken and the corresponding command is sent to the robot. The rationale behind this behavior is to improve the robustness of the detection phase. Since consecutive sliding windows are

shifted by one sample, it is likely that a gesture occupies more than one window, depending on the sampling frequency of the recording device. Hence, it is likely that, when a gesture is performed, it is recognized by the classification algorithm in several consecutive windows, whose number depends on the duration of the gesture and the device sampling frequency. On the other side, if a gesture is detected in a window only, it is likely to be a false positives and should be discarded.

#### IV. CLASSIFICATION STAGE

The classification stage of the proposed architecture consists in a machine learning module that receives as input the statistical features computed on the current sliding window and provides as output a decision whether a gesture is recognized or not. In this section, we overview the general setting of supervised machine learning, and we briefly describe four classic algorithms that have been tested in our experimental evaluation.

Broadly speaking, in a classification task, the goal is to predict the category  $y \in \mathcal{Y}$ , sometimes also named class, label, or *output*, of a given observed example  $x \in \mathcal{X}$ , which instead represents the *input* of the system. In a supervised setting, we are given a collection  $\mathcal{D}$  of  $m$  input/output pairs  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^m$  and we aim to fit a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  so that we can then predict the category  $\hat{y}$  of a given, novel (not previously seen) example  $\hat{x}$ . Clearly, different machine learning approaches exploit different definitions for function  $f$ , as well as different algorithms and techniques to learn such a function. To identify the best approach for the proposed HRI pipeline, we have compared classic machine learning approaches, namely  $K$ -nearest neighbors, support vector machines, random forests and neural networks. However, other different classification algorithms can be considered. In all these cases, we have considered a setup where each input instance  $x$  is described by a vector of real numbers, which are the characteristics, or *features* of that instance.

The parameters chosen to implement each approach are discussed in Sec. V-C, with reference to the considered datasets.

##### A. $K$ -Nearest Neighbors

Based on the concept of distance between examples, the  $K$ -Nearest Neighbors (KNN) classifier is not properly a learning algorithm. In fact, given a test example  $\hat{x}$  to be classified, the algorithm looks for the  $K$  examples in the training set that are the closest ones to  $\hat{x}$  according to a chosen metric (e.g., the Euclidean distance). The prediction is then performed via a majority voting procedure among the classes of the  $K$  nearest neighbors. Although very simple, this algorithm can work well in practice, when the distance computed on the feature vectors is highly discriminative of the target class. The parameter  $K$  defines the size of the neighborhood to consider for classification.

##### B. Support Vector Machines

Support Vector Machines (SVM) are another classic method for supervised classification in machine learning. Considering

binary classification, SVMs are trained to learn a discriminative function that best separates positive and negative examples, with the maximum possible margin [31]. Given a collection of  $m$  training samples, such a discriminative function is computed as:

$$f(x) = \sum_{i=1}^m \alpha_i \mathcal{K}(x_i, x) + b \quad (1)$$

where  $\alpha_i$  are the learnable model parameters, and the kernel function  $\mathcal{K}(\cdot, \cdot)$  aims to measure the similarity between examples. Function  $\mathcal{K}$  is, in general, a nonlinear function, so that the classifier can model nonlinear dependencies between features and class. Decision function  $f$  only depends on those training examples whose corresponding coefficient  $\alpha_i$  is different from zero: these are called *support vectors*. In general terms, the decision function  $f$  defines a hyperplane, which constitutes the decision boundary used to classify the data points. Data points falling on either side of the hyperplane can be attributed to different classes. This idea can be extended to the case of multiclass classification, breaking down the multiclassification problem into multiple binary classification problems. Two approaches can be selected to this end: in the one-versus-one approach, a binary classifier is set per each pair of classes, regardless of the other classes; in the one-versus-rest, a binary classifier is set per each class, to distinguish it from the rest of data.

A commonly used kernel function for SVM classifiers is the Radial Basis Function (RBF). The RBF kernel function for two points  $x_1$  and  $x_2$  computes the similarity or how close they are to each other, as a function of their Euclidean distance. It is specified by means of two parameters:  $C$  and  $\gamma$ . While  $C$  sets a trade off between misclassification of training examples and simplicity of the decision surface,  $\gamma$  defines the width of the radial functions.

##### C. Random Forests

A Random Forest [32] (RF) consists in a collection of individual Decision Trees (DTs) [33], whose predictions are combined typically through a voting process. A DT inductively learns a set of explainable classification rules by imposing conditions of the values of the features describing the examples. When creating a RF, each DT is trained from a distinct set of  $n$  examples randomly sampled from the original training set, and by testing only  $m$  out of  $M$  features at each node in the tree. The process of combining the outcome of individual classifiers into a single prediction is usually named as an *ensemble* approach, and it is known to typically improve the performance of the overall system, as well as to reduce overfitting. Each individual DT in the RF produces a class prediction and the class with most votes becomes the prediction of the model, as in classic ensemble approaches.

##### D. Neural networks

A Neural Network [34] (NN) is a nonlinear function transforming a set of input variables in a set of output variables via a set of adjustable parameters. In particular, a NN is a

combination of nonlinear basis functions. Each basis function is itself a nonlinear function (called activation function) of a linear combination of the inputs, and the coefficients in such combination are adaptive weights, that can be learned to fit the training data. Several layers of adaptive weights can be stacked to form a deep network.

## V. EXPERIMENTAL IMPLEMENTATION

In this section we present the experimental setting that was considered to implement the pipeline shown in Fig. 1. We first introduce the validation use case and describe different sets of features extracted from the raw signals; then we illustrate the adopted training and evaluation procedure, and finally the performance metrics employed at prediction time.

### A. Use case

To validate the proposed pipeline, we considered the gesture-based HRI approach introduced in [2]. In particular, the scenario consists in recognizing the following  $N_{gest} = 5$  gestures, depicted in Figure 2:

- 1) *Up*: sharp movement upwards in a plane parallel to the sagittal one,
- 2) *Down*: sharp movement downwards in a plane parallel to the sagittal one,
- 3) *Circle*: movement in a circular shape in a plane parallel to the frontal one,
- 4) *Left*: sharp movement to the left, from sagittal plane to frontal one, in a plane parallel to the transverse one,
- 5) *Right*: sharp movement to the right, from sagittal plane to frontal one, in a plane parallel to the transverse.

