# ChildBot: Multi-Robot Perception and Interaction with Children

**Niki Efthymiou**[1,2] · **Panagiotis P. Filntisis**[1,2] · **Petros Koutras**[1,2] · **Antigoni Tsiami**[1,2] ·
**Jack Hadfield**[1,2] · **Gerasimos Potamianos**[1,3] · **Petros Maragos**[1,2]

**Abstract** In this paper we present an integrated robotic system capable of participating in and performing a wide range of educational and entertainment tasks, in collaboration with one or more children. The system, called ChildBot, features multimodal perception modules and multiple robotic agents that monitor the interaction environment, and can robustly coordinate complex Child-Robot Interaction use-cases. In order to validate the effectiveness of the system and its integrated modules, we have conducted multiple experiments with a total of 52 children. Our results show improved perception capabilities in comparison to our earlier works that ChildBot was based on. In addition, we have conducted a preliminary user experience study, employing some educational/entertainment tasks, that yields encouraging results

Niki Efthymiou
Email: nefthymiou@central.ntua.gr

Panagiotis P.Filntisis
Email: filby@central.ntua.gr

Petros Koutras
Email: pkoutras@cs.ntua.gr

Antigoni Tsiami
Email: antsiami@cs.ntua.gr

Jack Hadfield
Email: jack103@windowslive.com

Gerasimos Potamianos
Email: gpotam@ieee.org

Petros Maragos
Email: maragos@cs.ntua.gr

[1] Athena Research and Innovation Center, Maroussi, 15125, Greece
[2] School of ECE, National Technical Univ. of Athens, 15773 Athens, Greece
[3] Dept. of ECE, University of Thessaly, Volos 38221, Greece

regarding the technical validity of our system and initial insights on the user experience with it.

## 1 Introduction

Recently, robotic systems entailing Human-Robot Interaction (HRI) at their core have been gaining momentum not only in research, but also in everyday life. Indeed, such robotic platforms have been entering many different aspects of human lives [25], e.g. rehabilitation [47,30], nursing or personal care [24], education [43], and entertainment [57]. In order to achieve a high level of naturalness that resembles human-to-human interaction, robots need to have the ability to perceive and understand the different modalities that people use for communication such as speech or body movements [56,36].

The majority of existing social robotics systems present two major deficiencies: First, they usually incorporate only specific modalities, forcing the users to adapt to the way the system perceives the environment, instead of the opposite. Secondly, they are developed and designed for specific applications and tasks.

A natural interaction involves the creation of smart adaptive integrated robotic systems capable of multitasking, and with a wide range of perceptual and actuation abilities. An HRI system that would be capable of multitasking, encourages users to design multiform interactive applications that can maintain the interest of the participants undiminished. This is especially important when the participants are children, in the context of education and entertainment (edutainment). In addition, systems leveraging multiple percep-
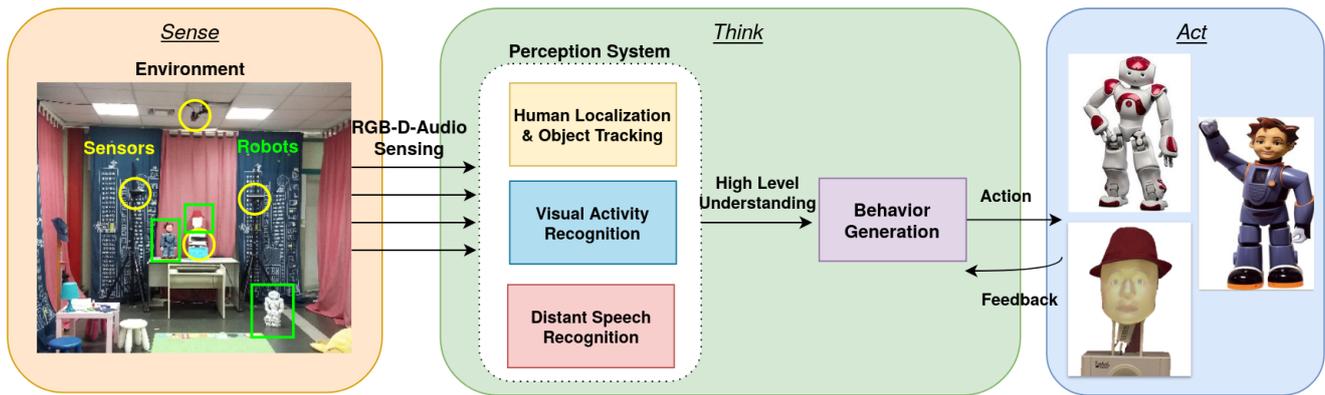
**Fig. 1** Schematic overview of the ChildBot system during Child-Robot Interaction. The multi-modal information of the child's action is received through a network of sensors placed around an interaction area. The perception system processes it and extracts high-level information about the context of action. Based on this, the behavior generation module decides and controls the robotic agents.

tual modalities allow the users to express themselves in their preferred communication means.

One of the difficulties we have to tackle in order to achieve the above, is the fact that commercial social robots have different capabilities, e.g., the NAO [4] robot is capable and adept in body movements, but incapable of facial expressions, while Furhat [3] can present a large variety of facial expressions, but can not move outside its installed space. In other cases, robots such as Zeno [5] are capable of movements and facial expressions, but they are not adept at both (Zeno's body movement lacks in comparison to NAO). Additionally, each social robot has different sensors, constraining the user to specific communication channels.

Motivated by the above, in this work we present an integrated robotic system that can be used for multiple edutainment applications, called ChildBot. To achieve this versatility, ChildBot incorporates: (i) multiple sensors and perception modules that allow the user to communicate with the robots via multiple channels, and (ii) multiple social robots, leveraging each other's strengths while diminishing their individual weaknesses.

An overview of the proposed system can be seen in Figure 1. ChildBot is developed using a *Sense-Think-Act* paradigm [22] and is indoors based, allowing us to employ external sensors that are arranged throughout a "smart space" where the interaction takes place. Robot-external sensing can both overcome common HRI problems such as occlusions, and also allows the fusion of different data streams, improving robustness and performance of the different perception modules. This way we also achieve perception of the interaction in a *robot-independent* fashion and bypass limitations of individual robotic systems sensing. In addition, this robot-agnostic architecture can easily accommodate new robots in the loop. The system coordinates a complex and continuous procedure that is required from the moment the child acts until the moment the robot responds

and vice versa. Specifically, multi-modal information flows from the sensors (Sense) to the perception modules. Then, high-level information about the context of this action is extracted, and the appropriate response/action is decided (Think). Finally, the system transmits the message to a robotic agent (Act).

To showcase the versatility and capabilities of the integrated system, we have designed five different edutainment use-cases. These use-cases are *indicative* and have been designed in order to exploit different components of the system, showcasing the large variety of applications that can be accommodated with ChildBot. The data collected by a pool of 52 children, while playing the aforementioned use-cases with the robots, allow us to objectively evaluate the performance of each module of the ChildBot system regarding its perception capabilities (accuracy-wise). Furthermore, our initial evaluation on the user experience shows encouraging results towards a complete well-designed subjective evaluation in the future.

ChildBot is an improved integrated extension of a set of preliminary conference publications by the authors on specific problems of multi-robot perception and interaction, and it presents a wide-application Child-Robot Interaction (CRI) system able to manage multitasking interaction autonomously and envelop a plethora of purposes, such as edutainment [59, 58, 15, 27]. The work presented here has integrated the previous works under a single and modular three-layer multi-robot architecture [22], includes improved perception modules, and is evaluated extensively on a larger corpus that contains spontaneous children data, more representative of CRI. To summarize, we highlight the most important contributions of the presented work:

– *An integrated system for HRI* has been designed and implemented by leveraging multiple robotic agents. The modular three-layer system architecture integrates multiple sensors, numerous perception modules and

different robotic agents, and culminates into a multi-application autonomous HRI system.

– *Perception modules for multimodal scene understanding* have been developed and adjusted in specific CRI conditions by incorporating novel approaches and extensive studies. Audio-visual active speaker localization, 6-DoF object tracking, visual activity recognition, and distant speech recognition are necessary for analyzing and tracking human behavior over time in the context of their surroundings. The perception modules of this system have been developed according to and sometimes exceed the state-of-the-art of the constituent technologies, as shown by our objective evaluations.

