



# Actuator-level motion and contact episode learning and classification using adaptive resonance theory

Vinzenz Bargsten<sup>1</sup> · Frank Kirchner<sup>1,2</sup>

Received: 19 October 2022 / Accepted: 2 August 2023  
© The Author(s) 2023

## Abstract

Several methods exist to detect and distinguish collisions of robotic systems with their environment, since this information is a critical dependency of many tasks. These methods are prevalently based on thresholds in combination with filters, models, or offline trained machine learning models. To improve the adaptation and thereby enable a more autonomous operation of robots in new environments, this work evaluates the applicability of an incremental learning approach. The method addresses online learning and recognition of motion and contact episodes of robotic systems from proprioceptive sensor data using machine learning. The objective is to learn new category templates representing previously encountered situations of the actuators and improve them based on newly gathered similar data. This is achieved using an artificial neural network based on adaptive resonance theory (ART). The input samples from the robot's actuator measurements are preprocessed into frequency spectra. This enables the ART neural network to learn incrementally recurring episodic patterns from these preprocessed data. An evaluation based on preliminary experimental data from a grasping motion of a humanoid robot's arm encountering contacts is presented and suggests that this is a promising approach.

**Keywords** Adaptive resonance theory · Robotics · Contact detection · Incremental learning

## 1 Introduction

The distinction of contact situations is an important dependency for many robotic tasks, be it the detection of an unexpected collision during motion of a manipulator arm, the determination of a contact during a grasping operation, or the assessment of a foothold of walking robots. Robotic systems that are used in unstructured environments, or that work closely with humans, or are operated autonomously, etc., have to be able to deal with a variety of contact situations. Unlike structured industrial processes, these situations cannot be fully predicted in advance. Detection based on static models and initial assumptions may not work in all situations. A continuous learning strategy enabling robotic

systems to recognize and distinguish episodes of motion and contact during operation and over its lifetime is therefore desirable. Rather than relying on dedicated sensors that translate into well-defined physical quantities, these learning strategies should seek to leverage the large amount of proprioceptive sensory signals wherever possible—even if the connection to the detection is not directly visible.

Traditional solutions for considering contacts in robotic manipulator applications involve reactive strategies such as various types of impedance control or end-effector force control, as well as hybrid approaches. While this is valid for preprogrammed tasks with known specific reference forces and motions, a purely reactive control strategy does not distinguish the source of the external force. For this purpose, dedicated collision detection methods are applied. In the simplest and most common case, two situations are distinguished. In the case of a robot manipulator arm, this can be the normal operation and an unexpected collision, which often leads to an emergency stop. Classically, this kind of collision detection is usually implemented as either model-based method relying on an accurate dynamic (physical) model or by using dedicated sensors measuring directly the forces. Manually tuned filters and thresholds are then used

---

✉ Vinzenz Bargsten  
vinzenz.bargsten@dfki.de  
Frank Kirchner  
frank.kirchner@dfki.de

<sup>1</sup> Robotics Innovation Center, DFKI (German Research Center for Artificial Intelligence), Bremen, Germany

<sup>2</sup> Workgroup Robotics, University of Bremen, Bremen, Germany

for the detection. An early work is described in [1], where a disturbance observer is used to determine a signal for the threshold-based detection of collisions. An extensive survey with a focus on dynamic model-based methods for collision detection or monitoring can be found in [2]. While the methods considered in that work are generally applicable to rigid body systems, it is clear that the filtering, thresholding, and error estimation need to be closely tailored towards a specific robotic manipulator, involving mostly manual work. An approach described in [3] tries to address this issue using proprioceptive motor side measurement data. It introduces dynamic thresholds and an adaptation of the filter coefficients to reduce the experimental effort for parameter tuning and parametric model identification.

With the emerging of the field of human-robot collaboration, where humans may actively interact with a robotic system also through physical contact during cooperative tasks, the class of intended contacts get into focus. To achieve this, it is necessary to divide the collisions into more than two categories. Cho et al. [4] further divide the rate of change in the torque measurement to distinguish collisions from external forces due to interaction. For less reliance on the absolute value of the measurement, detection methods have been applied in the frequency domain instead of directly monitoring the measurement signal in time domain. In particular, Kouris et al. [5] show a threshold-based method to distinguish between hand guidance and collisions based on the existence of high-frequency components in the frequency spectrum of the external forces. The external forces are measured directly with a dedicated force sensor, using a manipulator arm in a laboratory environment. In [6], the same author extended the method proposing a “distinction method based on the derivative of the spectral norm of the external force/torque signal in the frequency domain” in order to further reduce the delay between impact and detection. In addition to the threshold-based methods, data-driven approaches such as [7] exist, where a fuzzy modeling of the expected force/torque is used. However, this is more to improve the prediction of the dynamic model for the expected joint torques than to learn and classify the different contact situations during runtime. Another data-driven approach is described in [8], where tactile sensation and reflex is combined with a classical model-based control method through a concept of robot pain sensation. Deviation from a nominal pain level is then used to influence a joint-level impedance controller.

From these samples of related work, it can be seen that the improvement of contact classification is still an ongoing effort. Limitations of model-based and thresholding approaches based on assumptions during the setup of a system have been identified in the literature. Extensive work has been done to improve preprocessing and filtering. If the aspect of adaptation is considered, the adaptation usually refers to the thresholds or the parameters of a parametric

model. One, sometimes two categories of contact situations are regarded and the detection is mostly based on the instantaneous measurement signals and derivatives. The release of the contact, i.e., the end of a contact situation, is largely not addressed.

However, the conditions can change when a robot system is operating autonomously, or needs to quickly adapt to handling of similar product variants without reprogramming in industrial production. The dependence on once tuned thresholds and models limits the ability of robotic systems to adapt to these new conditions. The motivation of this work is therefore to improve the ability of a robotic system to adapt and to let it learn incrementally from its own measurement data during its operation. In particular, the aim of this work is to achieve a more fine-grained and automated discrimination of the episodes arising in repetitive tasks in terms of motion and contact. It should enable a robotic system to learn and recognize typical episodes itself, while performing moving and manipulating actions and use data from what almost any robotic system has—the actuators.

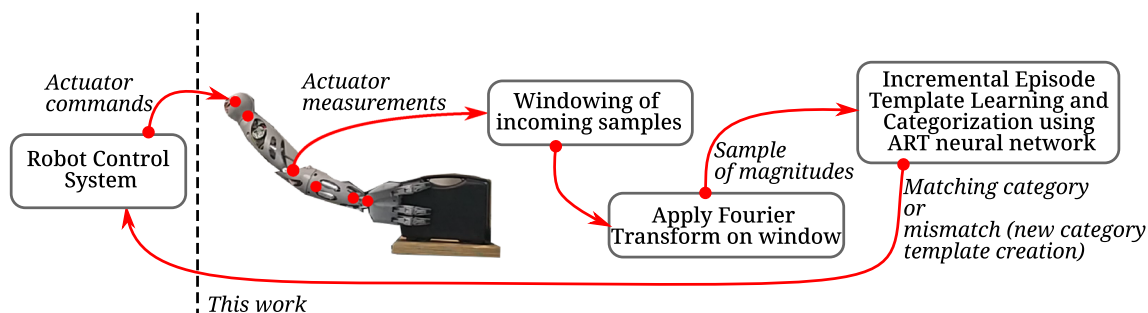
The operating principle of incremental learning allows an adaptation to new situations and also offers the possibility of a finer classification of contact and motion episodes. The novelty and contribution of this work lie in the combination of an incremental learning approach and the use of low-level actuator measurements of a robotic system in the frequency domain as a source of information to build an online classification system.