Examples can be seen in the multimedia attachment to [2]. These gestures are meant to be performed with the right arm, with the subject wearing an IMU on the right wrist.

As our input device we considered a commercial multi-purpose smartwatch, namely the Samsung Gear S device. Data used in the analysis presented in this paper consist of triaxial inertial measurements recorded by the smartwatch, namely  $x \in \mathbb{R}^{10}$ , and include timestamp, angular velocities and linear accelerations (raw and with automatically compensated gravity). Data are accessed by means of a Tizen interface and are provided on an uneven sampling grid. On average, approximately 25 samples per second are provided.

As regards the overall software architecture, data are recorded via Tizen and sent via Wi-Fi to an external computer, for ease of implementation. The classification stage is implemented in Python, using the `scikit-learn` library [35] and `Tensorflow 2.0.0` [36].

ROS can be, then, used for implementing the remaining of the pipeline [37]. It is an open-source hardware-independent middleware widely used in robotics and consists in a set of software libraries and tools that allow communication with robots, actuators, sensors and other devices commonly used in robotic applications. Since it supports Python, it can deal with data acquisition, data processing, gesture recognition and communication with the robotic platform. To this end, a publish-subscribe pattern can be used that can efficiently handle communication events, such as arriving messages and inform the robot about the detection of a gesture.



Fig. 2. Gestures considered in this study: from left to right, *Up*, *Down*, *Circle*, *Left* and *Right*. Animated examples can be seen in the multimedia attachment to [2].

TABLE II  
COMPOSITION OF THE TRAINING AND TEST DATASET. THE DURATION IS EXPRESSED IN AVERAGE NUMBER OF SAMPLES, WITH ITS STANDARD DEVIATION.

	Label	Instances	Duration ([samples])
Training	<i>Up</i>	60	$29.5 \pm 2.9$
	<i>Down</i>	60	$29.9 \pm 3.3$
	<i>Circle</i>	60	$32.6 \pm 4.8$
	<i>Left</i>	60	$25.9 \pm 2.5$
	<i>Right</i>	60	$28.9 \pm 2.5$
	$N_{\text{train}}$	300	
Test		Instances	Subject
	$N_{\text{test}}$	200	90 by $S_1$ 10 each by $S_2, \dots, S_{12}$

### B. Feature extraction

Starting from the data provided by the wrist-worn device, synthetic features are extracted and passed to the classification stage as input data. We considered three sets of features, which were compared as possible different inputs for the classification algorithms. The first two feature sets, denoted in the following as  $\mathcal{F}_1$  and  $\mathcal{F}_2$ , were defined as follows. Set  $\mathcal{F}_1$  includes standard statistics computed in the domain of time, for each inertial quantity: mean value, standard deviation, maximum value, and minimum value. These are customarily used in many gesture recognition applications (for example, see work in [38]). Following [24], [26], set  $\mathcal{F}_2$  includes mean absolute value, count of slope sign changes, count of zero crossings and waveform length, as follows:

- mean absolute value:  $\frac{1}{L} \sum_{k=1}^L |x_k|$
- slope sign change:  $(x_k - x_{k-1})(x_{k+1} - x_k) < 0$
- zero crossing:  $x_k x_{k+1} < 0$
- waveform length:  $\sum_{k=2}^L |x_k - x_{k-1}|$

where  $L$  is the number of samples in a sliding window.

Set  $\mathcal{F}_3$  is obtained from the union of  $\mathcal{F}_1$  and  $\mathcal{F}_2$  and hence consists of 8 features.

As a result, for each feature set, each input sequence for classification (both for training and test data sets) is represented by 36 features for  $\mathcal{F}_1$  and  $\mathcal{F}_2$  and 72 features for  $\mathcal{F}_3$ , since the above mentioned features were computed for each of the nine inertial quantities measured by the smartwatch.

### C. Training and evaluation

All the considered classification algorithms were trained by considering a labelled data set of 300 gestures, which included 60 trials per gesture. Such gestures were all performed by the same subject, denoted in the following as  $S_0$ , standing in

TABLE III  
ARCHITECTURES FOR NEURAL NETWORKS CONSIDERED IN THE CLASSIFICATION STAGE. LEARNING RATE WAS SET TO 1E-3, THE NUMBER OF TRAINING EPOCHS TO 500, AND THE BATCH SIZE TO 100 FOR ALL THE CONSIDERED ARCHITECTURES.

	Layer	Type	Nodes	Activation function
<b>Model 1</b>	Input	Dense	100	ReLU
	Hidden	Dense	50	ReLU
	Hidden	Dense	25	ReLU
	Output	Dense	6	Softmax
<b>Model 2</b>	Input	Dense	100	ReLU
	Hidden	Dense	50	ReLU
	Output	Dense	6	Softmax
<b>Model 3</b>	Input	Dense	50	ReLU
	Hidden	Dense	25	ReLU
	Output	Dense	6	Softmax

steady state. The data set contains inertial samples that refer to the execution of the gestures, only: samples between two consecutive gestures were manually excluded. The composition of the training data set is detailed in Table II. In order to tune the hyperparameters of the different classifiers, an inner  $k$ -fold cross-validation on the training set was conducted.

Based on preliminary experiments, in the analyses reported in this paper, for KNN we tested values for  $K = 5, 7, 10, 15$ . Since the two hyperparameters,  $C$  and  $\gamma$ , have to be jointly chosen, we exploited a grid search with an internal cross-validation<sup>1</sup> as customary in this kind of applications. We finally selected  $C = 0.1$  and  $\gamma = 10^{-6}$ . As for RF, we considered a total of 100 trees in our model. For NN, we considered three different architectures, summarized in Table III. Each network was trained using the Adam optimizer [39] with batch normalization, and the hyperparameters (learning rate, momentum and number of training epochs) were set with Bayesian optimization<sup>2</sup>.

