

# Social robotics in therapy of apraxia of speech

José Carlos Castillo<sup>a,\*</sup>, Diego Álvarez-Fernández<sup>a</sup>, Fernando Alonso-Martín<sup>a</sup>,  
Sara Marques-Villarroya<sup>a</sup>, Miguel A. Salichs<sup>a</sup>

<sup>a</sup>*Departamento de Sistemas y Automática, Universidad Carlos III de Madrid,  
28911-Madrid, Spain*

---

## Abstract

Apraxia of speech is a motor speech disorder in which messages from the brain to the mouth are disrupted, resulting in an inability for moving lips or tongue to the right place to pronounce sounds correctly. Current therapies for this condition involve a therapist that in one-on-one sessions conducts the exercises. Our aim is to work in the line of robotic therapies in which a robot is able to perform partially or autonomously a therapy session, endowing a social robot with the ability of assisting therapists in apraxia of speech rehabilitation exercises. Therefore, we integrate computer vision and machine learning techniques to detect the mouth pose of the user and, on top of that, our social robot performs autonomously the different steps of the therapy using multimodal interaction.

*Keywords:* Apraxia of Speech, Social Robotics, Robotic Therapy, Machine Learning, Face Recognition, Human-Robot Interaction

---

## 1. Introduction

Apraxia Of Speech (AOS) is a neurological disorder that causes that messages from the brain to the mouth are disrupted, and the person cannot move his/her lips or tongue to say sounds properly. This condition is caused by a

---

\*Corresponding author. Tel: +34 916246218

*Email addresses:* [jocastil@ing.uc3m.es](mailto:jocastil@ing.uc3m.es) (José Carlos Castillo), [liqueric@gmail.com](mailto:liqueric@gmail.com) (Diego Álvarez-Fernández), [famartin@ing.uc3m.es](mailto:famartin@ing.uc3m.es) (Fernando Alonso-Martín), [smarques@ing.uc3m.es](mailto:smarques@ing.uc3m.es) (Sara Marques-Villarroya), [salichs@ing.uc3m.es](mailto:salichs@ing.uc3m.es) (Miguel A. Salichs)

5 damage in the left hemisphere of the brain generated by strokes, Alzheimer’s  
or brain traumas, among others. The severity of the apraxia depends on the  
nature of the brain damage. AOS is also known as acquired apraxia of speech,  
verbal apraxia, and dyspraxia [1].

The focus of intervention is on improving the planning, sequencing, and co-  
10 ordination of muscle movements for speech production. The muscles of speech  
often need to be “retrained” to produce sounds correctly and sequence sounds  
into words. Exercises are designed to allow the person to repeat sounds over  
and over and to practice correct mouth movements for sounds [2]. Currently,  
there are three different interventions for AOS rehabilitation: (i) Intervention  
15 based on motor control: these exercises consist of producing phonemes and  
sequences of phonemes through accurate, controlled and conscious movements.  
The aim of these therapies is to automate such movements to be subsequently  
performed unwittingly; (ii) Intervention based on augmented systems: these  
methods include several input channels to improve the therapy results. Au-  
20 dio and images are traditionally mixed to help remembering how to pronounce  
difficult and long words; and (iii) Interventions based on melodies: these ther-  
apies are adopted in patients that preserve an auditive comprehension of the  
language. In these cases, the patient has to imitate different melodies proposed  
that remark the stressed syllables of the words, establishing the rhythm of the  
25 melodies [3]. These interventions are frequently planned as intensive, one-on-  
one speech-language therapy sessions for both children and adults. Thus, the  
repetitive exercises and personal attention needed to improve AOS are difficult  
to deliver in group therapy [4].

In recent years, robots are gaining popularity in rehabilitation therapies,  
30 mainly in traumatology, where the robot holds the user’s weight or helps moving  
a determined limb. Robots have proved to be effective in assisting the therapist  
to provide safe and intensive rehabilitation training for the stroke subjects. Nev-  
ertheless, in the general setting of these systems, a therapist is still responsible  
for the non-physical interaction and observation of the patient by maintaining a  
35 supervisory role of the training, while the robot carries out the actual physical

interaction with the patient. In most applications, rehabilitation robots have been mainly employed in lower and upper limbs therapy [5, 6, 7, 8].

Rehabilitation using robotics is generally well tolerated by patients, and has been found to be an effective adjunct to therapy in individuals suffering from motor impairments, especially due to stroke. Therefore, we believe that robotics can be introduced to other rehabilitation areas such as AOS. To the extent of our knowledge, this proposal is innovative as robotic technologies have not been applied to this field so far. We propose following the first kind of intervention presented in this section in which the user repeats exercises to practice mouth movements, in our case, we take inspiration from mouth poses associated to the five vowels in the Spanish tongue. Here, “a” is pronounced like the “a” in the word “father” (/a/); “e” is pronounced like the “a” in the word “date” (/e/), except that it is shorter and crisper; “i” is pronounced like the “ee” in the word “see” (/i/); “o” is pronounced like the “o” in the word “no” (/o/); and “u” is pronounced like the “e” in the word “new” (/u/). Thus, we propose using some of these sounds because their pronunciation imply different poses of the mouth, associated to a range of muscular movements.

We believe that a social robot could help in AOS therapy offering a new and eye-catching way of assisting in the exercises. The robot adds to the therapy some new resources such as a screen to stimulate the patient, offering a visual reinforcement to the exercises. Additionally, the Human-Robot Interaction (HRI) capabilities of a social robot could enhance the traditional therapy, maximizing the human resources whilst keeping a personalized treatment. That is, a therapists could take care of more patients having robots develop parts of the treatments.

We propose using Machine Learning techniques for vowel pose recognition and identification. The input information is collected by an RGB-D device, a *Microsoft Kinect<sup>TM</sup>*, and with this information the system obtains mouth poses which are used in the exercises to guide the users. Interaction is performed through a multimodal system that integrates body expressions, voice interaction as well as a Graphical User Interface (GUI), all of these modalities are developed

to give instructions to the patient as well as encourage him/her during the exercise.