In a preprocessing step, the time-domain actuator sensory data are transformed into frequency domain. To enable the system to incrementally learn from these data, a neural network based on adaptive resonance theory (ART) is employed. These ART-based neural networks inherently support incremental learning of new data without catastrophic forgetting. In this combination, recurring episodic patterns in the data, such as resulting from the motion of a manipulator arm, can be used to classify different motion and contact situations. An overview of the procedure is shown in Fig. 1. This method can serve as a basis for a higher level behavioral architecture, which could adapt its controllers, or to use a visual system or similar for closer inspection upon creation of new categories in order to rate the new situation.

Both parts, the preprocessing of the data and the classification approach, are described in Sect. 2. Section 3 then presents an experimental evaluation, followed by the discussion and conclusions in Sect. 4.

## 2 Proposed contact learning and classification method

Robot motion in particular, but also human motion [9], often contains recurring episodes, i.e., motions (or applied forces),



**Fig. 1** Overview of the data flow of the proposed method. The intended procedure is as follows: The robotic system executes a repetitive task, initially cautiously and being monitored manually or by additional systems. For some cycles of successful executions, the proprioceptive data are continuously preprocessed into the frequency domain and fed into ART network for an initial training. The ART network then learns category templates for the recurring episodic patterns in the data. As the

which are executed repeatedly in a similar way. The aim of preprocessing is to extract and encode the corresponding data patterns from the raw sensor data, such as the motor current measurements. There are different options for the inclusion of the signal history in the classification. These can be time-delayed samples, derivatives, or the implicit inclusion by filtering and transformation. In this work, the transformation to frequency domain has been selected, because frequency-domain analysis for collision detection (e.g., in [6]) has proven to be successful in the thresholding approaches discussed earlier. This way, the time duration a sample covers is implicitly determined through a windowing function and the following Fourier transformation. However, instead of using dedicated force sensors for the classification task, proprioceptive data coming from the actuators of a robotic system are used. This preprocessing part is briefly described in the next Section. Section 2.2 then presents the ART-based neural network approach applied to classify the resulting frequency-domain data.

## 2.1 Frequency-domain proprioceptive data

To transform the raw data samples into frequency domain, the method short time Fourier series (STFT) [10, 11] is used. It analyzes the time-domain signal by calculating the discrete Fourier transform (DFT) for a window moving over the original signal with time. To determine the DFT of the windowed signal  $X(m, \omega)$  at time  $m$  and frequency  $\omega$  from the time-domain signal  $x[n]$  sampled at time  $n$ , the following transform is executed:

$$X(m, \omega) = \sum_{n=-\infty}^{\infty} x[n]w[n-m]e^{-j\omega n}, \quad (1)$$

patterns repeat, the creation of new category templates stops. Only existing templates are matched and updated with new data. Upon encounter of patterns not matching any previously seen episodes, a high-level control system has to decide on an action such as to trigger a closer inspection. The new patterns can be learnt as a new category template by the ART network for future reference or discarded. The focus of this work is the incremental learning of actuator-level data

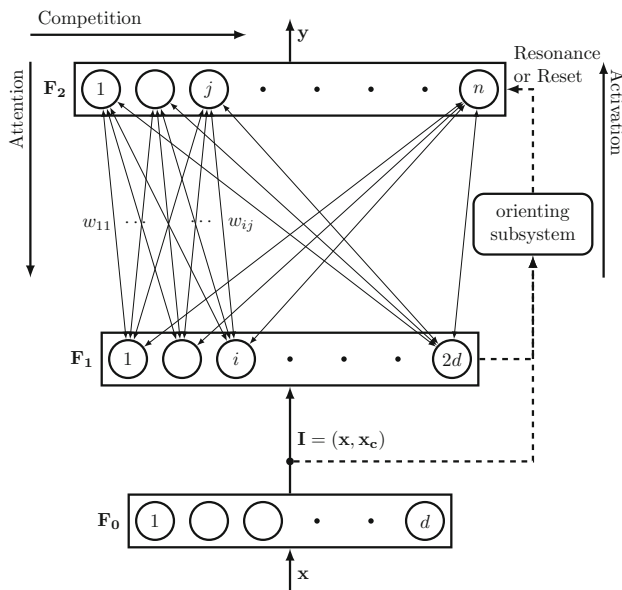
where  $w[n-m]$  is the sliding window function and  $e^{-j\omega n}$  is a modulation shifting the spectrum in order to select the frequency  $\omega$ .

The effects for this work are as follows: (1) it removes the hard dependence on the absolute values of the raw data (of the motor current measurements, for instance), (2) by applying the window function, the signal history is implicitly included in each transformed sample, and (3) the data are compressed (depending on the window size and overlapping). The implementation provided by the *Python SciPy* package, `scipy.signal.stft`, is used. On the downside, with STFT a number of parameters is introduced into the problem. Particularly, the window size, the type of window, and the number of overlapping samples have to be selected.

## 2.2 ART-based classification

Classification is one of the classical application fields of machine learning methods as they can significantly reduce the manual effort of modeling or the development of specific rule-based or (traditional) statistical classification systems. However, many of the currently popular artificial neural networks (ANNs) are based on an offline or batch training with prior definition of the number of classes. The incremental integration of new information is challenging as old information must be stably retained and on the other hand the network must have a certain *plasticity* to incorporate the new information. This trade-off is also known as *plasticity-stability-dilemma* [12]. While strategies have been discussed for the retraining of these kind of neural networks, solutions to this problem are particularly challenging due to the fundamental network properties and also due to the nature of the often used training based on back-propagation [12–16].

ART is a theory of cognitive and neural sciences and was introduced by Carpenter in Grossberg [17–19]. ART-based artificial neural networks are self-organizing and address



**Fig. 2** Basic structure of a FuzzyART network with complement coding of the input samples

the plasticity-stability dilemma as they inherently allow the online learning, in particular to update category templates or to create additional category templates to represent new information. While early implementations were limited to unsupervised training of binary data, various variants and extensions were developed to handle real valued input data and supervised training. An extensive survey is given in [20]. One of the most widely used and extended variants is *FuzzyART* [21], which uses fuzzy set theory in the metric for the similarity and resonance computation. This results in a rectangular-like shape (for two-dimensional data) of the category representation and is also referred to as hyper-rectangular shape. In particular, as shown in Fig. 2, a FuzzyART network is composed of an input layer  $\mathbf{F}_0$ , a feature representation layer  $\mathbf{F}_1$ , a category representation layer  $\mathbf{F}_2$ , and an orienting subsystem. The operation is as follows: an input sample's features are linearly scaled to the range  $[0,1]$ ; to reduce proliferation of the category templates (due to weight erosion), the sample  $\mathbf{x}$  is then complement-coded in the  $\mathbf{F}_0$  layer, giving the encoded feature vector  $\mathbf{I} = (\mathbf{x}, \mathbf{1} - \mathbf{x})^T$ ; then the encoded input sample  $\mathbf{I}$  is presented by the  $\mathbf{F}_1$  layer to the  $\mathbf{F}_2$  layer. The nodes of the  $\mathbf{F}_2$  layer then compete for the best value according to the activation function,

$$T_j = \frac{\|\mathbf{I} \wedge \mathbf{w}_j\|}{\alpha + \|\mathbf{w}_j\|} \quad (2)$$

where  $\|\cdot\|$  refers to the  $L_1$  norm and  $\mathbf{w}_j$  are the weights of the category template—composed of the stored weights  $w_{ij}$ , which in this ART variant are combined bottom-up and top-down weights. The operator  $\wedge$  is the fuzzy set AND operator,

i.e., an intersection, and is defined as the element-wise minimum,

$$\mathbf{a} \wedge \mathbf{b} \equiv (\min(a_0, b_0), \dots, \min(a_i, b_i), \dots)^T \quad (3)$$