To test the algorithms under analysis, we considered streams of inertial data containing gestures. In other words, consecutive gestures were separated by a non-constant number of samples in which the user moved in a natural manner, as described below. The goal was to replicate real life scenarios in which the subject is free to move while interacting with a robot. Occurrences of gestures were manually annotated to serve as ground truth. To deal with online classification, and following [2], input data were provided to the algorithms by considering sliding windows of length  $L = 35$  samples, with one sample shift. As a result, the test set includes  $N_{\text{test}} = 200$  gestures performed with the subjects moving in a natural manner between consecutive gestures (e.g., walking, waving, drinking). These gestures were performed by 12 subjects, namely  $S_1, S_2, \dots, S_{12}$ , different from the one contributing to the training set. They were not told the movements to execute between gestures, but were left free to chose. An equal number of the five gestures under consideration was included. The composition of the test set is detailed in Table II.

All the algorithms were implemented in a multiclass configuration. This implies that the detection of a gesture is accepted

only if the confidence of the classifier is above a certain threshold that can be tuned in order to balance between false positives and false negatives. To this end, different thresholds have been considered. In the case of RF, performance has been computed considering three different thresholds, namely  $\sigma_1, \sigma_2$  and  $\sigma_3$ , on the predicted probabilities for each class, to be used as a confidence level for gesture recognition. Specifically, for each input sequence,  $N_{\text{gest}}$  predicted probabilities  $P_i, i = 1, \dots, N_{\text{gest}}$  are given as output, each representing the likelihood for that sequence to contain one of the considered gestures. For each threshold  $\sigma_j$ , gesture  $i$  is, then, detected if  $\max P_i > \sigma_j$ , with  $j = 1, 2, 3$ . In a similar manner, for SVM a one-versus-one approach for multiclass classification is selected and five thresholds are considered, denoted as  $\theta_k, k = 1, \dots, 5$ . A score for each sample in relation to each gesture is provided and the threshold  $\theta_k$  is set with respect to such score. The results achieved for RF and SVM with the different thresholds  $\sigma_j$  and  $\theta_k$  ( $j = 1, 2, 3$  and  $k = 1, \dots, 5$ ) are reported in Section VI<sup>3</sup>. A threshold is considered also in the case of KNN: it was set to  $\tau = 1$  for any value considered for  $K$ , and, hence, is not reported in Section VI. Alternatively, when dealing with a specific case study, the optimal threshold can be set with a cross-validation step.

Finally, as introduced in Section III-B, a gesture was detected only if the same outcome was predicted by a classifier throughout several consecutive sliding windows. Given the average duration of the considered gestures (reported in Table II) and the approximate sampling frequency of the smartwatch of about 25 Hz, we set such number of consecutive sliding windows to 25.

#### D. Performance metrics

A set of established metrics have been computed to compare the performance of the considered algorithms [40], [41]. The selection of these metrics was guided by the need to guarantee high recognition performance in general, while taking into account specific HRI constraints as well. We name *Recall* the percentage of gestures in the test set that have been correctly detected by the classifier, whereas *Precision* is the percentage of predicted gestures that are correct<sup>4</sup>. These two measures account for different kinds of error: namely, *Recall* considers false negatives since it takes into account those gestures that are not recognized by the system, while *Precision* considers false positives, which are gestures that the system wrongly recognizes. The latter are much more dangerous in HRI, where a gesture starts a command to the robot, as highlighted by constraint C5 in Section III-A. Hence, in the need for a compromise between performance on *Precision* and *Recall*, it is important that *Precision* is as high as possible. Low *Precision* means that the risk of providing unintended commands to the robot is high, which is clearly unacceptable. Conversely, a false negative in HRI means that the user has

<sup>3</sup>We set  $\sigma_1 = 0.55, \sigma_2 = 0.60$  and  $\sigma_3 = 0.65$ , and  $\theta_k = 4.240 + k \cdot 0.005$ , with  $k = 1, \dots, 5$ .

<sup>4</sup>In general, in multiclass problems, *Precision* and *Recall* are more frequently computed on a per-class basis, while in our application scenario we prefer to report classifier-level metrics in order to focus on the global performance of the system.

<sup>1</sup>We used the `GridSearchCV` class of `scikit-learn`.

<sup>2</sup>We used the library `Keras Tuner`.



performed a gesture to command the robot, but this command was not received by the robot. Although annoying for the user, this circumstance is less dangerous than the previous one, since it does not have direct consequences in terms of unintended behaviour of the robot.

*Precision* and *Recall* are usually combined in the  $F_\beta$  score, defined as:

$$F_\beta \text{ score} = (1 + \beta^2) \frac{\text{Recall} \cdot \text{Precision}}{\text{Recall} + \beta^2 \text{Precision}} \quad (2)$$

The parameter  $\beta$  quantifies the importance of *Recall* over *Precision* and is typically set to 1. Since, in our context, it is more important that *Precision*, rather than *Recall*, is large, we set  $\beta = 0.5$ , to attribute more importance to *Precision* over *Recall* [40], [41]. Thus, in our analysis we consider  $F_{0.5}$  score, following constraint C5.

In order to better highlight the difference between the two main error categories (i.e., missed or misclassified gestures) and comply with constraint C5, we introduce two additional metrics. The *Misclassification Gesture Rate* (*MGR*) is the percentage of real gestures that are assigned to the wrong category, whereas the *Undetected Gesture Rate* (*UGR*) is the percentage of real gestures that are not detected. We hereby remark that  $\text{Recall} + \text{MGR} + \text{UGR} = 1$ . Following the line of reasoning above, it is more important to minimize *MGR* than *UGR*, since the former accounts for unintended swaps among commands to the robot, while the latter refers to the need, for the user, to repeat a gesture.

Finally, to guarantee real-time gesture recognition (see constraint C2 in Section III-A), algorithms have been compared also in terms of computational burden. To this end, the classification time has been considered. With respect to Figure 1, this amounts to considering the time required from the selection of a sliding window (dashed rectangles on the left) to the output of classification algorithms for that window (boxes “*Gesture recognized?*” in the figure).