– *Spontaneous children data during CRI* have been collected and used for system evaluation. Indicative use-cases have been defined and implemented in order to showcase the large range of applications that ChildBot can be used for. The data collected have allowed for an extensive objective evaluation of the ChildBot capabilities during real use-case scenarios, as well as a preliminary user experience study with promising results.

## 2 Related work

Many research projects have aimed at developing robots both in ambient assisted living environments [38], as well as in well-defined and constrained environments e.g., in bathrooms for assistive bathing [71]. Some robotic agents act as companions to improve quality of life, assist with mobility, or complete household tasks [19,42], while others are designed to help people live independently and serve themselves when they face difficulties due to disabilities or old age [33,60]. Nevertheless, the intervention of robots in human life remains a controversial issue [54,20,50,66].

Regarding educational CRI applications, many previous works focus on the theoretical exploration of different social robot behaviors in the learning experience, without delving deeply into the technical aspects, but mostly using off-the-shelf solutions for environment perception. An immediate result of this is the fact that the interaction space is constrained. In [34], a study took place that involved children playing an educative mathematics scenario with a NAO robot. In [52], Saerbeck et al. studied the effect of social robot behavior on the learning performance of the subjects, in the context of a language learning task. Similar studies can be found in [26], while in [49] a humanoid robot was employed to interact with autistic children.

Notable works that have also focused on the robot perception aspect of CRI include the ALIZ-E project [8], where a complete framework was developed for multimodal CRI,

and the NAOTherapist platform [46] for upper-limb rehabilitation sessions for children with physical impairments. A similar platform was built in the INSIDE project [39], where a multimodal perception system was developed to allow autonomous interaction of a NAO robot with ASD children. Other similar projects include EMOTE [53], where the perception system focused mostly on visual communication, L2TOR [9], where a NAO robot capable of multimodal perception assumed the role of a second-language tutor, and the EASEL [62] educational CRI project. Esteban et al. [16] built a multi-sensor system for autonomous interaction of a NAO robot with autistic children to perceive different features during an interaction such as gaze estimation, action recognition, and object tracking. The capabilities of the system were sufficient for the presented tasks, but limited for a more generic interaction, and the system lacked in real use-case evaluation. In [37], Marinoiu et al. introduced an action and emotion recognition system by exploiting 2D and 3D pose estimation methods and evaluated it in a large-scale dataset of robot assisted therapy sessions of children with autism. The ANIMATAS [1] project is also worth mentioning, focusing on training researchers to advance human-machine interaction.

A general review of the perception methods used for HRI in social robots until 2014 is presented in [68]. Three important issues associated with perception systems are highlighted there: the need for developing perception systems in real environments with real data, the requirement of creating good representations in accordance with the context of the interaction, and the demand of combining an efficient perception system with reasonable robot responses in order to create pleasant HRI experiences. Some works that attempt to develop perception systems for HRI are Zaraki et al. [70], where low-level and high-level features are combined to detect a range of human-relevant features that appear during a real use-case procedure, and Valipour et al. [61], where a novel paradigm for incrementally improving visual perception of a robot during an HRI experience is proposed.

Aiming to increase performance, flexibility, and robustness, ChildBot consists of multiple robots and multimodal perception modules designed for and adapted to children, and thus allows interaction in a relatively large space for a variety of edutainment tasks. Parts of the ChildBot system have been based on our preliminary previous works, where an early design and development was presented. More specifically, a preliminary setup and evaluation of a basic architecture in a few use-cases has been presented in [59], while [58] has focused mainly on multi-party interaction via speech. In [15] the techniques for multi-view fusion of action recognition have been explored more in depth, and finally, [27] has focused on the development of tracking algorithms, essential in interactive tasks between a child

and a robot. In all these previous works, a limited and specific functionality of the system has been investigated, and evaluation has been carried out employing data acquired in a strictly controlled annotation procedure, and not spontaneous data from real interaction.

The current paper integrates all these different modules from previous works under the same, unified system, using a three-layer architecture. The modules can now work both in isolation and in synergy, and due to the system modularity, the addition or removal of a component is an easy task. Moreover, apart from the integration, a lot of effort has been dedicated to improving performance of the perception modules. For this reason, more sensors have been included, and ablation studies have been carried out, in order to validate the plausibility of the employed modules. Most importantly, contrary to all other previous works, the modules evaluation has been performed employing real-time spontaneous children data (see Sec. 4.2), which are more challenging.

## 3 System Overview

### 3.1 Perception System

This section focuses on the ChildBot perception system which provides a global and effective supervision of a progressing CRI. A core audio-visual processing technology acts as eyes, ears, and brain for the robots in use by incorporating different recognition and tracking modules. By analyzing and tracking human behavior in the context of a structured CRI experience, we target to establish common ground and share goals with the children.

The main overview of the robot-agnostic perception system can be seen in Figure 2, consisting of three main modules: *Audio-Visual Active Speaker Localization and 6-DoF Object Tracking*, *Visual Activity Recognition*, and *Distant Speech Recognition*. Four Kinect V2 sensors capture a detailed raw data representation of the environment and feed it into the perception system. The Kinect V2 sensors have been placed at different positions and viewing angles, in order to sufficiently cover the entire environment, tackle occlusion problems (self or from objects), as well as offer multiple viewpoints for visual perception. The raw data that are recorded from the sensors include both RGB and Depth, as well as audio from the microphone array (4 microphones) of each Kinect. The spatial arrangement of the sensors is presented in Figure 3(a). Subsequently, we present an overview of each perception module:

**Audio-Visual Active Speaker Localization and 6-DoF Object Tracking:** In order to have a natural interaction between robots and humans, robotic awareness of active speaker localization, as well as important object detection and tracking are essential. An effective audio-visual method for active speaker localization in HRI scenes has been

developed to track the children's body, by leveraging audio information in addition to visual information. Moreover, a module for object recognition that can detect multiple toys based on their colors and size of color regions has been incorporated in the system. Finally, a 3D tracking method has been designed for providing both the 3D location and the orientation of rigid objects.

**Visual Activity Recognition:** A visual frontend has been developed for recognizing hand gestures that accompany everyday communication, as well as more general body movements that convey specific meanings. The multiview visual activity recognition module is able to successfully recognize the child's activity, while wandering around the room and interacting with the robots and objects. The gesture recognition version of the module aims at identifying hand gestures that deliver a conceptual message during the interaction, such as waving at the robot hello or asking the robot to come closer. On the other hand, the action recognition version targets child body movements that form complex meanings, such as pantomimic movements.

**Distant Speech Recognition:** A multisensory distant speech recognition (DSR) system in Greek has been developed to enable CRI via speech. As close-talking microphones are not convenient for children and restrict their movements, we take advantage of the multiple microphone arrays located around the room recording audio, while at the same time the children can move freely and communicate hands-free with the robots. In order to make the DSR more robust and exploit the distributed microphone arrays we experiment with adaptation and fusion.

The high-level understanding that is obtained by the perception modules is then fed to the Dialog Manager, along with extra input from a Touch Screen, which is used during the interaction as an extra means of communication. According to its input, the Behavioral Generator then decides on the action that the multi-robot system should do, and forwards its decision to the actuators. The actuators in their turn respond with information back to the system.

In order to create a detailed picture of the ChildBot system, we proceed by describing extensively each perception module. Following this, we present the system architecture, the intercommunications, and the dialog management module in order to describe our complete system for CRI.