The winner of the competition is activated by the orienting subsystem and has to overcome the resonance criterion

$$M_j \geq \rho, \quad (4)$$

in order to be successfully selected, where  $0 \leq \rho \leq 1$  is the vigilance parameter controlling the granularity of the categories. The match value  $M_j$  is computed by the match function,

$$M_j = \frac{\|\mathbf{I} \wedge \mathbf{w}_j\|}{\|\mathbf{I}\|} \quad (5)$$

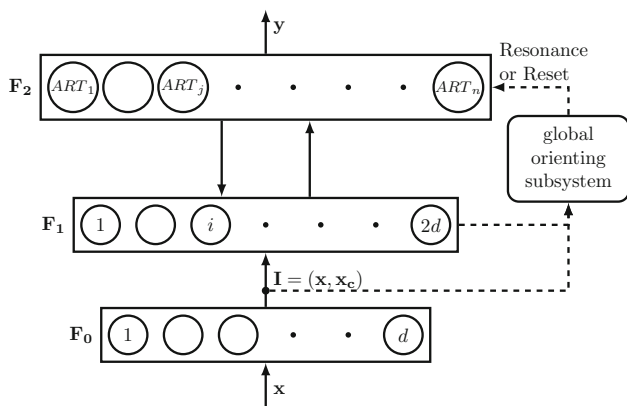
In this case, the category weights  $\mathbf{w}_j^{\text{old}}$  are updated using the input sample  $\mathbf{I}$  according to

$$\mathbf{w}_j^{\text{new}} = (1 - \beta)\mathbf{w}_j^{\text{old}} + \beta (\mathbf{I} \wedge \mathbf{w}_j^{\text{old}}), \quad (6)$$

where  $0 \leq \beta \leq 1$  denotes the learning rate.

Otherwise, if the resonance criterion was not satisfied, the winning node is reset and the category with the next highest activation is tested for resonance. If no resonance occurred at all, a new category is created based on the data from the input sample.

Due to the complement coding, an advantage is the match region shrinks as the category grows—resulting in a stabilizing evolution of the category size. A disadvantage of the (original) *FuzzyART* is the representation of the template data in hyper-rectangles. This leads to a limited classification performance, degrading with the ability to cover the actual data clusters' shapes with the hyper-rectangular FuzzyART category templates. Thus, many ART variants evaluated different metrics, resulting in different shapes of the category templates and match regions. Recently, da Silva et al. introduced the variant *DistributedDualVigilanceFuzzyART* (DDVFA) [22, 23]. In particular, DDVFA is an extension of the ART network by an additional layer of nested ART networks as shown in Fig. 3. Each node in the global ART's category representation field  $\mathbf{F}_2$  does not represent the data directly, but instead is a local ART. Each local ART contains a group of category templates of the actual input data. This allowed to model more complex shapes by a group of FuzzyART category templates. Thereby, the classification performance can significantly improve, especially if the data of a class do not fit a single hyper-rectangular (or, respectively, for other variants) shape and compete with current offline methods for two- and three-dimensional datasets. Instead of one, two vigilance parameters have to be set: the global vigilance



**Fig. 3** The structure of the *DistributedDualVigilanceFuzzyART* [22], which embeds local FuzzyART networks as  $F_2$  nodes, thereby representing a cluster by a group of hyper-rectangles (comparable to a pixel graphic)

parameter  $\rho_{glob}$ , controlling the resonance among the nested ARTs, and the local vigilance parameter  $\rho_{loc}$ , controlling the resonance inside a nested ART. This results in the constraint  $\rho_{loc} \geq \rho_{glob}$  for a useful operation.

The method was re-implemented for the contact learning and classification task envisioned in this work using the Python programming language. To scale the individual samples, multiple options have been implemented. In order to improve the balance between frequency ranges of rather high magnitudes and ranges with typically lower magnitudes, the absolute of the STFT magnitudes can be scaled logarithmically before linearization to the [0,1] range. A second option deals with the linearization itself. By default, the linearization uses the global minimum and maximum (or estimated bounds) of the features for all samples. This requires a priori knowledge of the feature value bounds. Alternatively, it is also possible to linearize each sample by using the minimum and maximum values of the sample itself. In this case, the features express a relation of the magnitudes in each window. The advantage is that no global minimum and maximum needs to be estimated and fixed beforehand. The disadvantage is that samples with similar relative magnitudes, but with different norm, become indistinguishable.

To summarize, in this work a combination of the two methods, STFT and ART, is implemented and applied on the actuator measurements obtained from a humanoid robotic arm and body actuators. The next section presents the respective experimental results.

### 3 Experimental evaluation

#### System Setup

The experimental system is the robotic system RH5v2 [24], composed of an upper body having 3 degrees of freedom

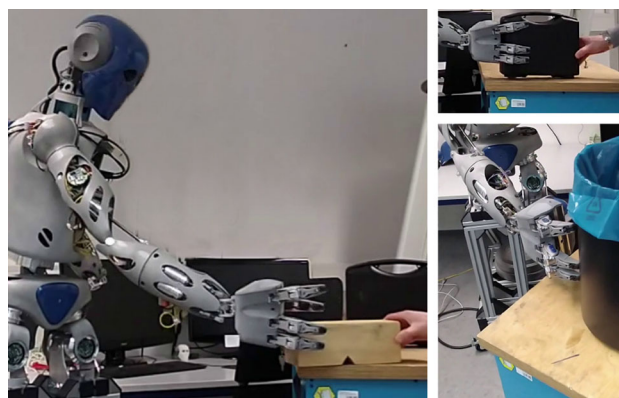
(DOF) and two arms with 7 DOF per arm. In contrast to typical commercial robotic manipulators, in this system access to the lowest level of actuator control is available and the actuator measurements can be obtained in a sufficiently high sampling rate. Two test scenarios have been investigated experimentally. In the first one, collisions are introduced into an otherwise free arm motion. This is used for a first analysis of the data and general applicability of the method. In the second one, the *normal* motion includes a contact episode with significant forces as a plug is inserted into a socket.

### 3.1 Collision contact

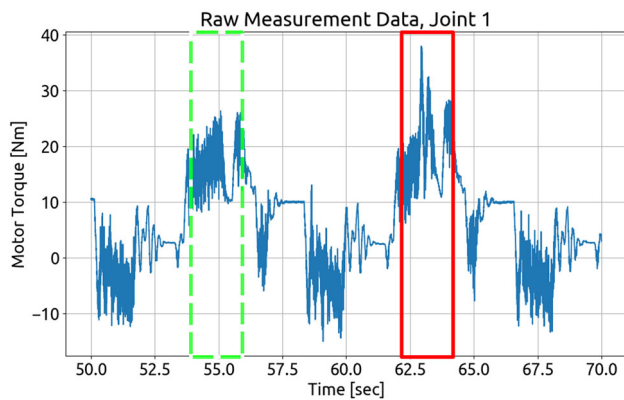
#### 3.1.1 Experimental Setup

The experimental data have been generated as follows: The robot’s upper body and right arm perform point-to-point motions, mimicking an approach motion to grasp an object from a table top. An external operator then introduces a number of collisions by placing objects on the table, see Fig. 4. The joint positions, velocities, and motor currents are measured. In addition, an end-effector force–torque sensor is used as ground truth to later on determine the impact times.