The rest of this manuscript is structured as follows: Section 2 provides the  
70 insights of current therapies for AOS, presents new robotic developments for  
physical therapy and cognitive rehabilitation and analyses some face detection  
and classification techniques related to our approach. Next, Section 3 presents  
the details of our proposal, describing its main phases. Section 4 present the  
experiments conducted to validate our work along with the robotic platform, the  
75 social robot Mini, and the metrics for evaluating the approach. This section also  
present the preliminary results from integrating and testing the AOS exercises  
in the social robot. Finally, Section 5 analyses the main contributions or our  
work and draws the main conclusions.

## 2. Related Work

80 The ability of speech is commonly affected after suffering Alzheimer's, de-  
mentia or a stroke. Traditional speech therapies focus on mitigating this prob-  
lem in case of cognitive impairment, or rehabilitating in case of cerebrovascular  
accidents. The recovery time in these cases is around three years [9] in which  
speech therapy yield positive results in most cases. Apart from this line of ther-  
85 apy, there are others such as music therapy that are usually applied to patients  
with neurological problems, generally elders. In the case of music, the therapy  
consists of patients emitting singing and emitting sounds from given melodies  
in order to improve pitch, variability and intelligibility of speech [10]. Sacks and  
Tomatino [11] demonstrated that music therapy helps reorganizing the brain  
90 function in patients with brain alterations.

Apart from traditional therapies, technology is being incorporated to health  
environments. More precisely, robotics is gaining importance, mainly in the  
fields of physical therapy and cognitive rehabilitation. In these cases, exercises  
are supervised by therapists who are in charge of selecting the tasks to perform  
95 and monitor the procedure [12]. The robot Paro is a good example of the ap-

plication of robots in cognitive therapy. It imitates a baby harp seal and has been used in therapy with elder people with dementia, increasing the willingness of patients to communicate and a steady increase in physical interaction, not only between patients and the robot but among patients as well [13, 14]. Another robotic platform used in cognitive therapy is Babyloid, a baby-like robot  
100 other robotic platform used in cognitive therapy is Babyloid, a baby-like robot designed to be taken care of [15]. This robot is intended for recreational therapy in which the robot becomes a pet instead of an animal. These proposals are mainly intended for interaction with elderly people with moderate cognitive impairments.

105 Other robotic platforms provide a higher degree of interaction in therapies with people mild cognitive problems. This is the case of Eldertoy [16], a robot developed to achieve both entertaining and gerontological capabilities. This robot offers different interaction channels to communicate with users: gestures, voice, touch-screen, and external actuators. Therapy with this robot is conceived through manipulation and display multimedia content. Therefore therapy  
110 specialists are furnished with a tool able to run games by using the sensors integrated in the platform. The robot Mini is another proposal for therapy with elders in early stages of Alzheimer’s or dementia [17, 18]. Mini is a plush-like desktop robot that offers functionalities related to safety, personal assistance, entertainment and stimulation. In this work, we aim to extend the capabilities  
115 of Mini to conduct speech therapy. More details about the robot design and features can be found in Section 3.

Another research area integrated in our work is Computer Vision. The literature offers several approaches for face detection and recognition [19, 20].  
120 Applications range from people recognition [21], surveillance [22] to emotion detection and regulation [23, 24]. Although there are several techniques to retrieve facial features, this problem is still challenging since most of the approaches are highly dependent on the face orientation. In this work we have integrated Stasm, an Active Shape Model-based approach coupled to a Support Vector Machine  
125 (SVM) classifier that retrieves facial features [20]. Out of these features, the

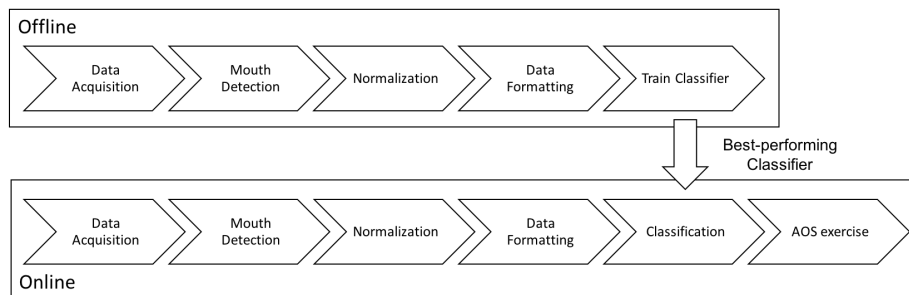


Figure 1: Proposed approach pipeline. The upper path corresponds to the offline analysis for assessing the best classifier. The lower path corresponds to the software running in the robot with the speech therapy exercise.

user’s mouth is represented with 18 3D points, which will be the input for the machine learning algorithm.

Apart from detecting the mouth, recognizing the mouth pose is crucial to have an algorithm that can be integrated into a speech therapy application.

130 Machine Learning have been widely applied in face recognition and recognition of facial expressions [25, 26]. Within the number of techniques, SVM, Adaboost, Linear Discriminant Analysis or, more recently, Deep Learning [27], among others, try to cope with known problems such as different poses, illumination, ages and occlusions that nowadays still pose a challenge. In our work we test several

135 classifiers integrated within Sci-Kit learn [28], an open source machine learning library written in Python language. It provides features classification, regression and clustering algorithms.