### 3.2 Perception Modules

*Audio-Visual Active Speaker Localization*

When analyzing and understanding an auditory or audiovisual scene that consists of multiple speakers, sound and speaker localization are necessary for tracking. In addition, the speaker's location is required for beamforming and for guiding the robot's attention/head in a multi-party scenario,
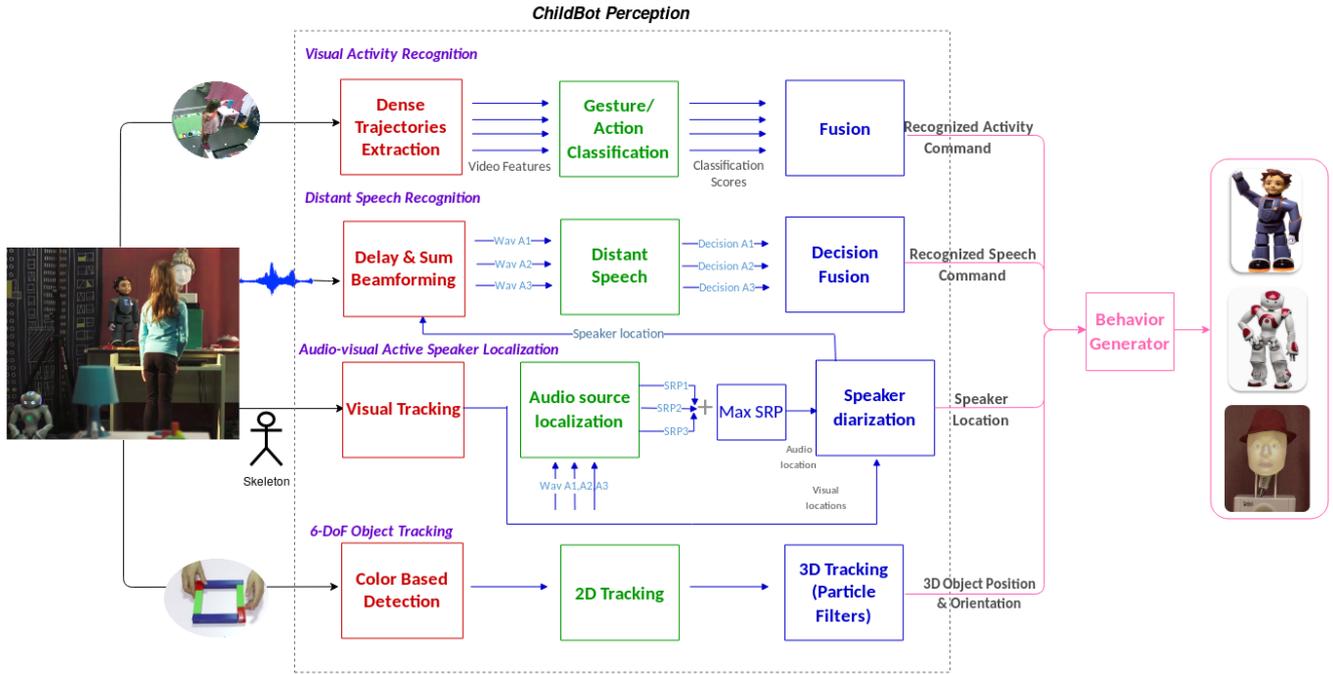
**Fig. 2** Overview of ChildBot perception modules including *Audio-Visual Active Speaker Localization and 6-DoF Object Tracking*, *Visual Activity Recognition*, and *Distant Speech Recognition*. "A" refers to microphone array and "SRP" to Steered Response Power. The modules are employed during CRI to monitor the multiple aspects of human behavior and then their outputs are fed to the robotic behavior generator module.

in order to achieve a natural and intuitive interaction. Although visual tracking can be more precise, it does not suffice when the active speaker/speakers have to be localized among other non-speaking persons in an audio-visual scene.

Various techniques for audio speaker localization [6] have been proposed in literature. Some of them have been specifically adapted to HRI setups [17, 11] for microphones mounted on robots. In our multi-robot case, microphones are external to the robots, and a fast algorithm is needed due to the real-time nature of our system. Thus, a real-time 3D audio localization SRP-PHAT (Steered Response Power - Phase Transform) system based on [14, 10] has been developed, which is robust to noise and errors. Regarding audio-visual speaker localization [23, 21, 40], several methods have been developed for RGB cameras, most of them employing Bayesian filtering techniques or fusion between audio and video features. In our case, visual tracking is accomplished via skeleton tracking, developed for Kinect sensors.

Our audio-visual active speaker localization that exploits the 3D skeleton and the microphone arrays is performed as follows: Person tracking is first achieved by retrieving the 3D skeletons from all persons present in the audio-visual scene. Auditory source localization via SRP-PHAT provides information concerning the speakers. The final active

speaker localization is performed by choosing the visual locations that are closest to the auditory ones. The speakers positions are then used for the robot's attention guiding by turning the robot's head towards the active speaker. An example of audio-visual active speaker localization can be seen in Figure 3.

### 6-DoF Object Tracking

In certain cases, children and robots may be expected to interact with various movable objects. Thus, aside from human localization which has been described above, the robot must possess an understanding of the configuration of these objects. We have developed a method for robustly tracking the 6-DoF poses of multiple objects in real-time. The main idea is to crudely detect the objects in a computationally cheap manner, and then use the detected positions to infer each object's 3D pose. The used objects are known beforehand, meaning that their shape and appearance models are predefined. The developed tracker consists of two stages: the first involves a tracking-by-detection scheme upon the color stream to locate the objects on the image plane, while the second performs an operation on the depth data to refine the first stage output and infer the remaining variables related
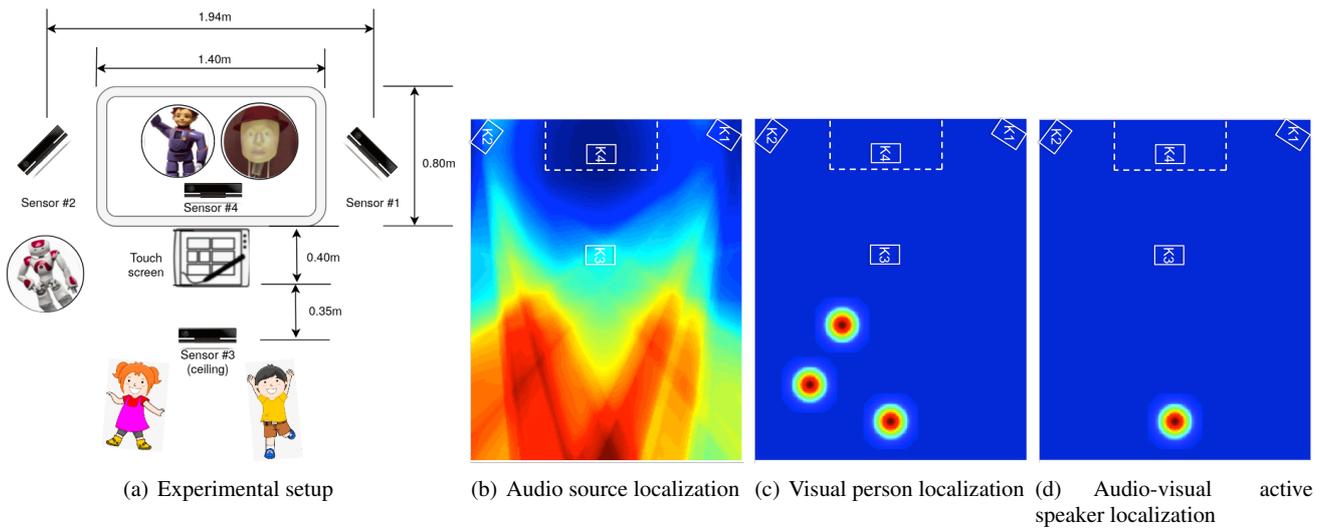
(a) Experimental setup  (b) Audio source localization  (c) Visual person localization  (d)    Audio-visual    active speaker localization

**Fig. 3** (a) Spatial arrangement of four sensors. (b-d) An example of audio-visual active speaker localization. The SRP output is shown with high values in red. Positions of the table and the four Kinects are also shown.

to the object rotations. The basic architecture is presented in Figure 4.

During the first stage, our approach uses a simple color histogram model to detect object regions, though depending on the object characteristics a variety of features could potentially be used. The histogram models are defined offline and remain unchanged during the entire tracking. The hue and saturation of the HSV color space were used to introduce sufficient robustness to brightness changes. Assuming the histograms are normalized, they define a probability distribution over the color space. Therefore, a probability map can be generated over the latest color image, which after thresholding and morphological filtering leads to a binary mask that contains the most likely object regions in the image. We choose to retain the region with the largest area, under the assumption that the remaining regions will correspond to noisy artifacts or irrelevant background objects.
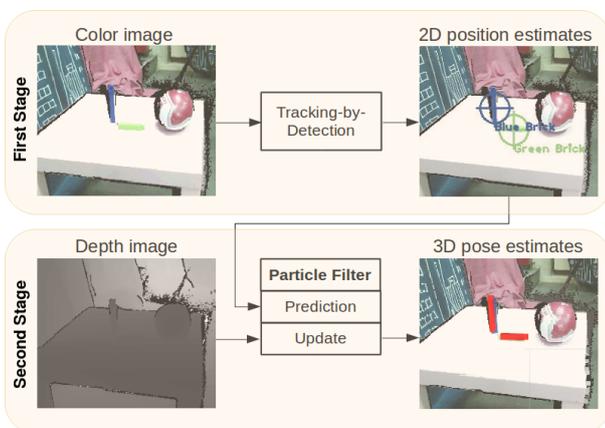


**Fig. 4** Overview of the implemented 6-DoF object tracking module (the bricks are tracked).

The center of the chosen region is taken as the object location, and a confidence score $s_k$ is produced for object $k$ .