From the 7 arm and 3 body actuators of the kinematic chain, the arm actuators of Joint 12 (Shoulder) and Joint 13 (Shoulder) have been selected for analysis in this experiment. In particular, the joint torques are calculated from the motor current measurements of these actuators, which are then used for the learning and classification task. The data have been sampled with 500 Hz, thus resulting in a spectrum up to 250Hz. The corresponding actuators are BLDC motors (rotational joint) with a gear ratio of 100:1. Exemplary raw data with a non-contact movement and a movement disturbed by pushing an obstacle are shown in Fig. 5. The short time Fourier transform is then applied on these raw data measurements, resulting in frequency magnitudes over time. These



**Fig. 4** Data collection experiments: Robot gripper approaches and pushes different types of objects such as a toolbox, a wood piece, and a bin



**Fig. 5** Exemplary raw data (zoomed in) of one actuator while multiple repetitions (period time ca. 8s) of an arm motion are executed; highlighted in red rectangle (solid line): collision at ca. 63s, highlighted in green rectangle (dashed line) for comparison: the same part of the motion in a repetition without collision (colour figure online)

**Table 1** Parameters used for the preprocessing and classification

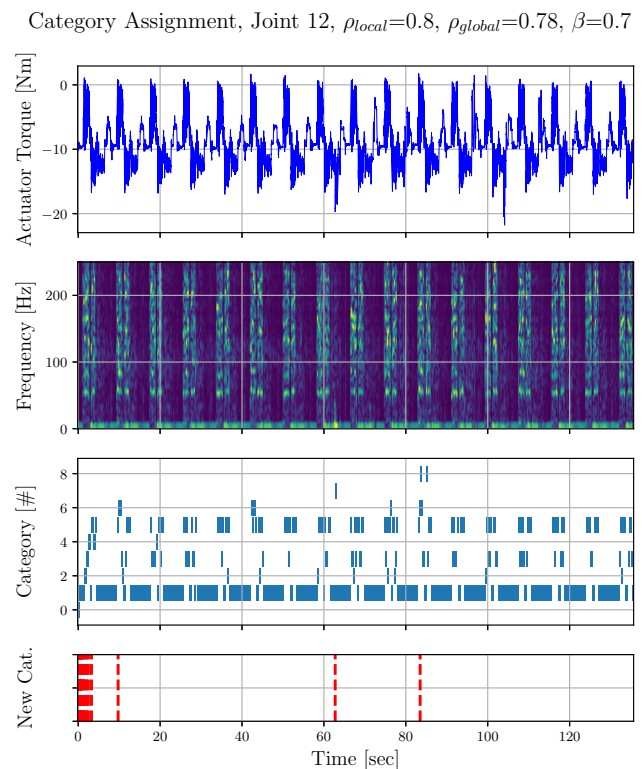
| Parameter name                 | Symbol               | Value                          |
|--------------------------------|----------------------|--------------------------------|
| Sampling frequency             | $f_s$                | 500Hz                          |
| Number of features             |                      | 513                            |
| Number of samples <sup>a</sup> |                      | 67760                          |
| STFT segment length            | $n_{\text{perseg}}$  | 100...500                      |
| STFT window type               | $\text{window}$      | blackmanharris                 |
| Overlap                        | $n_{\text{overlap}}$ | 0; $n_{\text{perseg}}/2$       |
| Local vigilance                | $\rho_{\text{loc}}$  | 0.8...0.97                     |
| Global vigilance               | $\rho_{\text{glob}}$ | 0.8...0.97 $\rho_{\text{loc}}$ |
| Learning rate                  | $\beta$              | 0.7                            |

<sup>a</sup>Time-domain samples of motor current per actuator

spectra are then scaled and used column-wise, i.e., over time, as feature vector for the input to the ART network. For this, the magnitudes are linearly scaled to the range [0, 1]. The values of the free parameters (see Table 1) have been determined through a grid search.

### 3.1.2 Results

The resulting category assignment in relation to the input frequency spectra for the selected measurements is shown in Figs. 6, 7. Figure 8 shows the  $L_2$  norm of the raw wrist force measurements to serve as ground truth. Depending on the choice of the vigilance parameters, this results in ca. 1–6 baselines of recurring category assignments resembling the episodic nature of the motion. Additionally, a number of outliers, i.e., new categories with few assignments, are visible. Table 2 gives an overview of the contact events and the outcome of the classification in terms of the detection of contact events as new category. In combination, the classifications of the spectrum from the three actuators result in a successful



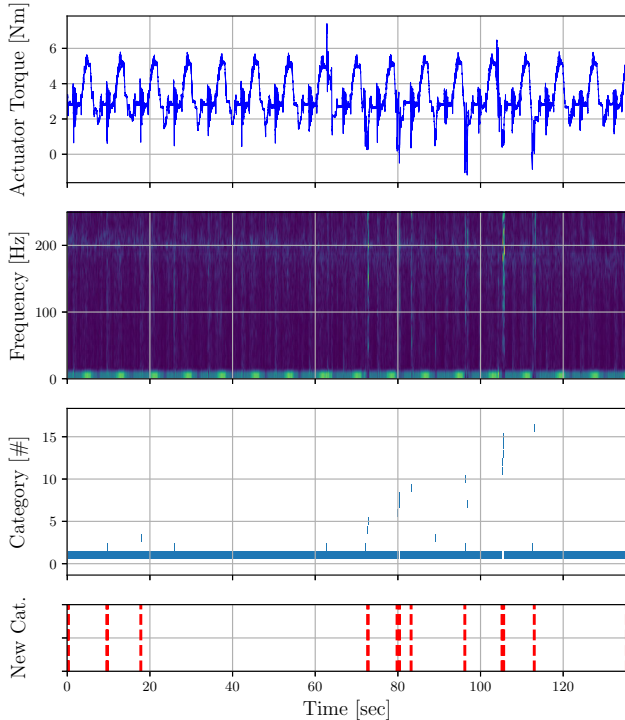
**Fig. 6** Experimental results from the collision experiment. Multiple repetitions are shown. A complete cycle of the grasping motion has an approximate period time of 8 s, resulting in the visible periodic patterns. From top to bottom: Raw actuator torque measurement, STFT frequency spectrum of the windowed raw data, timeline of the assigned categories, and indication of a mismatch leading to creation of a new category. For this joint, after some collision-free cycles, the ART neural network indicates a mismatch at times ca.  $t = 63$  s,  $t = 82$  s, meaning the data are different to what it has seen up to that point

detection. For the actuator of Joint 12 and 13, the full frequency range has been used. The next section discusses the influence of vigilance parameters, windowing parameters, and the frequency range.

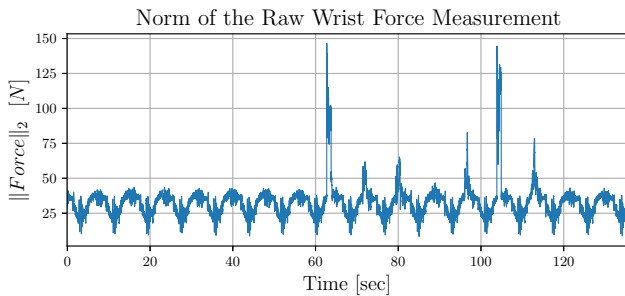
#### Analysis of parameter influence

The effect of the local and global vigilance parameters has been analyzed for the experimental data of this work. The resulting number of categories created with increasing values of the vigilance parameters is shown in Fig. 9. In addition, the number of samples assigned to the individual categories, also depending on the vigilance parameters, has been investigated. The outcome is shown in Fig. 10. To obtain the histogram representation, the following 7 bin sizes (intervals) have been set to determine the distribution of categories vs. the assigned samples:  $[0, 2^1 - 1]$ ,  $[2^1, 2^2 - 1]$ ,  $\dots$ ,  $[2^5, 2^6 - 1]$ ,  $[\geq 2^6]$ . The vigilance parameters have to be chosen such that the desired granularity of the classification is achieved. The extremes thus are assignment of all samples to one category vs. one category per sample. In addition, if the relevant information is mostly present in part of the feature vector, the

Category Assignment, Joint 13,  $\rho_{local}=0.93, \rho_{global}=0.90, \beta=0.7$



**Fig. 7** Experimental results from the collision experiment. Multiple repetitions are shown. A complete cycle of the grasping motion has an approximate period time of 8 s, resulting in the visible periodic patterns. From top to bottom: Raw actuator torque measurement, STFT frequency spectrum of the windowed raw data, timeline of the assigned categories, and indication of a mismatch leading to creation of a new category. For this joint, after some collision-free cycles, the ART neural network indicates a mismatch at times ca.  $t = 72$  s,  $t = 80$  s,  $t = 97$  s,  $t = 104$  s,  $t = 113$  s, meaning the data are different to what it has seen up to that point