## VI. ANALYSIS OF THE EXPERIMENTAL DATA

In this section, we analyze and discuss the experimental data collected implementing the proposed pipeline in the use case introduced in Section V-A. In particular, besides comparing the performance of the considered four classical machine learning approaches with different feature sets, we will focus on the generalization capabilities of the classifiers across different subjects, with varying levels of familiarity with the selected gestures and inclusion in training set. The idea behind this analysis is to assess the amount of user’s contribution, meant in terms of participation in the training set or gained acquaintance with the gestures, that is needed to achieve satisfactory performance for gesture recognition. To this end, we consider three different conditions: *a)* subjects with different familiarity with the considered gestures, whose data were not included in the training set; *b)* a subject with previous experience on the use of the selected gestures, but whose data were not included in the training set; *c)* the same subject as in the previous condition, whose data contributed, in small part, to the training set.

*a) Multiple subjects, expert and novel, not included in the training set:* Table IV reports the performance of the classification stage, for the different algorithms considered, in the case of subjects with and without prior familiarity with the considered set of gestures. With reference to Table II, numbers in Table IV refer to 200 gestures performed by all the subjects in the test set (from  $S_1$  to  $S_{12}$ ). We remark that these subjects were not represented in the training set and had received different amount of instructions on how to execute the gestures. Specifically, subject  $S_1$  had previous experience with the considered set of gestures, while all the others had not, and received little training soon before recording the test set. The table shows that satisfactory results can be achieved with RF (with features  $\mathcal{F}_3$ ) for *Precision* and *MGR*, which are the most relevant metrics for HRI applications, as discussed in Section V-D. Nevertheless, performance achieved for the other metrics are extremely low. The performance achieved with the other algorithms are quite poor, too.

*b) Expert subject, not included in the training set:* To improve the performance of gesture recognition, we limited the test set to gestures performed by an expert subject, who had prior experience in using them, namely  $S_1$ . Indeed, although the considered gestures consist in simple and natural movements, they are meant as sharp movements, as described in Section V-A, and require some experience to get familiar with movement speed and initial and final positions. Table V reports the results achieved in this condition. As for the previous analysis, the subject contributing to the training set was not included in the test set. The table shows that a notable increase in performance is achieved. However, classification performance is still quite low for all the considered algorithms. In particular, while very good performance can be achieved in terms of *Precision* and *MGR*, *Recall* is still quite low. Satisfactory *Recall* is achieved with NNs (Model 1 and features  $\mathcal{F}_1$ ), at the expenses of low *Precision* and  $F_{0.5}$  score. There is no algorithm, among those considered, that returns satisfactory results for both *Recall* and *Precision*.

*c) Expert subject, included in the training set:* As a further attempt to improve the performance of gesture recognition, we increased the training set adding some gestures by the same subject in the test set. In particular, with reference to Table II, the training set was augmented adding 100 new gestures (20 per type) by subject  $S_1$ . Regarding the test set, we considered 90 gestures performed by  $S_1$  (i.e., the same as condition *b)*). Table VI reports the performance achieved in this condition. The best classification performance are achieved with SVM with features  $\mathcal{F}_1$ , which proves successful for all the considered metrics. All the thresholds  $\theta_k$ , with  $k = 1, \dots, 4$  return similar results, while  $\theta_5$  is less performing. Considering  $F_{0.5}$  score, which is the weighted combination of *Recall* and *Precision*, the threshold  $\theta_1$  returns the highest value. RF provides slightly better performance for *Precision* and *MGR* than SVM, but  $F_{0.5}$  score for RF is quite low. Good performance, in terms of *Recall*, *UGR* and *MGR*, is also achieved with NNs with features  $\mathcal{F}_1$ , although *Precision* and  $F_{0.5}$  score are quite low. In summary, it is possible to argue that better performance can be achieved when using SVM with threshold  $\theta_1$  and features  $\mathcal{F}_1$ .



TABLE IV

CLASSIFICATION PERFORMANCE FOR GESTURES PERFORMED BY MULTIPLE SUBJECTS, EXPERT AND NOVEL, NOT INCLUDED IN THE TRAINING SET. (*R*: Recall; *P*: Precision; *UGR*: Undetected Gesture Rate; *MGR*: Misclassification Gesture Rate;  $F_{0.5}$ :  $F_{0.5}$  score.)

	Features $\mathcal{F}_1$					Features $\mathcal{F}_2$					Features $\mathcal{F}_3$				
	<i>R</i>	<i>P</i>	<i>UGR</i>	<i>MGR</i>	$F_{0.5}$	<i>R</i>	<i>P</i>	<i>UGR</i>	<i>MGR</i>	$F_{0.5}$	<i>R</i>	<i>P</i>	<i>UGR</i>	<i>MGR</i>	$F_{0.5}$
<b>RF</b> ( $\sigma_1$ )	29.1	72.2	62.8	8.1	33.0	20.9	31.1	55.6	23.5	22.4	27.0	67.9	62.8	10.2	30.7
<b>RF</b> ( $\sigma_2$ )	20.9	91.1	77.0	2.1	24.7	16.8	63.5	74.0	9.2	19.7	21.0	90.7	76.8	2.2	24.8
<b>RF</b> ( $\sigma_3$ )	17.3	91.9	81.2	1.6	20.7	8.2	69.6	88.3	3.5	9.9	12.8	<b>92.6</b>	86.2	<b>1.0</b>	15.4
<b>KNN</b> ( $K = 5$ )	44.4	56.5	33.7	21.9	<b>46.4</b>	21.9	62.3	66.8	11.3	25.2	29.6	68.2	59.7	10.7	33.4
<b>KNN</b> ( $K = 7$ )	43.4	60.7	38.3	18.3	46.0	14.3	60.9	77.6	8.1	16.9	25.5	70.4	65.8	8.7	29.2
<b>KNN</b> ( $K = 10$ )	40.3	67.5	44.9	14.8	43.8	11.2	62.9	82.7	6.1	13.4	20.9	75.9	73.0	6.1	24.5
<b>KNN</b> ( $K = 15$ )	37.2	75.3	52.0	10.8	41.4	5.6	55.0	90.3	4.1	6.8	14.3	71.8	80.6	5.1	17.0
<b>SVM</b> ( $\theta_1$ )	51.5	35.1	15.3	33.2	47.1	44.4	26.0	10.7	44.9	38.9	50.0	27.9	<b>9.7</b>	40.3	43.2
<b>SVM</b> ( $\theta_2$ )	48.5	33.5	17.9	33.6	44.5	42.3	24.5	<b>9.7</b>	48.0	37.0	47.4	28.4	10.7	41.9	41.9
<b>SVM</b> ( $\theta_3$ )	45.9	34.1	20.9	33.2	42.9	37.8	23.6	12.2	50.0	33.7	41.8	27.4	11.2	46.9	37.9
<b>SVM</b> ( $\theta_4$ )	43.4	44.0	31.6	25.0	43.5	34.2	23.0	11.7	54.1	31.2	37.2	26.6	14.3	48.5	34.5
<b>SVM</b> ( $\theta_5$ )	39.8	70.3	50.0	10.2	43.6	32.1	24.0	15.8	52.1	30.1	32.7	28.4	24.5	42.8	31.7
<b>NN</b> (Model 1)	<b>51.9</b>	35.3	28.8	19.3	37.3	35.2	25.4	35.0	29.8	26.7	47.2	32.0	27.3	25.5	33.9
<b>NN</b> (Model 2)	50.1	29.6	26.6	23.3	31.9	41.5	34.7	35.2	23.3	35.4	49.4	33.2	26.1	24.5	35.2
<b>NN</b> (Model 3)	46.8	29.2	27.2	26.0	31.3	38.2	29.2	36.8	25.0	30.5	47.2	30.3	25.9	26.9	32.4