### 3. Materials and Methods

This section presents the details of the proposed approach that allows a social

140 robot, equipped with a 3D camera, to conduct an AOS exercise autonomously. Figure 1 shows the main steps of our proposal, which is roughly divided into two operation modes. First, we need to asses the classifier that performs best for our kind of data. In this process we acquire information from users, pre-process it and train a set of classifiers to select the best-performing one. This classifier

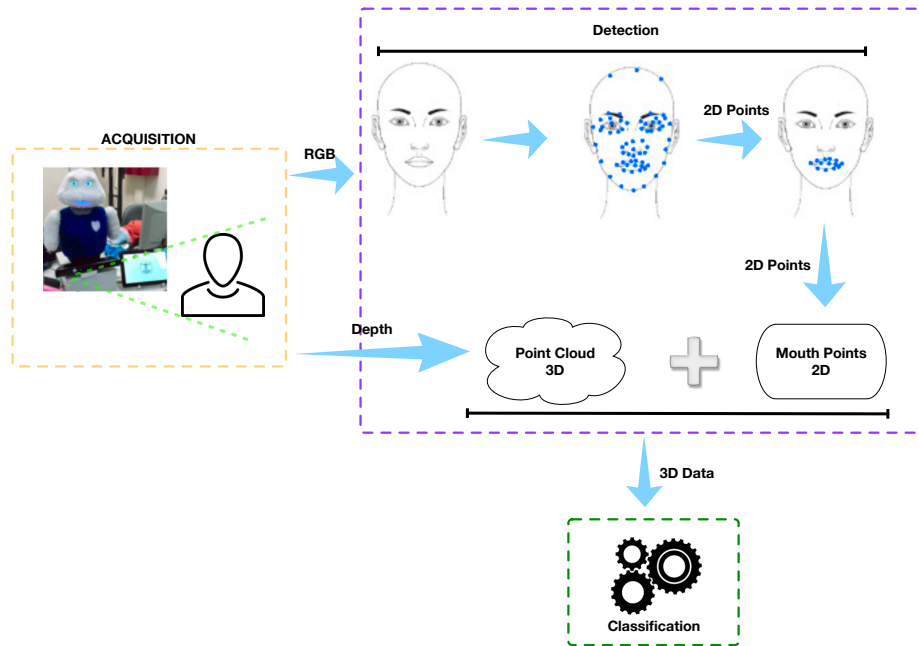


Figure 2: Mouth detection pipeline using Stasm.

145 is used next online, thus integrated in the robotic platform where the speech therapy application uses the mouth pose detected to conduct the exercises. Note that the four first steps are the same in both approaches.

### 3.1. Mouth Detection from RGB-D data

The system described here uses a Microsoft Kinect<sup>TM</sup>, which provides *RGB* images and depth data synchronized both in terms of time and field of view. 150 After information acquisition is performed, the system extracts face features in *2D* using the Open Source library *Stasm* [20]. Then, those features are translated into *3D* points which are finally classified to recognize the mouth pose (see Fig. 2).

155 The mouth detection process is composed by two main steps. Data acquisition is performed through a middleware specifically designed to work with

*RGB-D* devices, *OpenNI*<sup>1</sup>. Two information flows are generated from the Kinect device: an *RGB* image stream and a point cloud containing depth information. Next, the system processes the *RGB* data to identify the mouth within a detected face using *Stasm*. This library characterizes a face with 77 points of which 18 belong to the mouth. These points are next matched to the depth information from the camera and formatted to be used in the next phase, mouth poses classification. A more detailed description of the mouth detection system can be found in a previous work [29].

### 3.2. Machine Learning tools for Mouth Pose Classification

In our approach for mouth pose recognition we aim to test a series of classification techniques integrated within Scikit-Learn [28]. For the classifier selection we take as a starting point a previous work [29] in which mouth detection was evaluated using WEKA [30], a well known data mining tool that allows pre-processing, classification, regression, clustering, association rules, and visualization of data. In this case, we wanted to take the study one step further and integrate the best-performing classifier in our robotic platform (WEKA was not directly integrated within the framework ROS). Thus, we compared the performance of the following classifiers:

- **k-Nearest neighbours** (*k*-NN) is a non-parametric method used for classification and regression in which the input consists of the *k* closest training examples in a feature space [31]. In our problem, the output is a mouth pose where a sample is classified by a majority vote of its neighbours with the object being assigned to the most common class among its *k* nearest neighbours
- **Support Vector Machine** is a supervised learning technique for classification and regression that builds a hyperplane or set of hyperplanes in a high- or infinite-dimensional space [32]. An SVM can perform linear and

---

<sup>1</sup>OpenNI website: <http://openni.ru/>



185 non-linear classification mapping the inputs into high-dimensional feature spaces.

- **C4.5** is an algorithm that generates a decision tree from a set of training data using the concept of information entropy [33]. At each node of the tree, the algorithm chooses the attribute of the data that most effectively splits the set of samples into subsets enriched in one of the classes. 190 The splitting criterion is the normalized information gain (difference in entropy).
- **Random Forest** is an ensemble learning method for classification and regression that construct multiple decision trees at training time and outputs the class that corresponds to the mode of the possible classes (mouth 195 poses). An advantage of Random Forest is that this technique mitigates the overfitting problem caused by traditional decision trees [34].

### 3.3. Assessing the best classifier: Offline analysis

Before addressing the logic of the speech therapy exercise, it was necessary establishing which classifier offered better performance with our input data. 200 This operation mode, depicted in the upper path of Figure 1, starts with a Data Acquisition phase in which the RGB-D device provides colour images and Point Clouds with the 3D representation of the scene. The next step, Mouth Detection works as described in Section 3.1, using Stasm to generate a 3D array of 18 points corresponding to the mouth detected in the input data.

205 Since the head position in the image varies as the user moves, it is important to normalize the data to establish a common frame for reference. Thus, the Normalization step computes the centroid of the mouth, setting it as the origin of coordinates for the 18 points (see Equations 1, 2 and 3). Each one of these points is defined by its  $\langle x, y, z \rangle$  components and therefore the normalization

210 for each of them is calculated with respect to the centroid as shown in Equations 4, 5 and 6.

$$x_{centroid} = \frac{1}{18} \sum_1^{18} x_i \quad (1)$$

$$y_{centroid} = \frac{1}{18} \sum_1^{18} y_i \quad (2)$$

$$z_{centroid} = \frac{1}{18} \sum_1^{18} z_i \quad (3)$$

$$x_{i_{Normalized}} = x_i - x_{centroid} \quad \forall i \in [1, 18] \quad (4)$$

$$y_{i_{Normalized}} = y_i - y_{centroid} \quad \forall i \in [1, 18] \quad (5)$$

$$z_{i_{Normalized}} = z_i - z_{centroid} \quad \forall i \in [1, 18] \quad (6)$$

These normalized points are formatted in tuples for the classifier. Each tuple composed by 54 values plus the class for each pose recorded. After the data is formatted, we trained the classifiers previously described in Section 3.2.