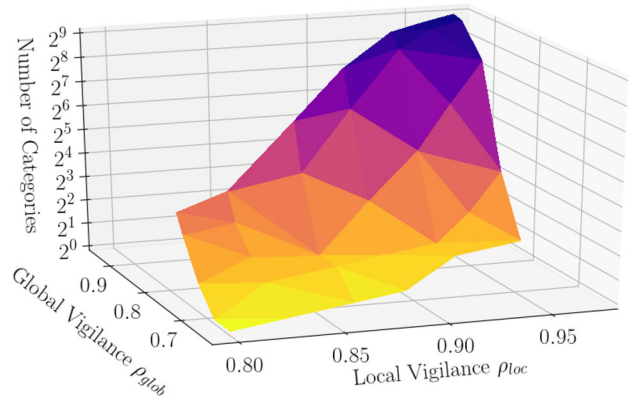
**Fig. 8** Experimental results from the collision experiment. This plot shows the  $L_2$  norm of the raw wrist force measurements to serve as ground truth. Multiple repetitions are shown, and the contact events can be seen in the spikes of the signal

classification performance may be improved when limiting to the respective range of the feature vector. In this work, this is for example selection of lower or upper quarter of the frequency range only for the classification. Moreover, the window size can be varied to focus on a longer or shorter time period in each sample.

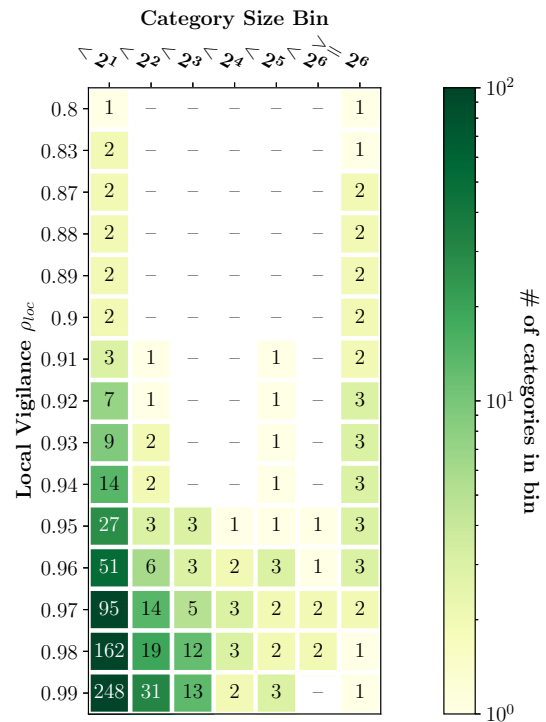
**Table 2** Detectability of the contacts

| Contact event      | Intensity | Joint |    |
|--------------------|-----------|-------|----|
|                    |           | 12    | 13 |
| at ca. $t = 63$ s  | Strong    | ✓     | ✗  |
| at ca. $t = 72$ s  | Light     | ✗     | ✓  |
| at ca. $t = 80$ s  | Light     | ✓     | ✓  |
| at ca. $t = 97$ s  | Medium    | ✗     | ✓  |
| at ca. $t = 104$ s | Strong    | ✗     | ✓  |
| at ca. $t = 113$ s | Medium    | ✗     | ✓  |

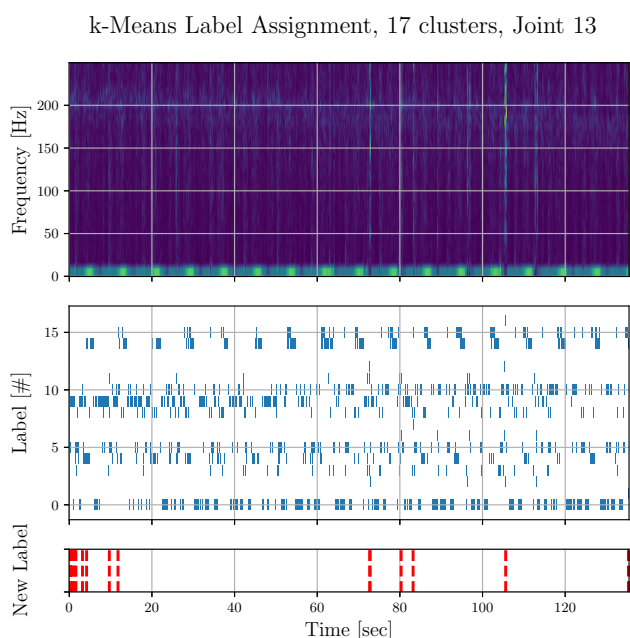
<sup>a</sup>Detected as distinct event/new category



**Fig. 9** Number of categories with increasing vigilance parameters



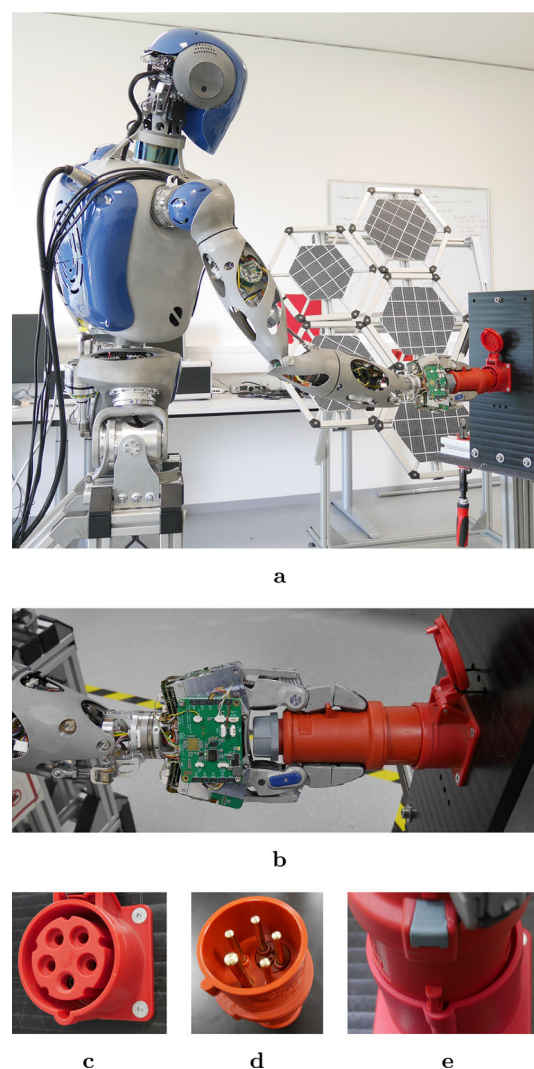
**Fig. 10** The figure shows a distribution of the number of categories by how many samples they would get assigned for different  $\rho_{loc}$ . Reading example: Setting  $\rho_{loc} = 0.93$  results in 9 categories representing only 1 sample, 2 categories representing 2 to 3 samples, and so on



**Fig. 11** For comparison, the resulting label assignment from the k-means algorithm based on the same data used in Fig. 7. Multiple repetitions are shown. A complete cycle of the grasping motion has an approximate period time of 8 s, resulting in the visible periodic patterns. From top to bottom: STFT frequency magnitudes of the windowed raw data, timeline of the assigned labels using the k-means algorithm (only non-incremental batch processing), and indication of first occurrence of a label. The number of clusters was set to 17 and has to be set beforehand

### Comparison with Batch-Clustering

Multiple classical unsupervised clustering algorithms available in the `scikit-learn` Python package have been tested for comparison. The previously obtained frequency-domain data shown in Fig. 7 are used as dataset to apply the offline-clustering algorithms. Prior to application of the clustering algorithms, the data have been scaled as suggested in the manual if necessary. Initially, two density-based unsupervised clustering algorithms, *density-based spatial clustering of applications with noise (DBSCAN)* [25] and *ordering points to identify the clustering structure (OPTICS)* [26], have been tested as they do not require the specification of the number of clusters beforehand. However, for the used dataset no meaningful clusters were found using the standard settings, and varying the  $\epsilon$  parameter. Therefore, the *k-means* algorithm [27] has been tested. It requires the number of clusters to be defined beforehand. The number of clusters has been increased until the resulting labels shown in Fig. 11 were obtained. As we can see, some of the contact events are detected and the results are—depending on the parameters—of a similar nature as the classifications obtained from the ART neural network. In comparison to ART, the disadvantage is the requirement to specify the number of clusters and the offline processing of the complete dataset, which prevents the