TABLE V

CLASSIFICATION PERFORMANCE FOR GESTURES PERFORMED BY AN EXPERT SUBJECT, NOT INCLUDED IN THE TRAINING SET. (*R*: Recall; *P*: Precision; *UGR*: Undetected Gesture Rate; *MGR*: Misclassification Gesture Rate;  $F_{0.5}$ :  $F_{0.5}$  score.)

	Features $\mathcal{F}_1$					Features $\mathcal{F}_2$					Features $\mathcal{F}_3$				
	<i>R</i>	<i>P</i>	<i>UGR</i>	<i>MGR</i>	$F_{0.5}$	<i>R</i>	<i>P</i>	<i>UGR</i>	<i>MGR</i>	$F_{0.5}$	<i>R</i>	<i>P</i>	<i>UGR</i>	<i>MGR</i>	$F_{0.5}$
<b>RF</b> ( $\sigma_1$ )	43.0	84.1	51.2	5.8	47.7	30.2	28.6	37.2	32.6	29.9	41.9	80.0	47.7	10.4	46.3
<b>RF</b> ( $\sigma_2$ )	34.9	90.9	61.6	3.5	39.8	24.4	63.6	61.6	14.0	27.8	33.7	90.6	62.8	3.7	38.6
<b>RF</b> ( $\sigma_3$ )	30.2	92.9	67.4	2.4	34.9	16.3	73.7	77.9	5.8	19.3	23.3	95.2	75.6	1.1	27.4
<b>KNN</b> ( $K = 5$ )	59.3	60.0	23.3	17.4	59.4	33.7	82.9	62.8	3.5	38.3	46.5	83.3	50.0	3.5	51.0
<b>KNN</b> ( $K = 7$ )	58.1	64.9	29.1	12.8	59.4	20.9	90.0	77.9	1.2	24.7	40.7	89.7	58.1	1.2	45.7
<b>KNN</b> ( $K = 10$ )	57.0	77.8	34.9	8.1	<b>60.2</b>	17.4	93.8	81.4	1.2	20.8	33.7	<b>100.0</b>	66.3	<b>0.0</b>	38.9
<b>KNN</b> ( $K = 15$ )	53.5	88.5	43.0	3.5	58.1	9.3	<b>100.0</b>	90.7	<b>0.0</b>	11.4	24.4	100.0	75.6	<b>0.0</b>	28.8
<b>SVM</b> ( $\theta_1$ )	68.6	31.9	<b>3.5</b>	27.9	55.8	57.0	27.8	4.7	38.3	47.1	68.6	33.5	4.7	26.7	56.7
<b>SVM</b> ( $\theta_2$ )	66.3	33.3	9.3	24.4	55.3	52.3	25.4	4.7	43.0	43.2	62.8	31.6	5.8	31.4	52.4
<b>SVM</b> ( $\theta_3$ )	60.5	39.1	20.9	18.6	54.5	46.5	24.4	5.8	47.7	39.4	53.5	29.9	7.0	39.5	46.2
<b>SVM</b> ( $\theta_4$ )	57.0	56.3	33.7	9.3	56.8	40.7	22.6	5.8	53.5	35.1	45.3	27.1	9.3	45.4	40.0
<b>SVM</b> ( $\theta_5$ )	51.2	83.0	47.7	1.1	55.4	40.7	23.8	5.8	53.5	35.6	41.9	28.1	12.8	45.3	38.1
<b>NN</b> (Model 1)	<b>83.6</b>	44.5	9.8	6.6	48.8	41.8	25.3	32.2	26.0	27.5	78.4	44.9	8.1	13.5	48.8
<b>NN</b> (Model 2)	76.8	32.5	10.0	13.2	36.6	57.8	37.9	19.9	22.3	40.5	80.7	43.4	5.4	13.9	47.6
<b>NN</b> (Model 3)	70.5	34.8	13.3	16.2	38.6	52.2	32.1	25.5	22.3	34.8	79.0	42.6	5.6	15.4	46.6

TABLE VI

CLASSIFICATION PERFORMANCE FOR GESTURES PERFORMED BY AN EXPERT SUBJECT, INCLUDED IN THE TRAINING SET. (*R*: Recall; *P*: Precision; *UGR*: Undetected Gesture Rate; *MGR*: Misclassification Gesture Rate;  $F_{0.5}$ :  $F_{0.5}$  score.)