#### 215 3.4. Online execution

The best-performing classifier identified in the previous section is integrated in the online execution mode of our system, described in the lower path of Figure 1. The fourth first steps are common to both offline and online execution as they are intended for data acquisition, detecting the points corresponding to the mouth and normalize them as well as formatting the data for classification. 220 In the online mode, the data formatted is then processed in a classifier which output is used in the AOS exercise to assess the user performance and guide him/her during the session.

When one of the three poses reaches a number of detections the system 225 selects that pose as the current one and interacts with the user, expressing

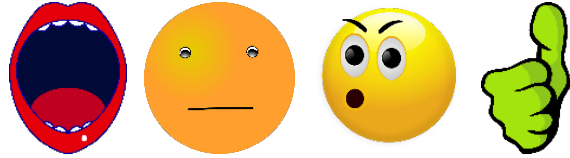


Figure 3: Images to give feedback about the mouth pose and the user’s performance. First, image to indicate how to make the “a” pose. Second, image to indicate how the closed mouth pose should be. Third, image to indicate how to make the “u” pose. Fourth, image to congratulate the user.

congratulations in case the pose detected was the one expected, or correcting the user if another pose is detected. A repertoire or corpus of utterances has been created for congratulate and correct the user (see Table 1) as a complement to the images shown in the tablet (see Fig. 3). Additionally, the robot expresses  
230 gestures with its body to help engaging the user in the exercise. When the user fails to perform the pose, a sad expression is performed whilst otherwise the robot shows a happy expression<sup>2</sup>.

In the current version of the system, the process of detecting a mouth pose and congratulate/correct the user is repeated three times although the system  
235 is flexible enough to change and adapt the exercise logic.

## 4. Results and Discussion

### 4.1. Robotic Platform: Mini

The system developed in this work was integrated in Mini, a desktop social robot designed and built at RoboticsLab research group from Carlos III university of Madrid (see Fig. 4). Originally, this robot was designed to interact with  
240 elder people with mild cognitive impairment [18]. Nevertheless, the capabilities of this platform allows other users and applications such as our current goal.

Mini is equipped with multiple HRI interfaces including Automatic Speech Recognition (ASR) [35], Voice Activity Detection [36], Emotion Detection [37],

---

<sup>2</sup>Instead of defining here what they are about, a demo video has been released in which those gestures are clearly demonstrated. The video link is presented at the end of Section 4.4.

Table 1: Set of utterances to convey messages during the exercise (approximate translation from Spanish).

Situation	Sentence	Details
Congratulating	Very good Keep Going You are doing great	General congratulations utterances
Correcting	Open a little less your mouth Open your mouth less Please, close your mouth a bit	Corrections in case the user is opening the mouth more than expected
	Open your mouth more Open your mouth a bit more Make a bigger mouth	Corrections in case the user is opening the mouth less than expected
	You are almost there Keep trying	Utterance to encourage the user during the exercise
Starting Exercise	Let's try to say 'a' correctly three times	Practicing with 'a' pose
	Let's try to say 'u' correctly three times	Practicing with 'u' pose
	Try to keep your mouth closed for three rounds	Practicing with mouth closed pose
	In 3, 2, 1... Now!	Start signal

245 a Text to Speech (TTS) system, a tablet and an RGB-D device. Moreover, Mini possesses 5 degrees of freedom to allow moving its arms, base and head. The interaction capabilities complete with touch sensors, two uOLED screens for the eyes, RGB LEDs in the cheeks and heart and a VUmeter as mouth to create the illusion of a talking robot. All of these interfaces are integrated in  
250 a Natural Dialogue Management System [38] which enables the robot to carry out natural interactions. Finally, these components are integrated using ROS framework [39].

#### 4.2. Metrics for evaluating the classifiers

Since our classifiers have to solve a multi-class problem, the metric selected  
255 for assessing the best one was the Macro-average F-score. Macro-average means that the metric is independently computed from each class and then take the average is calculated. This metric uses the Precision and Recall for each class

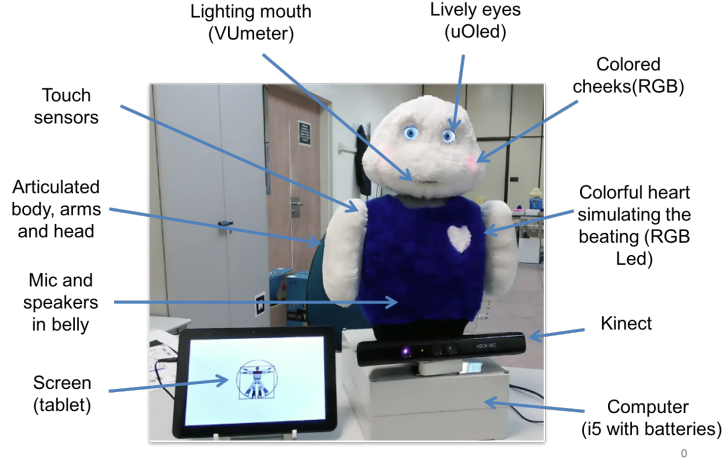


Figure 4: Mini, the social robot involved in the experiments [40]. Apart from its plushy shape, the robot is equipped with a series of sensors and actuators for HRI.

and calculates the mean Precision and Recall of all classes as shown in Equations 7 and 8, respectively:

$$\overline{Precision} = \frac{\sum_{i=1}^N TP_i}{\sum_{i=1}^N TP_i + FP_i} \quad (7)$$

$$\overline{Recall} = \frac{\sum_{i=1}^N TP_i}{\sum_{i=1}^N TP_i + FN_i} \quad (8)$$