**Fig. 12** The experimental setup used for the plug insertion and removal task. **a** shows the overall setup, **b** robot gripper with grasped plug, **c** the socket, **d** the plug, and **e** a misalignment

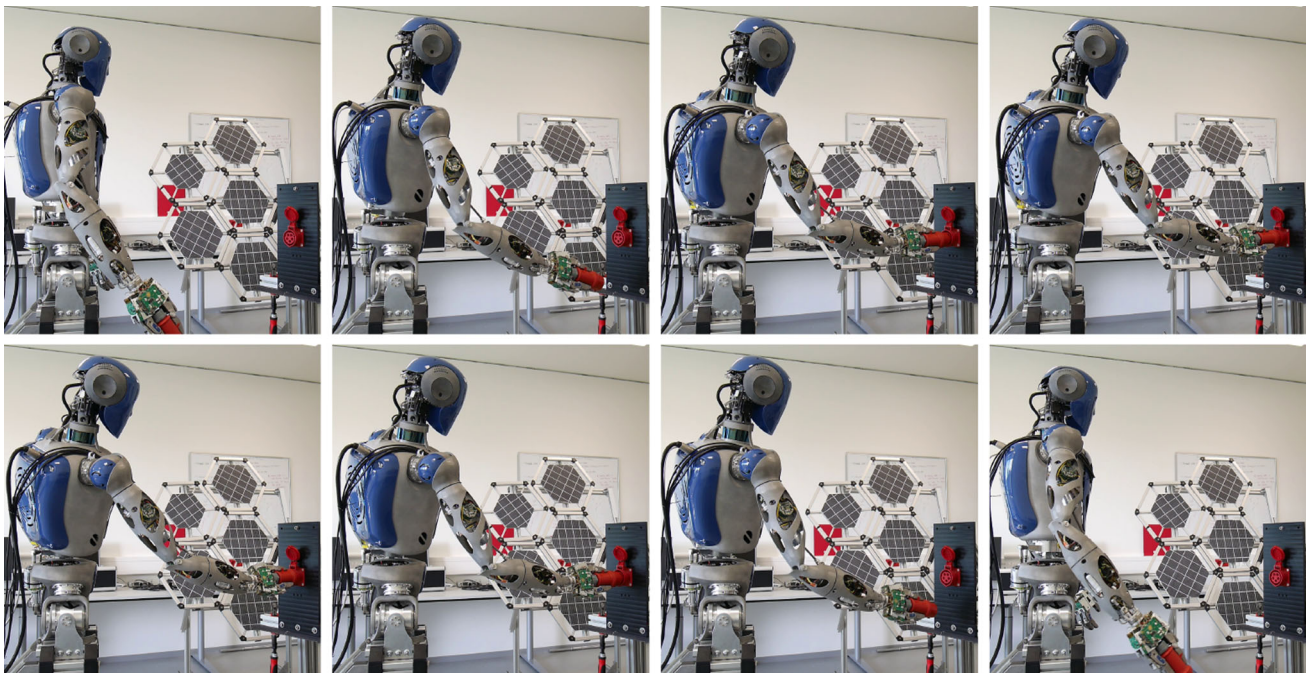
operation in a continuous way. For an extensive comparison of the employed DDVFA ART variant, other ART variants, and classical clustering algorithms, the reader is referred to [23].

## 3.2 Plug insertion and removal

### 3.2.1 Experimental Setup

In a second experimental setup, a *peg-in-hole*-type task was performed by the robot. The setup is shown in Fig. 12. In particular, the robot is used to insert a 5 pin power plug (3 L + N + PE, 6 h according to standard IEC 60309) into a socket and remove it again. The alignment tolerance is approximately  $\pm 1$  mm. The arm and body motion has been taught in by waypoints; object detection has not been used as the focus is





**Fig. 13** Motion of the robot body and right arm as performed for the plug insertion and removal experiment. The motion is determined by previously taught-in points

on the classification of the actuator-level data. Snapshots of the motion are shown in Fig. 13. It requires significant forces to insert the plug into the socket. So, the difference to the previous experiment is that the repetitive task includes these contact episodes within its normal operation. This makes it considerably harder to distinguish the contacts of the normal operation from collisions, which is the main challenge of this experiment.

Due to the forces needed to insert and remove the plug, the plug will slightly change the alignment within the gripper during the plugging and unplugging phase. Thus, a misalignment of the plug's and socket's circumference and notch is to be expected. Obviously, if the misalignment becomes too large, it prevents a successful insertion of the plug. The proposed method is employed to distinguish this situation from the previously executed successful insertions.

Again the joint positions, velocities, and motor currents are measured. In addition, an end-effector force-torque sensor is used as ground truth.

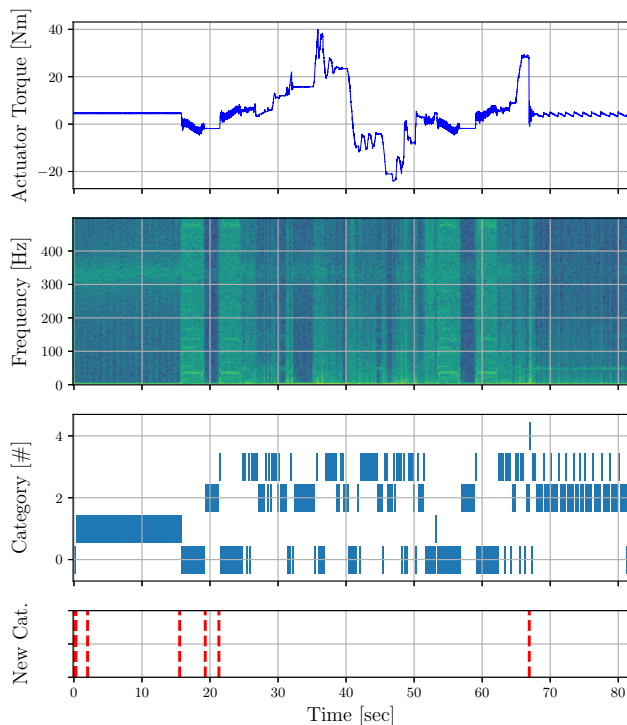
In a training phase the system is taught via programming by demonstration to plug and unplug. During a successful run, the system learns from the new data to improve its category templates. For the analysis in this work, the data have been acquired and stored. However, the operation of the ART network is fully incremental. An online operating component has been developed and can be used once the hyperparameters are chosen. It receives a continuous stream and after 1–2 initial cycles of plug insertion and removal (in this experi-

ment) the category template creation settles. The presented method only covers the contact/motion episode classification. It is therefore necessary to use it in a larger framework to react on mismatches after the initial learning phase.

### 3.2.2 Results

The motor current measurements of a body actuator have been selected for the learning and classification task. The data have been sampled with 1000 Hz, thus resulting in a spectrum up to 500 Hz. The resulting actuator measurement, frequency spectrum, as well as a timeline of category assignment and creation is exemplarily shown for a single actuator measurement in Fig. 14 and Fig. 15. In the data shown, one cycle of insertion and removal is successfully performed. In the second cycle, a slight misalignment has been introduced. The plug then does not slide into the socket but first gets stuck on the circumference of the socket. Only as the forces increase, the plug suddenly releases and slips into the socket. The motion is then stopped.

To show the effect of different scaling, Fig. 14 uses a logarithmic scaling of the frequency-domain data, whereas in Fig. 15 only the linear scaling is applied. In addition, Fig. 16 shows the  $L_2$  norm of the raw wrist force measurements to serve as ground truth. The parameters have been obtained through a grid search.