	Features $\mathcal{F}_1$					Features $\mathcal{F}_2$					Features $\mathcal{F}_3$				
	<i>R</i>	<i>P</i>	<i>UGR</i>	<i>MGR</i>	$F_{0.5}$	<i>R</i>	<i>P</i>	<i>UGR</i>	<i>MGR</i>	$F_{0.5}$	<i>R</i>	<i>P</i>	<i>UGR</i>	<i>MGR</i>	$F_{0.5}$
<b>RF</b> ( $\sigma_1$ )	61.6	98.1	38.4	<b>0.0</b>	66.6	36.0	81.6	58.1	5.9	40.6	45.3	92.9	52.3	2.4	50.5
<b>RF</b> ( $\sigma_2$ )	52.3	<b>100.0</b>	47.7	<b>0.0</b>	57.8	29.1	<b>100.0</b>	70.9	<b>0.0</b>	33.9	30.2	96.3	68.6	1.2	35.0
<b>RF</b> ( $\sigma_3$ )	40.7	<b>100.0</b>	59.3	<b>0.0</b>	46.2	24.4	<b>100.0</b>	75.6	<b>0.0</b>	28.8	27.9	<b>100.0</b>	72.1	<b>0.0</b>	32.6
<b>KNN</b> ( $K = 5$ )	74.4	52.9	18.6	7.0	68.8	25.6	53.7	69.8	4.6	28.6	40.7	62.5	54.7	4.6	43.8
<b>KNN</b> ( $K = 7$ )	73.3	56.3	19.7	7.0	69.1	14.0	50.0	84.9	1.1	16.3	33.7	64.4	64.0	2.3	37.3
<b>KNN</b> ( $K = 10$ )	73.3	64.9	23.2	3.5	71.4	11.6	62.5	87.2	1.2	13.9	26.7	79.3	73.3	<b>0.0</b>	30.8
<b>KNN</b> ( $K = 15$ )	64.0	75.3	36.0	<b>0.0</b>	65.9	7.0	85.7	93.0	<b>0.0</b>	8.5	18.6	94.1	81.4	<b>0.0</b>	22.2
<b>SVM</b> ( $\theta_1$ )	<b>96.5</b>	91.2	<b>2.3</b>	1.2	<b>95.4</b>	66.3	47.5	7.0	26.7	61.4	88.4	67.9	4.6	7.0	83.3
<b>SVM</b> ( $\theta_2$ )	94.2	95.3	5.8	<b>0.0</b>	94.4	65.1	51.9	8.2	26.7	61.9	88.4	74.5	5.8	5.8	85.2
<b>SVM</b> ( $\theta_3$ )	94.2	96.4	5.8	<b>0.0</b>	94.6	62.8	54.0	10.5	26.7	60.8	84.9	76.0	9.3	5.8	83.0
<b>SVM</b> ( $\theta_4$ )	94.2	97.6	5.8	<b>0.0</b>	94.8	62.8	55.1	10.5	26.7	61.1	80.2	77.5	14.0	5.8	79.7
<b>SVM</b> ( $\theta_5$ )	83.7	98.6	16.3	<b>0.0</b>	86.3	59.3	58.0	17.4	23.3	59.0	77.9	83.8	18.6	3.5	79.0
<b>NN</b> (Model 1)	94.6	62.9	3.4	2.0	67.3	55.0	41.0	24.4	20.6	42.9	90.8	63.4	6.3	2.9	67.0
<b>NN</b> (Model 2)	94.8	54.9	<b>2.3</b>	2.9	59.8	56.1	43.8	26.4	17.5	45.5	88.9	59.8	4.9	6.2	63.3
<b>NN</b> (Model 3)	94.6	62.9	3.4	2.0	67.1	58.9	42.7	17.6	23.5	45.0	87.0	50.6	5.5	7.5	54.9

Building upon this result, we tested this classifier considering gestures performed during intense motion conditions or moderate physical activity. To this end, subject  $S_1$  was asked to perform  $N_{\text{test}}^{\text{IM}} = 20$  gestures during running and leg skipping sessions and arm circumductions. These movements were intended as spurious noise for the classification algorithms. Nevertheless, SVM, with threshold  $\theta_1$  and features  $\mathcal{F}_1$ , has returned quite good performance also in this extreme motion condition:  $R = 100\%$ ,  $P = 57.1\%$ ,  $UGR = 0.0\%$ ,  $MGR = 0.0\%$ ,  $F_{0.5} \text{ score} = 87.0\%$ . While *Recall* is very high, the effect of intense motion can be seen in lower *Precision*, due to an increased number of false positive gestures.

#### A. Discussion

Tables IV, V and VI report the classification performance achieved considering a possible implementation of the proposed pipeline. To this end, we considered four standard classification algorithms and implemented them with different parameters. The achieved results show that, to achieve satisfactory classification performance, it is needed that the user receives some sort of training on how to execute the selected gestures. In particular, even if natural and easy movements are considered, the user should practise them, in order to get familiar with range and speed of motion and other specific features. To this end, including, in the classification algorithm, some amount of user specific training is beneficial. It is noteworthy that this is not in contrast with constraint  $C4$  of Table I. Indeed, in Table VI we considered the case that the algorithms were trained with data from a subject not in the test set in large part and, in small part, with data from the same subject in the test set. This is customary in devices running, for example, algorithms for voice or handwritten text recognition, which usually require some fine training by the user in charge.

As Table VI shows, if the user is familiar with the selected set of gestures and algorithms are trained also with her/his data, classification performances become highly satisfactory and allow an efficient use of gestures in HRI. This is confirmed also by the classification performance achieved in the case of intense motion condition: although this represents a quite unusual condition in HRI, these results show that gestures can be used without restrictions to any other tasks the user might be carrying out while interacting with a robot.

#### B. Analysis of computational burden

The computational burden of the proposed pipeline was assessed to verify its applicability to real-time gesture recognition (constraint  $C2$  in Table I). To this end, we computed the average classification time required by each of the considered algorithms. The measured times, averaged over 5,000 sliding windows, are reported in Table VII.

We hereby make two observations. First, according to the decision rule presented in Section III-B and Section V-C, a gesture is detected only if it is recognized in 25 consecutive sliding windows. In other words, the recognition of a gesture requires that 25 windows elapsed. Hence, since times in Table VII refer to classification for a single time window, the recognition of a gesture implies a delay that is 25 times

TABLE VII  
AVERAGE COMPUTATIONAL TIME FOR GESTURE CLASSIFICATION ON A SINGLE TIME WINDOW.