260 Where  $N$  is the total number of classes,  $TP_i$  corresponds to the True Positives achieved in class  $i$ ,  $FP_i$  are the False Positives for class  $i$ , and  $FN_i$  are the False Negatives in class  $i$ . Then, the Macro-average F-score, is computed as the harmonic mean of these two values as shown in Equation 9.

$$Macro - average\_F - score = \frac{2 \times \overline{Precision} \times \overline{Recall}}{\overline{Precision} + \overline{Recall}} \quad (9)$$

### 4.3. Experiments

265 In our experiments users sit in front of the robot, at 0,5 meters (see Fig. 5, left). A previous study indicated that at a range of 0,5 meters the detector

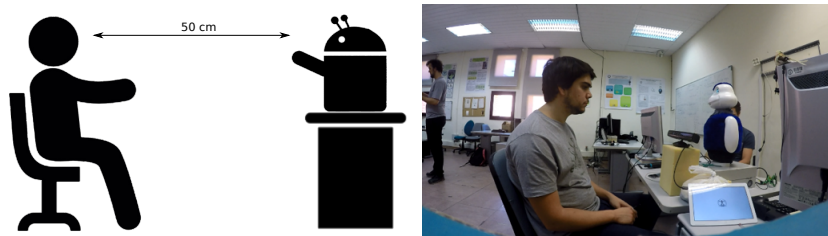


Figure 5: Experimental setup. The user was sitting in front of the robot at 0,5 meters, a natural distance for interaction that ensures clear images of the face.

performance in order to locate the mouth accurately in the face was high, suffering a degradation with the distance that at 2 meters was too poor to achieve reliable detections [29]. Moreover, as shown in Figure 5, right, this distance  
 270 allows a natural interaction with the robot, for instance with its tablet that is usually placed between the robot and the user. Note that in Figure 5 right, the tablet is on the left side of the robot. In this specific case, the tablet was placed there just for illustrative purposes. Also, the Kinect camera changed the usual location (see Fig. 4) to allow a better acquisition of face images.

275 The following sections detail the experiments conducted to assess the performance of our system and its feasibility for speech therapy. Also, the aim of this set of tests was to select the best classifier for our data. We first tested the performance with the most different poses “a” and “u”, as described in experiment 1. For our second experiment, we added a new “neutral” pose that  
 280 corresponded to the mouth closed and carried out experiment 2. 14 users were involved in our experiments and the method for dividing the datasets was 1-fold cross-validation with 60% of the instances for training and 40% for test.

#### 4.3.1. Experiment 1: Training 2 poses

285 This experiment is meant to assess the feasibility of the classifiers described in Section 3.2 to distinguish between two mouth poses. Although this set of poses may seem reduced, they are different enough as to implement a range of mouth movements that could be useful in SOA therapy. Moreover, recognising the mouth is not an easy task, leading to similar representations of different

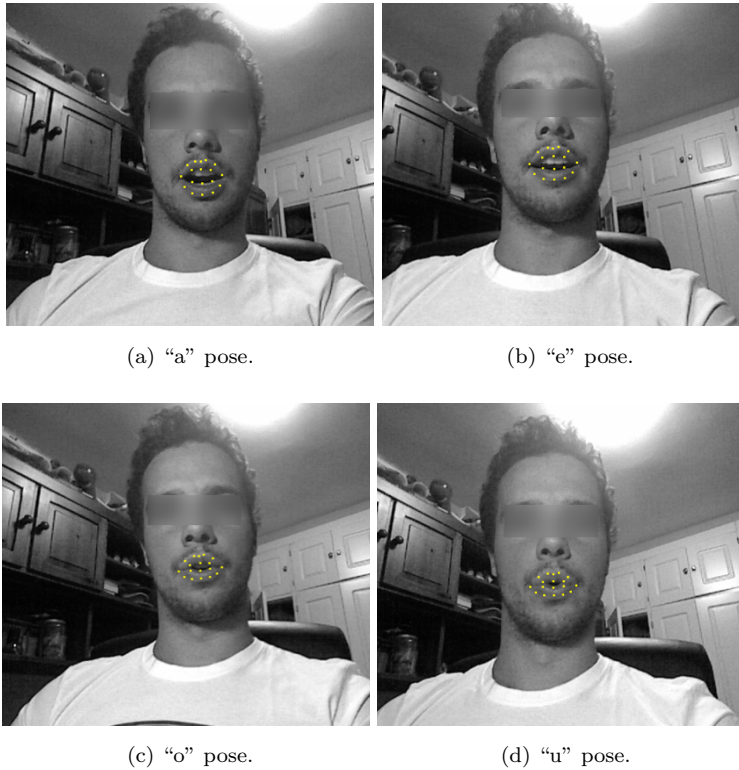


Figure 6: Mouth points as detected by Stasm. We can see in this example how there is low variability between some mouth poses.

poses as shown in Figure 6. The two first images corresponds to the poses  
 290 associated to "a" and "e" and the two last ones corresponds to "o" and "u".

In this test the dataset was composed of 1200 instances for the pose "a"  
 and 1425 for the "u" pose (see Table 2). After the test, two classifiers, C4.5  
 and SVM, showed promising results, with a Macro-average F-Score of 0,82 and  
 0,81, respectively, as shown in Table 3. Additionally, Table 4 offers the confusion  
 295 matrix for the best classifier in this experiment, C4.5, in which we can observe  
 an accuracy that starts to be competitive for our speech therapy application.

Table 2: Datasets summary for experiments 1 and 2

Pose	Instances	
	Experiment 1	Experiment 2
“a”	1200	4623
“u”	1425	5250
Mouth closed	N/A	3624
Total	1625	13497

Table 3: Results for experiments 1 and 2. Macro-average F-Score for the four classifiers tested with two mouth poses.