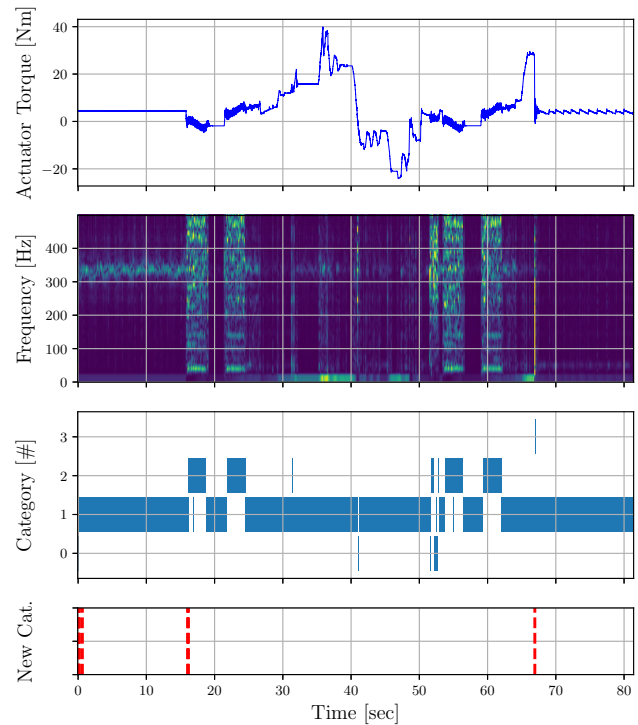
Category Assignment, Joint 18,  $\rho_{local}=0.8$ ,  $\rho_{global}=0.78$ ,  $\beta=0.7$ 

**Fig. 14** Experimental results from the plug insertion, one cycle of successful insertion and one partial cycle with misalignment of plug and socket at ca.  $t = 67$  s. From top to bottom: raw actuator torque measurement, STFT frequency spectrum of the windowed raw data (logarithmic scaling), timeline of the assigned categories, and indication of a mismatch leading to creation of a new category. At ca.  $t = 67$  s, the ART neural network indicates a mismatch, meaning the data are different to what it has seen up to that point

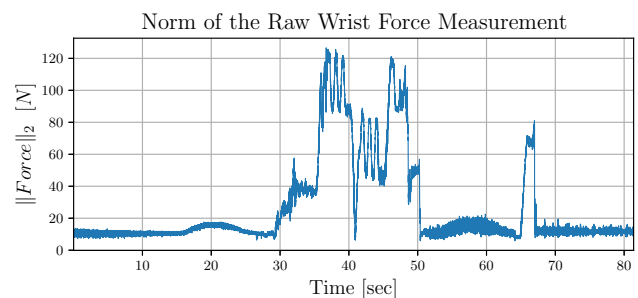
From this experiment, we can conclude that the method also works in a repetitive task involving contacts to detect unseen episodes.

## 4 Discussion and conclusion

The aim of this work was to achieve a more fine-grained and automated discrimination of the episodes arising in repetitive tasks in terms of motion and contact. The proposed method builds up on an incremental learning neural network and the preprocessing low-level actuator into frequency-domain data to allow an online operation. The presented experimental results show that this method can be used to recognize and distinguish recurring episodes in the data. In comparison to [6], the proposed method works without additional end-effector sensors directly on the actuator measurements. Moreover, using the incremental learning approach, an arbitrary number of categories can be learnt during an initial phase. As the experiments have shown, also episodes with contacts can be distinguished. A sub-objective was to analyze

Category Assignment, Joint 18,  $\rho_{local}=0.8$ ,  $\rho_{global}=0.76$ ,  $\beta=0.7$ 

**Fig. 15** Experimental results from the plug insertion, one cycle of successful insertion and one partial cycle with misalignment of plug and socket at ca.  $t = 67$  s. From top to bottom: raw actuator torque measurement, STFT frequency spectrum of the windowed raw data, timeline of the assigned categories, and indication of a mismatch leading to creation of a new category. At ca.  $t = 67$  s, the ART neural network indicates a mismatch, meaning the data are different to what it has seen up to that point. For comparison with Fig. 14, the samples here have been scaled linearly only before classification



**Fig. 16** Experimental results from the plug insertion, one cycle of successful insertion and one partial cycle with misalignment of plug and socket at ca.  $t = 67$  s. This plot shows the  $L_2$  norm of the raw wrist force measurements to serve as ground truth

the way in which the parameters influence the classification for the selected input data. Instead of thresholds for the input data, the vigilance parameter effectively allows to control the resolution. It determines how fine-grained the classification will be or how similar the input data must be to be assigned to the same category. A particular interest is whether the

number of categories stabilizes with the continuous input of episodic movements, while at the same time new situations can be recognized. The experimental results show that values for the vigilance can be found through grid search, such that the number of categories stabilizes and does not grow indefinitely with more data captured. This also means previous episodes have been mapped to already known category templates. Minor numerical differences in the data do not lead to continuous creation of new categories, which would otherwise increase the number of categories rapidly with new data samples. Thus, the system is robust with respect to small differences resulting from sources such as measurement noise. After an initial phase, a pattern of category assignments matching approximately the pattern of the frequency spectrum became visible. If disturbances result in a different distribution of frequency amplitudes, new category templates were learnt subsequently. Some of the new category assignments are visible as outliers, since they have been placed where no new disturbance was introduced. In some of the cases, by closer inspection of the frequency spectrum corresponding differences can be found in the data. This leads to the assumption that the assignment of new categories at these points was correct. A potential effect to consider is the indirect influence of the contact on the actuators due the resulting oscillations of the structure.

It is clear that depending on the configuration of the kinematic chain, a single actuator measurement is limited in its contribution to a classification. It could be shown that the classification of the measurements of several actuators can be combined to improve the recognition capabilities, though. Due to the use of joint actuator measurements, the applicability of the method is limited to those configurations where the forces are not transmitted exclusively through the structure of the robot. In principle, ART-based classification can work with arbitrary sensor data, possibly coming from different sensor modalities. If necessary, the use of additional sensors is therefore a possible option to overcome this limitation.

While ART allows incremental learning, some prior knowledge of the input data is required. In particular, these are the minimum and maximum values of the features in order to scale the samples into the range of  $[0, 1]$  and to allow the complement coding. In this work, this information has been determined from the full experimental dataset. Additionally, due to the nature of incremental learning, initially very different samples may be connected by intermediate samples only at a later state. As a result, for a single data cluster multiple categories may be created. To correct this, compression or merging procedures can be used to reduce the number of categories—during times where the system is not actively operating and has free resources.

For a useful application, the classification needs to be integrated in a larger framework to, on the one hand, trigger actions based and the classification results. On the other hand,

it needs to receive information to map the category templates to a more high-level knowledge, and obtain labels for the categories. The ART network can be used to detect a previously unseen episode, learn its signature, and communicate this event to the high-level robot behavioral control system. This could initiate a fast response of the robotic system, to reduce the speed or force, and a slow response to use a vision system to inspect the particular area more closely. Finally, communicating the result back to the ART network—directly or even after an offline processing at a later point—would allow to mark this category appropriately as critical or as ordinary during the operation.

The ART neural network only has few parameters to tune, with the vigilance parameters having by far the most influence. However, the parameters of the preprocessing are as well crucial. In particular, if the frequency range of the input data is large, changes in small ranges may be ignored if the sample otherwise matches well with the category template. In this case, a contact might be well visible in lower frequencies, but may still be matched to a non-contact category since the higher frequencies are matching well to a non-contact category. Increasing the vigilance parameter value then is an option, but could introduce sensitivity to noise. As an outlook, having a committee of ARTs with different focus on the features, i.e., streaming the data through an additional attentional layer, could mitigate this trade-off. At least it was possible through manual focusing, by selecting a frequency range empirically, to improve the classification for some of the actuators.