Once the object locations have been detected on the image plane, the tracker's second stage consists of estimating the 6-DoF poses with the help of the newest depth image. The developed tracker employs particle filters and is closely based on the algorithm proposed in [67], where the hidden states are augmented with a set of binary variables that model the occlusions at each pixel. Using the camera inverse perspective mapping and the depth image, we transform the $k$-th object's 2D position estimates $\mathbf{p}_k$ into 3D estimates $\mathbf{P}_k$. The input vector for each object $k$ is then $\mathbf{u}_k = (\mathbf{P}_k - \mathbf{r}_k) \cdot s_k$, where $\mathbf{r}_k$ is the particle's position estimate from the previous time step, and $s_k$ is the confidence score produced by the tracking-by-detection module. Using a Rao-Blackwellisation technique [41], only the pose variables need to be sampled, while the occlusion variables can be marginalized out analytically. In order to prevent collisions in the object estimated configuration, the observation model is weighted by a factor that depends on the existence of mesh intersections in the particle estimates. If no intersections exist, this factor is set to 1, otherwise it is set to 0.01.

*Visual Activity Recognition*

For understanding nonverbal communication, an efficient multi-sensor visual activity recognition frontend has been developed by experimenting with Dense Trajectories features [63] along with different encoding methods and fusion schemes for visual information processing. Dense Trajectories have been chosen over convolutional neural network pretrained features, as the actions included in the state-of-the-art databases are not similar to those of children

**Fig. 5** Example of the extracted Dense Trajectories from different sensor perspectives while the child is performing the swimming pantomime (see Sec. 4.1).

and the fine-tuning of pretrained networks does not perform adequately, since in our case we have limited data from real-world CRI [15].

The main goal of this investigation is to establish a robust framework for tackling different tasks, such as generic body movements performed by kids, with limited training data. We have implemented two different versions of the module in the ChildBot system that work independently, one for gesture recognition and one for action recognition. Although the pipeline for both versions is the same, they are trained, tested, and enabled separately for hand gestures and more general body movements respectively. An example of the extracted Dense Trajectories during a pantomime performance (see Sec. 4.1) is presented in Figure 5.

In a more detailed view of the system, the recorded RGB information from each of the four RGB cameras is sampled frame by frame for the various children visual actions. Feature points are sampled for each frame on a grid and are tracked through time based on dense optical flow [18]. Multiple spatial scales are used for the sampling and the tracking independently, while the trajectories are pruned to a fixed length to avoid drifting. The computed features include the Histograms of Optical Flow (HOF) [35] and the Motion Boundary Histograms (MBH) [63] on both axes (MBHx, MBHy).

Afterwards, the features are encoded employing either the zero-order statistics Bag-of-Visual-Words (BoVW) [44] or the first-order statistics Vector of Locally Aggregated Descriptors (VLAD) [32]. Videos are classified based on their BoVW representation, using non-linear Support Vector Machines (SVMs) with the $\chi^2$ kernel [64]. In addition, the above different types of descriptors are combined with the Trajectory descriptor [63] and the Histograms of Oriented Gradients (HOG) [35], by computing distances between their corresponding BoVW histograms and adding the corresponding kernels.

The encoded features that result from VLAD are classified employing linear SVMs. After the feature extraction, we follow three different approaches - in multiple levels - for the fusion of the RGB information acquired by the multiple sensors: i) feature fusion, ii) encoding fusion, and iii) score fusion. We modify the general frameworks of BoVW and VLAD in order to deal with our proposed multi-view approach for visual activity recognition.

*Feature Fusion:* In this method, the visual information is fused at an early stage where only low-level $D$-dimensional feature descriptors $\mathbf{x}_m^i \in \mathbb{R}^D$ have been extracted, i.e., local descriptors alongside dense trajectory $m = 1,...,M_i$, from each different sensor $i = 1,...,S$. The codebook generation approach, which is based on the k-means algorithm, is modified in order to deal with the multi-view data. Given a set of feature descriptors $\mathbf{x}_m^i$, our goal is to partition the feature set into $K$ clusters $\mathbf{D} = [\mathbf{d}_1,...,\mathbf{d}_K]$, where $\mathbf{d}_k \in \mathbb{R}^D$ is the centroid of the $k$-th cluster. These $\mathbf{d}_k$ are shared between the features of all sensors. Using the notation of [44], if descriptor $\mathbf{x}_m^i$ is assigned to cluster $k$, then the indicator value $r_{m,i,k} = 1$ and $r_{m,i,\ell} = 0$ for $\ell \neq k$. The optimal $\mathbf{d}_k$ can be found by minimizing the objective function:

$$\min_{\mathbf{d}_k, r_{m,i,k}} \sum_{k=1}^{K} \sum_{i=1}^{S} \sum_{m=1}^{M_i} r_{m,i,k} \|\mathbf{x}_m^i - \mathbf{d}_k\|_2^2. \tag{1}$$

Then the encoding procedure is employed for both the BoVW and the VLAD method, resulting to a representation $\mathbf{s}_{n_j}^i$ for each trajectory $n_j$ of the $j$-th video captured by sensor $i$. The global representation $\mathbf{h}$ of the multi-view video using a sum pooling scheme is given by:

$$\mathbf{h} = \sum_{i=1}^{S} \sum_{n_j=1}^{N_j} \mathbf{s}_{n_j}^i \tag{2}$$

Finally, for the BoVW approach, an $L2$ normalization scheme [45] is applied, while for the VLAD the intra-normalization strategy proposed in [7] is followed.

*Encoding Fusion:* In this approach, a different global vector $\mathbf{h}^i$ is created by encoding the dense trajectory features using a different codebook $\mathbf{D}^i$ for each sensor $i$. For the BoVW encoding the multi-view fusion is applied by adding the $\chi^2$ kernels:

$$K(\mathbf{h}_j, \mathbf{h}_q) = \sum_{i=1}^{S} \sum_{c=1}^{N_c} \exp\left(-\frac{1}{A_c} L\left(\mathbf{h}_j^{c,i}, \mathbf{h}_q^{c,i}\right)\right), \tag{3}$$

where $\mathbf{h}_j^{c,i}$, $\mathbf{h}_q^{c,i}$ denote the BoVW representations of the $c$-th descriptor for the $j$-th and $q$-th video respectively captured by sensor $i$, and $A_c$ is the mean value of $\chi^2$ distances $L(\mathbf{h}_j^{c,i}, \mathbf{h}_q^{c,i})$ between all pairs of training samples from a specific sensor $i$. On the other hand, for the VLAD encoding a simple concatenation of the vectors that correspond to the different sensors is applied as follows: $\mathbf{h} = [\mathbf{h}^1,...,\mathbf{h}^S]$.

*Score Fusion:* For a given sensor $i$ a different SVM is trained for all employed classes and obtains the probabilities $\mathbf{P}^i$ as described in [12]. Then a softmax normalization is applied to each sensor's SVM probabilities. For the fusion of the different sensor output probabilities an average fusion is employed: $\mathbf{P} = \frac{1}{S} \sum_{i=1}^{S} \mathbf{P}^i$. Finally, the class with the highest fused score is selected, following an one-against-all approach.

### Distant Speech Recognition

In order to ensure a natural communication between humans and robots in an HRI system, it is essential to incorporate a speech recognition module. A considerable amount of distance between the robot and the users imposes the need to employ a distant speech recognition system (DSR) [65, 51] that will have to efficiently address challenging problems, such as noise and reverberation [31]. Especially when children are the end users that play and interact with the robots, the problem of speech recognition becomes much more challenging because of the special characteristics of children voices and the difficulty to acquire quality data.