Classifier	Experiment 1:	Experiment 2:
	Macro-average F-Score	Macro-average F-Score
Random Forest	0.57	0.47
C4.5	<b>0.82</b>	<b>0.95</b>
k-NN	0.54	0.93
SVM	0.81	0.63

#### 4.3.2. Experiment 2: Training 3 poses

In this experiment a new pose was added to the dataset, mouth closed, to complement the cases for the AOS exercise. Therefore, a new class, mouth closed, was added to our dataset with 3624 instances. For the previous classes new instances were added as well, having in total 4623 instances for the “a” pose and 5250 instances for the “u” pose. Finally, the classifiers were retrained with the new data. Results show that C4.5 is again the best classifier, with k-NN offering competitive performance (see Table 3). Therefore, since C4.5 showed the best behaviour in both experiments, this classifier is the one selected for the online execution.

In this experiment, the results improved with respect to the previous one as shown in the confusion matrix for C4.5 classifier (see Table 5). Here the recognition rate for the “a” pose reached 95%, in case of the “u” pose the rate is 93% and finally for mouth closed pose the rate is 99.67%.



Table 4: Experiment 1: Confusion matrix for C4.5 classifier identifying 2 poses.

	Predicted A	Predicted U
Actual A	86%	14%
Actual U	12%	88%

Table 5: Experiment 2: Confusion matrix for C4.5 classifier identifying 3 poses.

	Predicted A	Predicted U	Predicted Mouth Closed
Actual A	95%	4%	1%
Actual U	5%	93%	2%
Actual Mouth Closed	0.03%	0.3%	99.67%

#### 4.4. Integration in the social robot

This section analyses the performance of the detection and classification integrated with the speech therapy application. In this case, we first tried to use the system trained with the dataset described in Experiment 2 (Section 4.3.2), but in this case classifying poses from seven untrained users. Table 6 shows the performance of our pre-trained C4.5 classifier when offered new data. We can see how in some cases as in mouth closed pose the performance drops to the point that the system is not usable.

At this point, we realised that we needed to perform some training with the new users' data, but in this case that training should not be as intensive as in previous experiments. Since the final application is speech therapy, we cannot expect that users will be willing to train the robot for long periods of time. In this case, trained the system with online detections from the users in periods ranging between 5 to 10 minutes for the three poses together. We believe this would not cause boredom or fatigue in the users as it only needs to be performed once per new user. With this new data, the performance of the system improves to levels comparable to experiment 1 (see Table 7) and, although not reaching the scores achieved in experiment 2 with cross-validation, these levels are high enough to ensure a good detection rate.

In the online execution, we experimentally set the score threshold to consider valid detection to 0,35 and a pose was output after six successful recognitions.

Table 6: Confusion matrix for the first test with untrained users.

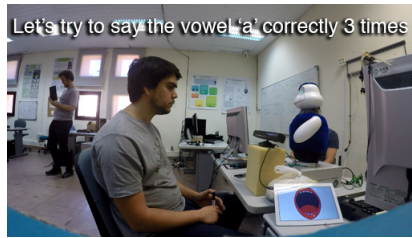
	<b>Predicted A</b>	<b>Predicted U</b>	<b>Predicted Mouth Closed</b>
<b>Actual A</b>	52%	32%	16%
<b>Actual U</b>	36%	40%	24%
<b>Actual Mouth Closed</b>	0%	90%	10%

Table 7: Confusion matrix for the test retraining the classifier with new users data.

	<b>Predicted A</b>	<b>Predicted U</b>	<b>Predicted Mouth Closed</b>
<b>Actual A</b>	85%	12%	3%
<b>Actual U</b>	8%	75%	17%
<b>Actual Mouth Closed</b>	0%	14%	86%

In most cases we noticed that the detection score was close to one, dropping to low values for missdetections. The number of successful recognitions to consider a valid pose directly impacts the execution time of our system since a bigger number would cause a slow response and a smaller number could lead to wrong  
335 detections. Therefore, six valid detections was considered as a good tradeof between time of response and accuracy.

Figure 7 offers an overview of the speech therapy proposal with the different phases where the robot guides the user along the exercise. First, the robot  
340 provides a simple explanation about the exercise (see Fig. 7a) using gestures, voice and the tablet to convey the messages. Next, the exercise starts and the user should start making the desired pose while the system is detecting and classifying the mouth pose (see Fig. 7b). Finally, after three successful detections, the robot congratulates the user (see Fig. 7c). There is also the  
345 possibility that the system does not detect the target pose. In this case, the robot corrects the user, explaining how to achieve the desired pose (see Fig. 7d). Along this exercise the robot uses voice, gestures and the tablet to give instructions and feedback to the user. A video has been uploaded with more details about the execution of the system can be found in <https://youtu.be/XRrIP3BcwCY>.  
350



(a) Starting the exercise. The robot shows a 'happy' expression while giving instructions to the user through voice. In parallel, the tablet shows the pose that the user should imitate.



(b) The robot tells the user to start performing the pose and in parallel captures and analyses the mouth features. The tablet keeps showing the pose that should be imitated.



(c) After three valid detections, the robot congratulates the user through voice, gestures and using the tablet.



(d) If the mouth pose detected does not match the desired one, the robot corrects the user using voice and gestures while the tablet shows the target pose.

Figure 7: Running example of the speech therapy proposal. The robot leads the user through the exercise, encouraging him to keep participating.

Table 8: Numbers summarizing the experimental conditions and results of our proposal

Recognized poses	<b>3 poses</b> (“a”, “u”, mouth closed)
Dataset Instances	<b>1625</b> for experiment 1 and <b>13497</b> for experiment 2
Input features per instance	<b>54 features</b> : 18 mouth points * 3 components (x,y,z)
Users involved in the dataset	<b>14 users</b> for the 2 datasets
Classifiers tested	<b>4 classifiers</b>
Metric for comparison	<b>Macro-average F-score</b>
Best algorithm for classifying poses	<b>C4.5</b> (0,81 and 0,95 of macro-average F-score in the experiments)
Users involved in the real tests	<b>7 users</b>

#### 4.5. Discussion

The results and the experimental conditions of the proposed approach are summarized in Table 8. These results show how our proposal provides high accuracy for mouth poses classification, up to 0,95 in the cross-validation test. It is worth remarking that the experimental phase in this paper is intended as a proof of concept and that we are currently working on testing it with real users. Additionally, we are aware that mouth poses could change when working with people with motor mouth problems and that this fact could affect the performance of the classifiers. In this regard, our plan is to add real users data and retrain the system when deploying it in real scenarios.