Model	RF	KNN	SVM	NN	NN	NN
				Model 1	Model 2	Model 3
Time (ms)	23	3	2	16	17	18

the one reported in the table, from the beginning of its execution. Second, for ease of computation, we implemented the classification stage on an external computer<sup>5</sup> running Python. Computation times reported in Table VII refer to this setting. In a real world operational setting, an hardware-driven software implementation should be considered, possibly relying on the computational capacity of the robot.

Table VII shows that all the algorithms are quite fast, with SVM being the fastest among those considered. Even considered the need to process 25 consecutive windows, the delay introduced by the classification operation does not affect the fluency of gesture-based interaction. As a result, the selection of the most suited algorithm for the proposed classification can be guided by classification performances only.

## VII. CONCLUSIONS

In this paper we considered the problem of forearm gesture recognition for HRI. The ultimate goal is that of providing commands to robots by means of intuitive gestures. To this end, we proposed a pipeline for gesture recognition specifically designed for HRI applications. As input data we consider wrist inertial movements, which can be recorded with any commercial device mounting an IMU. Gesture detection and classification is performed by a classification stage that relies on machine learning algorithms. To this end, we compared the performance of several machine learning algorithms used in classification problems: RF, KNN, SVM, and NN. An extensive evaluation was performed, including also an analysis of computational burden for real time gesture recognition. Different performance metrics were introduced, to provide a thorough assessment of algorithms and to highlight the specific needs of gesture recognition in HRI context. In particular, we highlighted that the cost of false positives and misclassifications (wrongly recognized gestures) is much higher than that of false negatives (gestures not detected at all), thus motivating the need for models with a high precision, even though at the cost of a lower recall.

An evaluation use case was selected to show an implementation of the proposed pipeline. It consisted of five gestures, recorded with a smartwatch. Twelve subjects were included in the validation, considering different confidence with the selected gestures. Firstly, we analyzed whether a general, meant as opposite to user tailored, training of classification algorithms could be sufficient to achieve satisfactory recognition performance. Then, to improve recognition performance, we considered the need of training algorithms also with a small amount of data from the subject in charge.

<sup>5</sup>The computer used for this analysis embeds an Intel i5 3.2 GHz CPU and 8 GB RAM, and it runs Ubuntu 16.04.7.

As concluding remarks, it is noteworthy that the proposed pipeline is general. Indeed, it can be used implementing other classification algorithms, or with other parameters. Additionally, it can be used to interact with any kind of robots and with other sets of gestures.

### ACKNOWLEDGEMENT

The authors would like to thank Anna Bellodi and Ilias Maite for their contribution to the preliminary analysis of the proposed approach.