In out setup, the microphone arrays distributed in space are employed for the DSR task. Children can use a set of utterances adopted for the specific context of the employed use-cases to communicate with the robots, thus our speech recognition system is grammar-based. A continuous system would require a large amount of children data to be collected; that was unfeasible in our case. Also, a grammar-based speech recognition system is adequate to fulfill the requirements of the considered use-cases. The employed language is Greek, and the set of utterances contains children possible answers in some games and some general purpose speech utterances.

The DSR system is able to detect and recognize the spoken utterances at any time, namely it is always-listening. Since speech is usually corrupted by reverberation, noise or other non-speech events, robustness is achieved via beamforming of the far-field signals and adaptation of the acoustic models.

More into the details, a sliding window of 2.5 sec duration with a 0.6 sec shift is used in order to process speech in time frames. A custom module has been developed and integrated in Robot Operating System (ROS), which allows raw audio processing from the Kinect microphones. Each speech frame is first denoised with a simple delay-and-sum beamforming applied on each available 4-channel Kinect array: The insertion of delays to the different microphone signals $a_n(t)$, allows us to align them appropriately, so as to enhance speech coming from a specific direction. For uniform linear arrays with $N$ microphones, which is also our case, if the desired direction is denoted by $\phi$, the time-delay to be applied

to each microphone is

$$\tau_n = \frac{(n-1)d\cos\phi}{c}, \tag{4}$$

where $c$ is the speed of sound and $d$ the space between microphones. The beamformed signal is denoted by:

$$y(t) = \frac{1}{N} \sum_{n=1}^{N} a_n(t - \tau_n) \tag{5}$$

The denoised signal is then fed to the DSR module where we enforce recognition of one of the pre-defined sentences. Regarding acoustic modeling, Gaussian Mixture Models and Hidden Markov Models built on cross-word tri-phone models have been trained using standard Mel-Frequency Cepstrum Coefficients-plus-derivatives features on the Logotypografia database. The Logotypografia [13] database contains clean, close-talk speech in Greek. Thus, we artificially distort the database by convolving the clean speech with room impulse responses and adding white Gaussian noise in order to match the far-field condition [51]. Maximum likelihood linear regression (MLLR) adaptation is employed to transform the means of the Gaussians on the states of the models, aiming to reduce the mismatch between the initial model and the adaptation data [69].

### 3.3 System Architecture

In this subsection we describe the backbone of the perception system: the hardware architecture, the interconnection and communication between the different modules, and how the flow of the interaction is managed.

The perception modules are integrated in the full perception system based on the following hardware architecture. The system runs on four distributed interconnected machines, three of which run the Linux operating system and the ROS, and one the Windows Operating System. Each of the three Linux machines is connected with a Kinect V2 sensor which provides raw data (i.e., color, depth, and audio). The Windows machine is also connected to a Kinect V2 sensor, and using the Microsoft SDK Kinect V2 API provides additional skeletal and tracking information. A touch screen is also connected to the Windows machine and sends feedback to the dialog module about the choices of the children. The main data processing of the perception modules takes place on each of the three Linux machines, while the multiview fusion is handled in one of the Linux machines.

Streaming of data and communication between the modules of the system flows via events that are transmitted through the TCP/IP broker which runs in the Windows machine and is provided by the IrisTK framework [55]. Under the IrisTK paradigm, we divide events in three classes:

– *Sense*: events that include information about what the sensors of the system perceive
– *Action*: events that order an actuator (i.e., a robot) to do something
– *Monitor*: background events that contain feedback information about the actions of the system (e.g., when a robot has ended speaking)

Similarly, the architecture of the system was designed based on the *Sense - Think - Act* principle [22], as shown in Figure 1. The multi-sensor setup of the system represents the Sense part, while the perception modules are classified into the Think principle. Finally, the multiple robots belong to the Act part of the architecture. This three-layer architecture allows for high-level modularity, in the sense that the different layers can be replaced/modified without affecting the others.

The broker, along with the dialog management module which will be described next, acts as a central unit that receives events from all system modules and distributes them accordingly to the appropriate modules. In addition to the high-level modularity offered by the three-layer architecture, this unit offers an extra fine-grained modularity allowing modules to be easily removed or added in the architecture by simply defining the sets of events that the module should perceive or send back to the broker.

The dialog manager is the central module of the system and models the flow of the interaction between the user and the system. The interaction is modeled using Harel states [29]- states that can be hierarchically structured, can be executed conditionally, and contain parameters that alter the flow and the transitions between them. In addition, states can be called as functions, which means that the flow of the execution will continue to the caller state, after the callee has finished his execution.

For the design and development of the statechart that models the dialog, we have included what we call "action states", i.e., states that act as a mediator between the core dialog flow and the robots. These action states contain the information that is needed to instruct the robots of the system to perform an action, and include the robot as an additional state parameter. As a result, the core dialog flow is decoupled from robot-specific details, and we avoid defining multiple similar states for different robots. This extension also gives us the capability to easily include new robots to the dialog flow by adding the robot-specific details in the action states and handling the event on the robot side. An example can be seen in Figure 6 where a state in the core dialog flow "calls" the "speak action state", including the robot that is needed to speak as an additional parameter.

From the three robots that our multi-robot system uses, the Furhat robot head and the Avatar are already integrated in the IrisTK framework. For the NAO and Zeno robots we
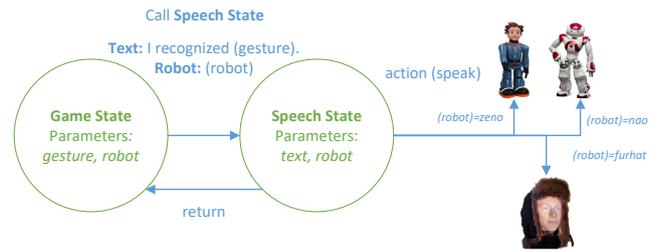


**Fig. 6** The "Speak Action State" employed for announcing a gesture during the interaction.

developed intermediate APIs that we use for communication between the robots and the dialog.

## 4 Use-cases for CRI

### 4.1 Tasks Description

A set of scenarios has been designed in order to highlight the capabilities of the system during an amusing and educative multimodal interaction between children and robots. As it has been explained extensively, our integrated system perceives various events that occur during the interaction, such as children speech and activities, along with children locations in the room, and tracking of objects. Each task focuses on different technologies and combines them appropriately to create a smooth interaction. Children are asked to complete the following tasks-games: i) "Show me the Gesture", ii) "Express the Feeling", iii) "Pantomime", iv) "Assembly Game", v) "Form a Farm".

In the first task, "Show me the Gesture", a child interacts with the robot via gestures and speech. The robot requests from the child to perform a gesture that usually denotes a meaning and tries to recognize it. It then asks the child for confirmation of the recognition. The different gestures of this game are: i) stating an agreement, ii) calling the robot to come closer, iii) asking the robot to sit down, iv) pointing an object in the room, v) asking the robot to stop, and vi) drawing a circle in the air. Except from the first gesture that is usually performed by nodding, the rest are hand gestures. The children are allowed to gesture spontaneously, as they would do when interacting with another human.

The "Express the Feeling" game motivates children to reveal their feelings using both their face and their body during an entertaining interaction with the robot. In this game, the child selects one of the cards that are depicted on the touch screen and expresses the chosen feeling. The emotions included in this game are happiness, sadness, fear, anger, surprise, and disgust. After the child's reaction, the robot also expresses the same feeling using its body and face.

"Pantomime" is a popular game, during which, one person mimes a handwork and the other figures out the depicting handwork. The child can use the whole body to mimic an activity and interact extensively with the robot. The robot

and the child repeatedly swap the roles of the mime and the guesser. The twelve activities used in this game are the following: i) cleaning a window, ii) driving a bus, iii) hammering a nail, iv) swimming, v) working out, vi) dancing, vii) reading a book, viii) digging a hole, ix) playing the guitar, x) wiping the floor, xi) dancing, and xii) ironing a shirt.

For the "Assembly Game" one or more children are asked to complete an assembly under robot supervision. Six interconnectable 3D printed bricks of different lengths are used to create rectangles and squares. The bricks are placed on a table in front of the child, with the robot standing close by. The child is responsible for the manipulation of the assembly subcomponents, while the robot provides instructions and feedback. If the child correctly completes a connection, the robot congratulates the child and proceeds to give the next instruction. If the child makes a mistake, however, the robot will attempt to recognize this mistake and will react accordingly. Aside from verbal instructions, the robot also looks and points at the bricks that it refers, for clarity.