Also, the set of mouth poses recognized may seem too small but for AOS therapy purposes their differences were considered enough for a first approach. It is our intention involving experts to evaluate the feasibility of our proposal both in terms of poses recognized and the dynamic of the exercise.

## 5. Conclusions

This manuscript introduced an approach for apraxia of speech therapy using a social robot. The system consists of two main phases: an *offline* one in which we train a set of classifiers after detecting and normalizing the mouth information from users; and an *online* one that runs in our social robot Mini.

370 This operation runs in realtime, integrating the best-performing classifier, and guides the user through an AOS exercise.

The experiments included up to three mouth poses (“a”, “u” and mouth closed) which we consider are enough for a first approach to therapy for their differences regarding the mouth positions. The classifiers trained on our dataset, 375 composed by information from 14 users in our experiments offline, to assess the best one. In these offline experiments, C4.5 was the best classifier for our data (achieving a 0,95 of Macro-average F-Score) and, therefore it was integrated within the final approach. In the online tests with the whole system integrated in the social robot we conducted additional experiments with 7 new users, the 380 first one running the system with untrained data which showed a performance decrease in the mouth poses classification. This motivated the second experiment in which we retrained the classifier adding a small set of samples from the new users. In this case the performance rose again to competitive values.

We believe that the results achieved in our experiments are promising and 385 thus we are intended to proceed with the next stage: testing the AOS exercise with real users and therapists.

To the extent of our knowledge, this proposal is innovative as robotic technologies have not been applied to this field so far.

### **Conflicts of interest**

390 The authors declare that there is no conflict of interest regarding the publication of this paper.

### **Acknowledgements**

The research leading to these results has received funding from the projects: Development of social robots to help seniors with cognitive impairment (ROB- 395 SEN), funded by the Ministerio de Economía y Competitividad; and RoboCity2030-III-CM, funded by Comunidad de Madrid and cofunded by Structural Funds of the EU.

## References

- [1] R. González Victoriano, L. Toledo Rodríguez, Apraxia del habla: Evaluación y tratamiento, *Revista Neuropsicología, Neuropsiquiatría y Neurociencias* 15 (1) (2015) 141–158.
- [2] American Speech-Language-Hearing Association, Apraxia of speech in adults, <https://www.asha.org/public/speech/disorders/ApraxiaAdults/>, Accessed: 2017-11-3.
- [3] A. Ygual-Fernández, J. Cervera-Mérida, Dispraxia verbal: características clínicas y tratamiento logopédico, *Rev Neurol* 40 (2005) 121–26.
- [4] National Institute of Deafness and Other Communication Disorders (NIDCD), U.S. Department of Health & Human Services, Apraxia of speech, [https://www.nidcd.nih.gov/health/apraxia-speech#apraxia\\_05](https://www.nidcd.nih.gov/health/apraxia-speech#apraxia_05), Accessed: 2017-10-31.
- [5] J. F. Veneman, R. Kruidhof, E. E. Hekman, R. Ekkelenkamp, E. H. Van Asseldonk, H. Van Der Kooij, Design and evaluation of the lopes exoskeleton robot for interactive gait rehabilitation, *IEEE Transactions on Neural Systems and Rehabilitation Engineering* 15 (3) (2007) 379–386.
- [6] N. Ho, K. Tong, X. Hu, K. Fung, X. Wei, W. Rong, E. Susanto, An emg-driven exoskeleton hand robotic training device on chronic stroke subjects: task training system for stroke rehabilitation, in: *Rehabilitation Robotics (ICORR)*, 2011 IEEE International Conference on, IEEE, 2011, pp. 1–5.
- [7] H. S. Lo, S. Q. Xie, Exoskeleton robots for upper-limb rehabilitation: State of the art and future prospects, *Medical engineering & physics* 34 (3) (2012) 261–268.
- [8] D. Copaci, A. Flores, F. Rueda, I. Alguacil, D. Blanco, L. Moreno, Wearable elbow exoskeleton actuated with shape memory alloy, in: *Converging Clinical and Engineering Research on Neurorehabilitation II*, Springer International Publishing, 2017, pp. 477–481.

- [9] O. Juncos, A. Rozas, Problemas del lenguaje en la tercera edad. Orientaciones y perspectivas de la logopedia, *Revista galego-portuguesa de psicología e educación: revista de estudios e investigación en psicología y educación* 8 (2002) 387–398.
- 430 [10] N. S. Cohen, The effect of singing instruction on the speech production of neurologically impaired persons, *Journal of Music therapy* 29 (2) (1992) 87–102.
- [11] C. M. T. Oliver Sacks, Music and neurological disorder, *International Journal of Arts Medicine* 1 (1) (1992) 10–12.
- 435 [12] K. Wada, Y. Ikeda, K. Inoue, R. Uehara, Development and preliminary evaluation of a caregiver’s manual for robot therapy using the therapeutic seal robot PARO, in: *RO-MAN, 2010 IEEE, IEEE, 2010*, pp. 533–538.
- [13] K. Takayanagi, T. Kirita, T. Shibata, Comparison of verbal and emotional responses of elderly people with mild/moderate dementia and those with severe dementia in responses to seal robot, PARO, *Frontiers in aging neuroscience* 6.
- 440 [14] W. L. Chang, S. abanovic, L. Huber, Use of seal-like robot PARO in sensory group therapy for older adults with dementia, in: *2013 8th ACM/IEEE International Conference on Human-Robot Interaction (HRI), 2013*, pp. 101–102.
- 445 [15] Y. Furuta, M. Kanoh, T. Shimizu, M. Shimizu, T. Nakamura, Subjective evaluation of use of Babyloid for doll therapy, in: *Fuzzy Systems (FUZZ-IEEE), 2012 IEEE International Conference on, IEEE, 2012*, pp. 1–4.
- [16] L. F. Cossío, J. M. L. Salvador, S. F. Martínez, Robotics for social welfare, *Journal of accessibility and design for all* 2 (1) (2012) 94–113.
- 450 [17] E. Salichs, A. Castro-González, M. Malfaz, M. Salichs, Mini: A social assistive robot for people with mild cognitive impairment, in: *New Friends*