### REFERENCES

- [1] A. Chaudhary, J. L. Raheja, K. Das, and S. Raheja, "Intelligent approaches to interact with machines using hand gesture recognition in natural way: a survey," *International Journal of Computer Science & Engineering Survey (IJCES)*, vol. 1, no. 2, pp. 122–133, 2013.
- [2] V. Villani, L. Sabattini, G. Riggio, C. Secchi, M. Minelli, and C. Fantuzzi, "A natural infrastructure-less human-robot interaction system," *IEEE Robot. Automat. Lett.*, vol. 2, no. 3, pp. 1640–1647, 2017.
- [3] V. Villani, L. Sabattini, N. Battilani, and C. Fantuzzi, "Smartwatch-enhanced interaction with an advanced troubleshooting system for industrial machines," *IFAC-PapersOnLine*, vol. 49, no. 19, pp. 277–282, 2016.
- [4] P. Neto, M. Simão, N. Mendes, and M. Safeea, "Gesture-based human-robot interaction for human assistance in manufacturing," *The International Journal of Advanced Manufacturing Technology*, vol. 101, no. 1–4, pp. 119–135, 2019.
- [5] H.-D. Yang, A.-Y. Park, and S.-W. Lee, "Gesture spotting and recognition for human-robot interaction," *IEEE Transactions on Robotics*, vol. 23, no. 2, pp. 256–270, 2007.
- [6] M. Burke and J. Lasenby, "Pantomimic gestures for human-robot interaction," *IEEE Transactions on Robotics*, vol. 31, no. 5, pp. 1225–1237, 2015.
- [7] D. Kim, J. Lee, H.-S. Yoon, J. Kim, and J. Sohn, "Vision-based arm gesture recognition for a long-range human-robot interaction," *The Journal of Supercomputing*, vol. 65, no. 1, pp. 336–352, 2013.
- [8] D. Xu, X. Wu, Y.-L. Chen, and Y. Xu, "Online dynamic gesture recognition for human robot interaction," *Journal of Intelligent & Robotic Systems*, vol. 77, no. 3–4, pp. 583–596, 2015.
- [9] M. Sigalas, H. Baltzakis, and P. Trahanias, "Gesture recognition based on arm tracking for human-robot interaction," in *Proc. IEEE/RSJ Int. Conf. Intelligent Robots and Syst. (IROS)*. IEEE, 2010, pp. 5424–5429.
- [10] G. Cicirelli, C. Attolico, C. Guaragnella, and T. D'Orazio, "A Kinect-based gesture recognition approach for a natural human robot interface," *Int. J. of Advanced Robotic Systems*, vol. 12, no. 22, 2015.
- [11] S. Mitra and T. Acharya, "Gesture recognition: A survey," *IEEE Tran. Syst. Man Cybern. C, pl. Rev.*, vol. 37, no. 3, pp. 311–324, 2007.
- [12] V. Villani, L. Sabattini, G. Riggio, A. Levratti, C. Secchi, and C. Fantuzzi, "Interacting with a mobile robot with a natural infrastructure-less interface," in *Proc. IFAC 20th World Congress Int. Federation Autom. Control IFAC*, ser. IFAC-PapersOnLine, vol. 50, no. 1. Elsevier, 2017, pp. 12753–12758.
- [13] V. Villani, B. Capelli, C. Secchi, C. Fantuzzi, and L. Sabattini, "Humans interacting with multi-robot systems: a natural affect-based approach," *Autonomous Robots*, vol. 44, no. 3, pp. 601–616, 2020.
- [14] S. S. Rautaray and A. Agrawal, "Vision based hand gesture recognition for human computer interaction: a survey," *Artificial intelligence review*, vol. 43, no. 1, pp. 1–54, 2015.
- [15] J. Suarez and R. R. Murphy, "Hand gesture recognition with depth images: A review," in *21st IEEE Int. Symp. Robot and Human Interactive Communication (ROMAN)*. IEEE, 2012, pp. 411–417.
- [16] M. Chandarana, A. Trujillo, K. Shimada, and B. D. Allen, "A natural interaction interface for uavs using intuitive gesture recognition," in *Advances in Human Factors in Robots and Unmanned Systems*. Springer, 2017, pp. 387–398.
- [17] D. Roggen, A. Calatroni, M. Rossi, T. Holleczeck, K. Förster, G. Tröster, P. Lukowicz, D. Bannach, G. Pirkl, A. Ferscha, et al., "Collecting complex activity datasets in highly rich networked sensor environments," in *7th Int. Conf. Networked Sensing Systems (INSS)*. IEEE, 2010, pp. 233–240.
- [18] H. Junker, O. Amft, P. Lukowicz, and G. Tröster, "Gesture spotting with body-worn inertial sensors to detect user activities," *Pattern Recognition*, vol. 41, pp. 2010–2024, 2008.
- [19] N. C. Krishnan and D. J. Cook, "Activity recognition on streaming sensor data," *Pervasive and mobile computing*, vol. 10, pp. 138–154, 2014.
- [20] H. Han and S. W. Yoon, "Gyroscope-based continuous human hand gesture recognition for multi-modal wearable input device for human machine interaction," *Sensors*, vol. 19, no. 11, p. 2562, 2019.
- [21] M. Kim, J. Cho, S. Lee, and Y. Jung, "Imu sensor-based hand gesture recognition for human-machine interfaces," *Sensors*, vol. 19, no. 18, p. 3827, 2019.
- [22] A. Carfi, C. Motolese, B. Bruno, and F. Mastrogiovanni, "Online human gesture recognition using recurrent neural networks and wearable sensors," in *27th IEEE Int. Symp. Robot and Human Interactive Communication (RO-MAN)*. IEEE, 2018, pp. 188–195.
- [23] M. E. Benalcázar, A. G. Jaramillo, A. Zea, A. Páez, V. H. Andaluz, et al., "Hand gesture recognition using machine learning and the myo armband," in *2017 25th European Signal Processing Conference (EUSIPCO)*. IEEE, 2017, pp. 1040–1044.
- [24] M. F. Wahid, R. Tafreshi, M. Al-Sowaidi, and R. Langari, "Subject-independent hand gesture recognition using normalization and machine learning algorithms," *Journal of computational science*, vol. 27, pp. 69–76, 2018.
- [25] Z. Zhang, K. Yang, J. Qian, and L. Zhang, "Real-time surface emg pattern recognition for hand gestures based on an artificial neural network," *Sensors*, vol. 19, no. 14, p. 3170, 2019.
- [26] S. Jiang, B. Lv, W. Guo, C. Zhang, H. Wang, X. Sheng, and P. B. Shull, "Feasibility of wrist-worn, real-time hand, and surface gesture recognition via semg and imu sensing," *IEEE Transactions on Industrial Informatics*, vol. 14, no. 8, pp. 3376–3385, 2017.
- [27] M. Georgi, C. Amma, and T. Schultz, "Recognizing hand and finger gestures with IMU based motion and EMG based muscle activity sensing," in *Int. Joint Conf. Biomedical Engineering Systems and Technologies - Volume 4*, ser. BIOSTEC 2015, 2015, pp. 99–108.
- [28] P. Trigueiros, F. Ribeiro, and L. P. Reis, "A comparison of machine learning algorithms applied to hand gesture recognition," in *7th Iberian Conf. information systems and technologies (CISTI)*. IEEE, 2012.
- [29] H. Sagha, S. T. Digumarti, J. d. R. Millán, R. Chavarriaga, A. Calatroni, D. Roggen, and G. Tröster, "Benchmarking classification techniques using the Opportunity human activity dataset," in *2011 IEEE Int. Conf. Systems, Man, and Cybernetics*. IEEE, 2011, pp. 36–40.
- [30] J. P. Wachs, M. Kölsch, H. Stern, and Y. Edan, "Vision-based hand-gesture applications," *Communications of the ACM*, vol. 54, no. 2, pp. 60–71, 2011.
- [31] M. Hearst, S. Dumais, E. Osuna, J. Platt, and B. Scholkopf, "Support vector machines," *IEEE Intelligent Systems and their applications*, vol. 13, no. 4, pp. 18–28, 1998.
- [32] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, Oct 2001.
- [33] J. Quinlan, "Induction of decision trees," *Machine learning*, vol. 1, no. 1, pp. 81–106, 1986.
- [34] C. M. Bishop, *Pattern recognition and machine learning*. Springer, 2006.
- [35] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al., "Scikit-learn: Machine learning in python," *Journal of machine learning research*, vol. 12, no. Oct, pp. 2825–2830, 2011.
- [36] M. Abadi, A. Agarwal, and P. B. et al., "TensorFlow: Large-scale machine learning on heterogeneous systems," 2015, software available from tensorflow.org. [Online]. Available: <https://www.tensorflow.org/>
- [37] M. Quigley, K. Conley, B. Gerkey, J. Faust, T. Foote, J. Leibs, R. Wheeler, and A. Y. Ng, "ROS: an open-source robot operating system," in *Proc. ICRA Workshop Open Source Software*, vol. 3, no. 3.2, 2009, p. 5.
- [38] J. Kim, S. Mastnik, and E. André, "Emg-based hand gesture recognition for realtime biosignal interfacing," in *Proceedings of the 13th international conference on Intelligent user interfaces*, 2008, pp. 30–39.
- [39] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [40] D. M. Powers, "Evaluation: from precision, recall and f-measure to ROC, informedness, markedness and correlation," *Journal of Machine Learning Technologies*, vol. 2, no. 11, pp. 37–63, 2011.
- [41] A. Tharwat, "Classification assessment methods," *Applied Computing and Informatics*, 2018.