The "Form a Farm" game is a multi-party game scenario involving two roles that can be interchanged and be equally played by both the robot and the children, aiming to entertain, educate, but also establish a natural interaction between all parties. The game involves two different roles. These roles, the picker and the guesser, can be equally played and interchanged between the children and the robot. The picker chooses an animal and utters characteristics of it. The guesser has to guess the picked animal. The interaction proceeds as follows: At first, the robot chooses a random animal and the human players take turns guessing the chosen animal. In case of a wrong guess, the robot reveals more characteristics of the animal (animal color, number of legs, animal class, e.g., mammals, reptiles). In case of correct identification of the animals, the robot asks the children to properly place the animal in a farm with some distinct segmented areas, which appears in a touch screen in front of them. In the second round, the roles are reversed: children discuss and pick an animal, and reveal one characteristic. The robot then tries to guess the picked animal. If the robot guesses correctly, the children are again asked to place the animal in the farm, otherwise they reveal more characteristics of the animal, one at at time. The game continues by interchanging the role of the guesser between children and robots in each round. The game features a total of 19 animals, and their characteristics belong to five different classes: color, size, species, number of legs, and a distinctive property, i.e. for the dog: "it's the human's best friend".

The aforementioned tasks aim to create a proper framework for multimodal communication between children and robots, as it happens between humans. This way, the tasks demonstrate the capabilities of the system, give some examples of how ChildBot can be used, and can be employed for



**Fig. 7** Data collection room and experimental setup.

system evaluation. Taking into consideration also the fact that the tasks are destined for children, the use-cases were designed under the supervision of psychologists and pilot studies with eight children were conducted, leading to the presented scenarios. Even though each task focuses on one of the system perception technologies, more than one modules are used in parallel. In Table 1, the used modules are summarized along with the eligibility of the robots to participate in each task.

### 4.2 Database

Real data obtained through HRI prove to be especially important during the process of developing a system, from the training to the evaluation stage. Such data contribute to an adaptation of the system to actual circumstances and human spontaneous behavior. Thus, an extensive data collection has taken place with the participation of a pool of 52 children, aged from six to eleven years old, in a specially designed room and in a school classroom.

Most of the data have been collected in a room that resembles a child's room and where the robotic agents and the sensors have been located, as it is presented in Figure 7. There, the data collection has been carried out in two phases. In the first one, children data have been recorded while performing certain actions and uttering certain phrases that are expected to arise throughout the interplay between them and the robots, in a strictly controlled way, when asked to do so. These data will be referred to below as *development data* since they have been used for the development of the system. In the second phase, the data have been collected during the experimental procedure where children interact with robots without any interruption or other people's intervention. The latter data will be referred to as *use-case related data*. Both types of data are equally important for CRI, as the first one is indispensable for training the perception modules on data that are relevant to use-cases, while the second one is es-

| | Distant Speech Recognition | Detect & Track | Speaker Loc. | Visual Act. Rec. | | Touch screen | Behavioral Generation | Robots | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | | Gesture | Action | | | NAO | Furhat | Zeno |
| Show me the Gesture | ✓ | | ✓ | ✓ | | | ✓ | ✓ | ✓ | ✓ |
| Pantomime | ✓ | | ✓ | | ✓ | ✓ | ✓ | ✓ | | |
| Assembly Game | | ✓ | ✓ | | | | ✓ | ✓ | | ✓ |
| Form a Farm | ✓ | | ✓ | ✓ | | ✓ | ✓ | ✓ | ✓ | ✓ |
| Express the Feeling | | | | | | ✓ | ✓ | | ✓ | ✓ |

**Table 1** Used ChildBot technologies in each use-case scenario and the eligibility of each robotic agent for participating.

sential for the testing of the behavioral monitoring software during CRI. Table 2 presents the most important recorded events during the two phases and the total number of their occurrences.

The information we collected during the data collection includes Full HD ($1920 \times 1080$) RGB and depth ($512 \times 424$) video streams from all four Kinect cameras, running at 30fps, as well as raw audio from the microphone array embedded in each Kinect sensor. By exploiting the Kinect v2 API we have also captured the following streams from the frontal Kinect sensor: (a) Skeletal information both in 2D (image) and 3D (world) coordinates; (b) Bounding boxes from face detection, facial landmarks and a facial 3D mesh.

For the development data, 28 children have participated by performing seven gestures and twelve pantomimes, and uttering 40 phrases from a vocabulary of 120 phrases. This phase is crucial for developing the perception models and adapting them to children, since they focus on speech, gestures, and actions relevant to the use-cases. Specifically, children are more spontaneous and expressive than adults and their speech is usually brief and low-voiced. Thus, in order to test the performance of ChildBot modules, it is necessary to have a plethora of children activities and utterances. Moreover, adults' data have been collected to augment the data related to the use-cases as well as to validate and highlight the need of children data for enhancing performance in the perception models.

As far as the use-case related data are concerned, 31 children with an average of 8.6 years old, 10 girls and 21 boys, had been chosen randomly from a set of volunteers that met our team in a dissemination event. All children, from six to eleven years old, spoke Greek and were able to read and write. Each child accompanied by his/her parents entered the specifically designed room and was introduced to the robots by a researcher. The child got familiarized with the room and the robots while the researcher explained the structure of the procedure and the tasks as they were presented above. Afterwards, the parents and the researcher exited the interaction space and the child played the individual games with the robots. After completing the individual interaction, a second child (which had completed the same interaction previously) entered the space and collaborated with the other child while playing the "Form a Farm" task.

| Collected Data | Event Type | # of Events |
| --- | --- | --- |
| Development Data | Utterances | 977 |
| | Gestures | 196 |
| | Pantomimes | 336 |
| Use-case Related Data | Utterances | 641 |
| | Gestures | 143 |
| | Pantomimes | 109 |

**Table 2** Statistics of the most important child activities during the data collection.

In cases where there was no second child available, an adult took its place. However, these data were removed from the subsequent evaluation. Finally, after completing the procedure, the children were asked to complete a questionnaire that included subjective statements regarding their experience. The questionnaire will be described and discussed in Section 5.2.

The above procedure has been approved by the Ethics Committee of the Athena Research Center, and also includes a consent form that had been sent by email to the parents before the experiments. In addition, all experiments have been supervised by an experienced child psychologist.

The use-case related data regarding the "Assembly Game" were collected in a Greek primary school where 21 students, 9-10 years old, participated either individually or in groups of five. For this task, a single Kinect camera and one robotic agent (NAO robot) have been chosen as a lightweight version of the system to accommodate the educational process. Such a version can be easily installed in a typical classroom and help the teacher to give a vivid lesson through a CRI experience (Figure 9).

## 5 CRI Evaluation

Each perception module of the ChildBot system has been evaluated by measuring its performance in efficient multimodal scene understanding, using the collected data. In addition, we have performed a preliminary user experience study in order to assess how the children interact and perceive the system, and collect insights towards carrying out a more complete subjective evaluation in the future.
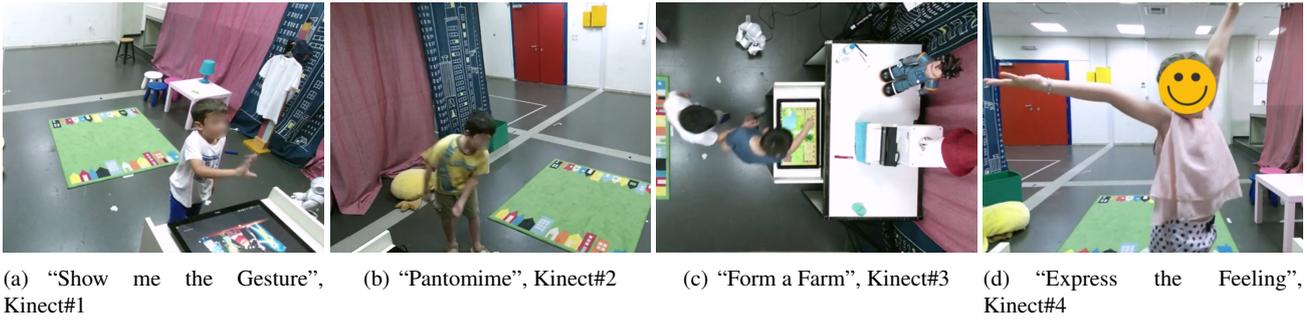
(a) "Show me the Gesture", Kinect#1

(b) "Pantomime", Kinect#2

(c) "Form a Farm", Kinect#3

(d) "Express the Feeling", Kinect#4

**Fig. 8** The four different use-cases that took place in our laboratory, each one presented from one the four different camera viewpoints of ChildBot.