2016. The 2nd international conference on social robots in therapy and education., 2016, pp. 29–30.
- 455 [18] M. A. Salichs, I. P. Encinar, E. Salichs, A. Castro-González, M. Malfaz, Study of scenarios and technical requirements of a social assistive robot for alzheimers disease patients and their caregivers, *International Journal of Social Robotics* 8 (1) (2016) 85–102.
- [19] S. Zafeiriou, C. Zhang, Z. Zhang, A survey on face detection in the wild: 460 past, present and future, *Computer Vision and Image Understanding* 138 (2015) 1–24.
- [20] S. Milborrow, F. Nicolls, Active shape models with SIFT descriptors and MARS, in: *Computer Vision Theory and Applications (VISAPP)*, 2014 International Conference on, Vol. 2, IEEE, 2014, pp. 380–387.
- 465 [21] W. Zhao, R. Chellappa, P. J. Phillips, A. Rosenfeld, Face recognition: A literature survey, *ACM Comput. Surv.* 35 (4) (2003) 399–458.
- [22] Z. Kalal, K. Mikolajczyk, J. Matas, Face-tld: Tracking-learning-detection applied to faces, in: *Image Processing (ICIP)*, 2010 17th IEEE International Conference on, IEEE, 2010, pp. 3789–3792.
- 470 [23] A. Fernández-Caballero, A. Martínez-Rodrigo, J. M. Pastor, J. C. Castillo, E. Lozano-Monador, M. T. López, R. Zangrizz, J. M. Latorre, A. Fernández-Sotos, Smart environment architecture for emotion detection and regulation, *Journal of Biomedical Informatics* 64 (C) (2016) 55–73.
- [24] J. C. Castillo, A. Castro-González, F. Alonso-Martín, A. Fernández- 475 Caballero, M. A. Salichs, Emotion detection and regulation from personal assistant robot in smart environment, in: *Personal Assistants: Emerging Computational Technologies*, Springer, Cham, 2018, pp. 179–195.
- [25] M. S. Bartlett, G. Littlewort, M. Frank, C. Lainscsek, I. Fasel, J. Movellan, 480 Recognizing facial expression: machine learning and application to spontaneous behavior, in: *Computer Vision and Pattern Recognition*, 2005.



- CVPR 2005. IEEE Computer Society Conference on, Vol. 2, IEEE, 2005, pp. 568–573.
- [26] G. Littlewort, M. S. Bartlett, I. Fasel, J. Susskind, J. Movellan, Dynamics of facial expression extracted automatically from video, *Image and Vision Computing* 24 (6) (2006) 615 – 625.
- 485 [27] Y. Sun, Y. Chen, X. Wang, X. Tang, Deep learning face representation by joint identification-verification, in: *Advances in neural information processing systems*, 2014, pp. 1988–1996.
- [28] 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* 12 (Oct) (2011) 2825–2830.
- 490 [29] J. C. Castillo, I. Pérez-Encinar, A. Conti-Morera, A. Castro-González, M. A. Salichs, Vowel recognition from RGB-D facial information, in: *Ambient Intelligence-Software and Applications-7th International Symposium on Ambient Intelligence (ISAmI 2016)*, Springer, 2016, pp. 225–232.
- 495 [30] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, I. H. Witten, The WEKA data mining software: an update, *ACM SIGKDD explorations newsletter* 11 (1) (2009) 10–18.
- [31] N. S. Altman, An introduction to kernel and nearest-neighbor nonparametric regression, *The American Statistician* 46 (3) (1992) 175–185.
- 500 [32] C. Cortes, V. Vapnik, Support-vector networks, *Machine learning* 20 (3) (1995) 273–297.
- [33] J. R. Quinlan, *C4. 5: programs for machine learning*, Elsevier, 1993.
- 505 [34] T. K. Ho, Random decision forests, in: *Document Analysis and Recognition, 1995.*, Proceedings of the Third International Conference on, Vol. 1, IEEE, 1995, pp. 278–282.

- [35] F. Alonso-Martín, M. A. Salichs, Integration of a voice recognition system in a social robot, *Cybernetics and Systems: An International Journal* 42 (4) (2011) 215–245.
- 510
- [36] F. Alonso-Martin, A. Castro-González, J. F. Gorostiza, M. A. Salichs, Multidomain voice activity detection during human-robot interaction, in: *International Conference on Social Robotics*, Springer, 2013, pp. 64–73.
- [37] F. Alonso-Martín, M. Malfaz, J. Sequeira, J. F. Gorostiza, M. A. Salichs, A multimodal emotion detection system during human–robot interaction, *Sensors* 13 (11) (2013) 15549–15581.
- 515
- [38] F. Alonso-Martín, A. Castro-González, J. F. Gorostiza, M. A. Salichs, Augmented robotics dialog system for enhancing human–robot interaction, *Sensors* 15 (7) (2015) 15799–15829.
- [39] M. Quigley, K. Conley, B. Gerkey, J. Faust, T. Foote, J. Leibs, R. Wheeler, A. Y. Ng, ROS: an open-source robot operating system, in: *ICRA workshop on open source software*, Vol. 3, Kobe, 2009, p. 5.
- 520
- [40] J. C. Castillo, D. Cáceres-Domínguez, F. Alonso-Martín, A. Castro-González, M. A. Salichs, Dynamic gesture recognition for social robots, in: *International Conference on Social Robotics*, Springer, 2017, pp. 495–505.
- 525