ART is a very transparent and explainable ANN. This work combines it with the proprioceptive data, which is anyways available in robotic systems and often goes unused. The results indicate that the method is a promising approach towards more autonomous operation of a system dealing with the classification of contacts and motion patterns. Further experiments and analysis are necessary for a more precise performance quantification and are therefore focus of ongoing and future research.

**Funding** Open Access funding enabled and organized by Projekt DEAL. This work was supported through grants of the German Federal Ministry for Economic Affairs and Energy (BMWi) in the *TransFIT* project <https://robotik.dfki-bremen.de/en/research/projects/transfit/>, grant numbers FKZ 50RA1701, FKZ 50RA1702, and FKZ 50RA1703.

## Declarations

**Conflict of interest** The authors have no competing interests to declare that are relevant to the content of this article.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the

source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## References

- Takakura S, Murakami T, Ohnishi K (1989) An approach to collision detection and recovery motion in industrial robot. In: 15th Annual conference of IEEE industrial electronics society, pp 421–426. <https://doi.org/10.1109/IECON.1989.69669>
- Haddadin S, De Luca A, Albu-Schaffer A (2017) Robot collisions: a survey on detection, isolation, and identification. *IEEE Trans Robot* 33(6):1292–1312. <https://doi.org/10.1109/TRO.2017.2723903>
- Makarov M, Caldas A, Grossard M, Rodriguez-Ayerbe P, Dumur D (2014) Adaptive filtering for robust proprioceptive robot impact detection under model uncertainties. *IEEE/ASME Trans Mechatron* 19(6):1917–1928. <https://doi.org/10.1109/TMECH.2014.2315440>
- Cho C-N, Kim J-H, Kim Y-L, Song J-B, Kyung J-H (2012) Collision detection algorithm to distinguish between intended contact and unexpected collision. *Adv Robot* 26(16):1825–1840. <https://doi.org/10.1080/01691864.2012.685259>
- Kouris A, Dimeas F, Aspragathos N (2016) Contact distinction in human-robot cooperation with admittance control. In: 2016 IEEE international conference on systems, man, and cybernetics (SMC). IEEE, Budapest, Hungary, pp 001951–001956. <https://doi.org/10.1109/SMC.2016.7844525>
- Kouris A, Dimeas F, Aspragathos N (2018) A frequency domain approach for contact type distinction in human-robot collaboration. *IEEE Robot Autom Lett* 3(2):720–727. <https://doi.org/10.1109/LRA.2017.2789249>
- Dimeas F, Avendaño-Valencia LD, Aspragathos N (2015) Human-robot collision detection and identification based on fuzzy and time series modelling. *Robotica* 33(9):1886–1898. <https://doi.org/10.1017/S0263574714001143>
- Kuehn J, Haddadin S (2017) An artificial robot nervous system to teach robots how to feel pain and reflexively react to potentially damaging contacts. *IEEE Robot Autom Lett* 2(1):72–79. <https://doi.org/10.1109/LRA.2016.2536360>
- Gutzeit L (2022) Hierarchical segmentation of human manipulation movements. In: Proceedings of the 26th international conference on pattern recognition. International conference on pattern recognition (ICPR-2022), August 21–25, Montreal, QC, Canada. IEEE. <https://doi.org/10.1109/ICPR56361.2022.9955634>
- Allen JB, Rabiner LR (1977) A unified approach to short-time Fourier analysis and synthesis. *Proc IEEE* 65(11):1558–1564. <https://doi.org/10.1109/PROC.1977.10770>
- Griffin D, Jae Lim (1984) Signal estimation from modified short-time Fourier transform. *IEEE Trans Acoust Speech Signal Process* 32(2):236–243. <https://doi.org/10.1109/TASSP.1984.1164317>
- Grossberg S (2020) A path toward explainable AI and autonomous adaptive intelligence: deep learning, adaptive resonance, and models of perception, emotion, and action. *Front Neurobot* 14:36. <https://doi.org/10.3389/fnbot.2020.00036>
- Lesort T, Lomonaco V, Stoian A, Maltoni D, Filliat D, Díaz-Rodríguez N (2020) Continual learning for robotics: definition, framework, learning strategies, opportunities and challenges. *Inf Fusion* 58:52–68. <https://doi.org/10.1016/j.inffus.2019.12.004>
- Lesort T (2020) Continual learning: tackling catastrophic forgetting in deep neural networks with replay processes. [arxiv:2007.00487](https://arxiv.org/abs/2007.00487) [cs]
- Luo Y, Yin L, Bai W, Mao K (2020) An appraisal of incremental learning methods. *Entropy* 22(11):1190. <https://doi.org/10.3390/e22111190>
- Mermillod M, Bugajska A, Bonin P (2013) The stability-plasticity dilemma: investigating the continuum from catastrophic forgetting to age-limited learning effects. *Front Psychol*. <https://doi.org/10.3389/fpsyg.2013.00504>
- Carpenter GA, Grossberg S (1987) A massively parallel architecture for a self-organizing neural pattern recognition machine. *Comput Vis Graph Image Process* 37(1):54–115. [https://doi.org/10.1016/S0734-189X\(87\)80014-2](https://doi.org/10.1016/S0734-189X(87)80014-2)
- Carpenter GA, Grossberg S, Rosen DB (1991) ART 2-A: an adaptive resonance algorithm for rapid category learning and recognition. *Neural Netw* 4(4):493–504. [https://doi.org/10.1016/0893-6080\(91\)90045-7](https://doi.org/10.1016/0893-6080(91)90045-7)
- Grossberg S (2013) Adaptive resonance theory: how a brain learns to consciously attend, learn, and recognize a changing world. *Neural Netw* 37:1–47. <https://doi.org/10.1016/j.neunet.2012.09.017>
- da Silva LEB, Elnabarawy I, Wunsch II DC (2019) A survey of adaptive resonance theory neural network models for engineering applications. [arxiv:1905.11437](https://arxiv.org/abs/1905.11437) [cs, stat]
- Carpenter GA, Grossberg S, Rosen DB (1991) Fuzzy ART: an adaptive resonance algorithm for rapid, stable classification of analog patterns. In: *IJCNN-91-Seattle international joint conference on neural networks*, vol ii, pp 411–416. IEEE, Seattle, WA, USA. <https://doi.org/10.1109/IJCNN.1991.155368>
- Brito da Silva LE, Elnabarawy I, Wunsch DC (2020) Distributed dual vigilance fuzzy adaptive resonance theory learns online, retrieves arbitrarily-shaped clusters, and mitigates order dependence. *Neural Netw* 121:208–228. <https://doi.org/10.1016/j.neunet.2019.08.033>
- Brito da Silva LE, Elnabarawy I, Wunsch DC (2019) Dual vigilance fuzzy adaptive resonance theory. *Neural Netw* 109:1–5. <https://doi.org/10.1016/j.neunet.2018.09.015>
- Boukheddimi M, Kumar S, Peters H, Mronga D, Budhiraja R, Kirchner F (2022) Introducing RH5 manus: a powerful humanoid upper body design for dynamic movements. In: *IEEE international conference on robotics and automation (ICRA-2022)*, Philadelphia, USA. IEEE. <https://doi.org/10.1109/ICRA46639.2022.9811843>
- Schubert E, Sander J, Ester M, Kriegel HP, Xu X (2017) DBSCAN revisited, revisited: why and how you should (still) use DBSCAN. *ACM Trans Database Syst* 42(3):1–21. <https://doi.org/10.1145/3068335>
- Ankerst M, Breunig MM, Kriegel H-P, Sander J (1999) OPTICS: ordering points to identify the clustering structure. *SIGMOD Rec* 28(2):49–60. <https://doi.org/10.1145/304181.304187>
- Lloyd S (1982) Least squares quantization in PCM. *IEEE Trans Inf Theory* 28(2):129–137. <https://doi.org/10.1109/TIT.1982.1056489>

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.