### 5.1 Perception module evaluation

*Audio-Visual Active Speaker Localization*

The evaluation results of audio-visual speaker localization are presented in Table 3. For audio-only speaker localization, the employed metrics are Pcor (Percentage correct) which is the percentage of correct estimations (deviation from ground truth less than 0.5m) over all estimations, RMSE (Root Mean Square Error) between the estimation and the ground truth, and RMSEf (RMSE for estimations with error less than 0.5m - i.e., 'fine errors') . For audio-visual speaker localization, since person locations are estimated by the Kinect skeleton, the problem is essentially transformed into an active speaker localization problem. Thus, evaluation is performed in terms of correct speaker estimation, where Pcor is used, denoting the correct speaker estimations over all estimations. Audio-only localization does not perform sufficiently well yielding a Pcor of 45%, but the average RMSE is 60cm, meaning that the average source localization error is 60cm which is not very large. If both audio and visual information are used, then the active speaker localization performance is boosted to a Pcor of 86%.

*6-DoF Object Tracking*

For object tracking, we have performed both an objective evaluation and a subjective evaluation, in order to assess the performance of the 3D visual tracking module.



**Fig. 9** Setup of the "Assembly Game" at a Greek primary school.

| Audio source localization | | | Audio-visual active speaker localization |
|---|---|---|---|
| Pcor | RMSE | RMSEf | Pcor |
| 45.51% | 0.60m | 0.35m | 85.58% |

**Table 3** Evaluation of the audio-visual active speaker module.

During the objective evaluation, because it is difficult to annotate and obtain ground truth poses for 3D tracking, we have placed two static objects on a table, along with obstacles in order to add occlusions, and have moved a camera around the objects with sudden movement bursts, in order to establish the robustness of the tracker. We have compared our method with an SDF (Signed Distance Function) tracker [48]. Our results have showed that, although the SDF tracker has produced a steadier output than our tracker in cases of partial occlusions and slow camera movements, when we have introduced sudden jolts and full occlusion, the SDF tracker has been unable to continue tracking and has failed, without recovering, even when the normal conditions (no occlusion-jolts) were restored. On the other hand, our tracker has been able to successfully track the object with low error, even under full occlusion and fast movements, and recover in rare cases where the tracking has been lost, without the need for reinitialization, proving its robustness. More information on these objective evaluations can be found in [27].

During the subjective evaluation, we used the NAO robot as a supervisor for the Assembly Game which is described in Section 4. From the 21 participants, 6 played the game on their own, while the remaining children played in groups of 5. The children were required to complete two different rectangles and one square, by choosing and connecting items from 6 different brick objects.

In Table 4 we can see the results from the experiments in the form of statistics for the different assemblies. We present the percentage of total and required connections that the system recognized within a time interval of 5 sec and 20 sec, which is referred to as identification time. The term "total connections" includes both correct connections that the

child completed and mistaken connections, while the term "required connections" refers only to the correct connections needed to complete the assembly.

| | Total Connections (Recall %) | | Required Connections (Recall %) | |
|---|---|---|---|---|
| Identification Time | 5s | 20s | 5s | 20s |
| Rectangle | 70.00 | 80.00 | 50.0 | 56.25 |
| Square | 39.39 | 57.58 | 43.24 | 59.46 |

**Table 4** Statistics about the performance of the 6-DoF object tracking employed in "the Assembly Game".

*Visual Activity Recognition*

Firstly, we have examined if the age group of the participants (adults or children) has an impact on the accuracy of visual activity recognition. For both visual activity recognition tasks, we trained separate models using as training set: a) children, b) adults, c) mixed (both adults and children) and as testing set: a) children, b) adults. In Table 5, it can be noticed that the use of children training data is imperative for achieving high accuracy in children activity recognition, irrespectively of the task. This result justifies our choice for collecting development data from children movements. On the other hand, recognition models trained on mixed age groups perform better for adult gesture recognition since the diversity with which children perform the gestures accommodates the generalization of the model. For the action recognition task, children employ a wider range of different movements that adults do not use as they act stereotypically, and the mixed age training models perform worse than the adults models.

Furthermore, Tables 6 and 7 summarize the evaluation of gesture and action recognition respectively, for several combinations of different features, encodings in both the single-view case and multi-view case along with the multi-level fusion. Also, the recognition models have been trained on children development data and tested on both development and use-case related data separately, using the leave-one-out cross-validation approach.

| | | Training group | | |
|---|---|---|---|---|
| | Testing group | Adults | Children | Mixed |
| Gesture Recognition | Adults | 92.19 | 62.08 | **95.10** |
| | Children | 56.25 | **83.80** | 80.09 |
| Action Recognition | Adults | **87.36** | 72.53 | 86.26 |
| | Children | 56.51 | **74.46** | 74.26 |

**Table 5** Evaluation of the activity recognition modules after the fusion of different camera scores using MBH features and BoVW encoding.

Specifically, Table 6 presents average accuracy results (%) for the 7 gestures and a background model. Results indicate that the best multi-view model outperforms the best single-view model by about 7%, which underlines the need of a multi-view system for unrestrained CRI. For the development data, it can be seen that the combination of different types of features performs better than HOF and MBH features individually. Among the single-sensor cases, Kinect#1 (right side view) performs best as most of the kids are right-handed and they stood at approximately the same location while performing the gestures. The best recognition accuracy of 85.19% is noticed for the fusion in the final step of the procedure with the VLAD encodings and the feature combination. As far as the use-case related data are concerned, the accuracy for the single streams is moderately lower than the previous ones, which reveals the difficulty that children faced while they were trying to perform the gesture spontaneously. Also, as the children stand at a completely different location, usually closer to the cameras, the best single stream result appears for Kinect#3 (floor plan view). Regarding the fusion of the different streams, recognition performance is slightly better for the encodings fusion than the scores fusion, and it approaches 74%. More generally, VLAD encodes more effectively the visual information than the BoVW, since it contains rich information about the distribution of the visual words. Finally, we have to note that, as nodding requires a gentle movement, it is usually confused with the background movement.

In order to verify the appropriateness of the proposed visual activity recognition system in more challenging tasks, we evaluate the visual activity recognition system for the pantomimes. Table 7 presents the average accuracy results (%) for the 12 pantomimes as well as the background model. The fusion of the single-view information enhances remarkably the performance of the recognition, as was observed in the gesture case too. The highest accuracy for testing on development data appears with VLAD encodings in scores fusion, since the visual information in these data is more consistent compared to use-case related data, e.g. similar time duration of a pantomime or similar children locations in the room. Regarding the single-view case, in both types of data, the right side view Kinect#1 appears to be the best perspective for the trained models. It can be noticed that in use-case related data MBH yields slightly better results than feature combination. Moreover, feature fusion, i.e. fusion of the information at an early step of the entire procedure results in the best performance regardless of the type of the encodings.

In conclusion, the accuracy of the visual activity recognition is lower in use-case related data since children act more spontaneously while they move around the room and interact freely with the robots. Furthermore, due to the fact that the variation of the visual information in the pantomime task is larger than in the gesture task, the early fusion of the

| | Development Data | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Single Camera | | | | Fusion | | | | | |
| Features | Kinect #1 | Kinect #2 | Kinect #3 | Kinect #4 | Features | | Encodings | | Scores | |
| | BoVW | | | | BoVW | VLAD | BoVW | VLAD | BoVW | VLAD |
| HOF | 70.83 | 70.37 | 69.21 | 63.43 | 71.76 | 74.07 | 77.78 | 81.48 | 75.93 | 81.94 |
| MBH | 76.85 | 67.82 | 68.29 | 65.28 | 76.39 | 76.85 | 81.02 | 81.48 | 82.87 | 83.80 |
| Traj.+HOG+HOF+MBH | **77.78** | 73.84 | 73.61 | 75.00 | 81.48 | 82.87 | 82.87 | 83.80 | 82.87 | **85.19** |
| | Use-case Related Data | | | | | | | | | |
| | Single Camera | | | | Fusion | | | | | |
| Features | Kinect #1 | Kinect #2 | Kinect #3 | Kinect #4 | Features | | Encodings | | Scores | |
| | BoVW | | | | BoVW | VLAD | BoVW | VLAD | BoVW | VLAD |
| HOF | 56.92 | 54.49 | 57.10 | 51.97 | 54.56 | 71.61 | 58.01 | 74.73 | 63.26 | 74.83 |
| MBH | 62.70 | 56.47 | 60.15 | 54.25 | 65.32 | 72.70 | 67.72 | 72.52 | 66.73 | 72.72 |
| Traj.+HOG+HOF+MBH | 57.96 | 54.08 | **67.03** | 59.16 | 61.51 | 69.85 | 63.38 | **73.95** | 64.82 | **73.35** |

**Table 6** Average classification accuracy (%) for the employed 8 gestures. Results on both development and use-case related data are shown for the different features, encoding and fusion methods of the activity recognition module.

| | Development Data | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Single Camera | | | | Fusion | | | | | |
| Features | Kinect #1 | Kinect #2 | Kinect #3 | Kinect #4 | Features | | Encodings | | Scores | |
| | BoVW | | | | BoVW | VLAD | BoVW | VLAD | BoVW | VLAD |
| HOF | 68.31 | 56.31 | 48.62 | 53.85 | 66.77 | 67.08 | 68.00 | 69.23 | 68.62 | 75.50 |
| MBH | 70.77 | 60.92 | 61.85 | 55.22 | 76.00 | 76.69 | 76.92 | 76.92 | 74.46 | 76.50 |
| Traj.+HOG+HOF+MBH | **73.85** | 63.38 | 60.00 | 61.45 | 75.08 | 76.92 | 77.23 | 77.85 | 75.08 | **79.00** |
| | Use-case Related Data | | | | | | | | | |
| | Single Camera | | | | Fusion | | | | | |
| Features | Kinect #1 | Kinect #2 | Kinect #3 | Kinect #4 | Features | | Encodings | | Scores | |
| | BoVW | | | | BoVW | VLAD | BoVW | VLAD | BoVW | VLAD |
| HOF | 46.34 | 46.19 | 25.50 | 47.70 | 63.08 | 61.02 | 49.87 | 56.17 | 52.59 | 57.99 |
| MBH | **61.42** | 46.28 | 31.59 | 45.57 | **70.25** | 67.97 | 57.70 | 59.04 | 62.18 | 62.49 |
| Traj.+HOG+HOF+MBH | 52.59 | 46.74 | 36.62 | 48.16 | 63.52 | **69.37** | 60.75 | 61.55 | 55.00 | 64.90 |

**Table 7** Average classification accuracy (%) for the employed 13 pantomimes. Results on both development and use-case related data for the different features, encoding and fusion methods of the activity recognition module are depicted.

features performs better for the pantomime while the scores fusion is satisfactory for the "Show me the Gesture" since only one camera is adequate to recognize the gesture.

*Distant Speech Recognition*

Two sets of data have been employed for the offline evaluation of the DSR task: the development set and the use-case related set, both consisting of children data. As stated before, the DSR system is grammar-based, namely depending on the context of the application, and there is a specific set of commands that the users adopt in order to communicate with the robot. Thus, we have designed one set for the "Show me the Gesture", the "Express the Feeling", and the "Pantomime" games and another one for the cooperative game, i.e. the "Form a Farm" game. The grammar size and other statistics concerning the two datasets can be found in Table 8.

The development data have been used for adapting speech models and testing them. Results are presented in Table 9 in terms of word and sentence accuracies, denoted by WCOR and SCOR respectively. Four different adaptation schemes have been tested for comparison: In the "No-adapt" case, the employed models have been trained on the Logotypografia database which contains adult data. The available children data included in the development

set have been used for testing. In the "Adults" case, speech models have been adapted to a small amount of adult data and tested both with adult and children data. In the "Children" case, data from 20 out of 28 participants of the development set have been used to adapt speech models globally, i.e. data from the Kinect arrays have been used to adapt a single model. The remaining 8 participants form the testing set. The adaptation and testing has been 4-fold cross-validated. Finally, in the "Mixed" case we have used both adult and children data to adapt the models and then test them separately on adult and children data.

Speech recognition achieves satisfactory performance for adults even without adaptation. However, adaptation indeed improves performance for all cases, even when it is performed in a different age group than testing. Results indicate that the best performance is obtained for adapting and testing on the same group, which was expected. The best achieved results are 98.87% for adults and 95.50% for children in terms of SCOR. The results concerning children underline the need and importance of collecting children data. Performance is boosted, from 75.3% to 97.8% for WCOR and from 71.2% to 95.5% for SCOR, when children data have been used for adaptation and testing, who are indeed the target group of the system.

All the above results were referring to the development set where data collection was controlled and guided. The

| | Development set | | | Use-case Related set | | |
|---|---|---|---|---|---|---|
| | # speakers | # utterances | grammar size | # speakers | # utterances | grammar size |
| Single Game | 28 | 642 | 75 | 31 | 426 | 157 |
| Cooperative Game | 28 | 335 | 58 | 9 pairs | 215 | 113 |

**Table 8** Data statistics for the children DSR task.

| | DSR-Adaptation scheme | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | No-adapt | | Adults | | Children | | Mixed | |
| Test | WCOR | SCOR | WCOR | SCOR | WCOR | SCOR | WCOR | SCOR |
| Adults | 97.54 | 91.25 | 99.58 | **98.87** | 96.73 | 93.20 | 99.50 | 98.43 |
| Children | 79.06 | 69.95 | 75.31 | 71.20 | 97.81 | **95.50** | 90.71 | 82.06 |

**Table 9** Evaluation of the DSR recognition on the development data.

| | No-adapt | | | Adapt-all | | |
|---|---|---|---|---|---|---|
| | WCOR | SCOR | LabelCOR | WCOR | SCOR | LabelCOR |
| Single Game | 56.68 | 29.52 | 55.12 | 59.64 | 43.77 | 55.12 |
| Cooperative Game | 72.95 | 61.02 | 63.16 | 78.00 | 67.69 | 70.51 |

**Table 10** Average word (WCOR), sentence accuracy (SCOR), and Label accuracy (LabelCOR) (%) for the children DSR task.

use-case related set contains data in the wild, namely, while children were actually playing with the robots, their speech data were collected. Although the children had received some instructions concerning the utterances they could use, it is obvious that they were not followed in most cases. Thus, after the data annotation, new grammar sets have been formed in order to incorporate new phrases.

The use-case related data results are depicted in Table 10, in terms of WCOR, SCOR and LabelCOR. LabelCOR refers to the percentage of correct recognition of the semantic content. For example, there can be various ways to express a negation: "no", "no, I don't know", "no, I did not find it", etc. All similar phrases in terms of content are given a specific label, and after speech recognition data post-processing calculates the score of the correct recognition of labels. Adaptation has been performed using the development set. We notice that adaptation boosts the performance, achieving a percentage of 59.64% for the gesture and pantomime games and 78% for the farm game in terms of WCOR. The lower percentage of gesture and pantomime games can be attributed to the larger grammar size, the distance between the speakers and the microphones, and the relatively large variations of speaker orientation.

## 5.2 User Experience Study

We have also performed a user experience study which included a pool of 52 children, from six to eleven years old, participating in the designed interaction described in Section 4.2. The purpose of this study is to collect objective statistics

and insights, as well as get a measure of the system's ability to accommodate a complete CRI.

**Objective Statistics** Regarding the individual tasks ("Show me the Gesture", "Express the feeling", "Pantomime"), where each of the 31 kids participated alone, all of them were able to complete successfully the games. The average duration of completion, including the introduction by the robots, is 9 minutes, with a variance of 2 minutes.

In 32% of the cases, human (verbal) intervention was required up to two times during the experimental flow, when the children became confused or had questions about the procedure. For example, some children asked for a confirmation about what to do or needed a prompt in order to act. Such possible deviations from the designed scenario have been overcome by enabling the dialog manager to recognize these cases (e.g., if the child is silent for a moderate period of time), and getting the robots prompt the children or ask them to repeat their utterance/activity. In cases where children were expected to say something or their speech was not recognized, robots requested for repetitions up to two successive times, while in case of a child's action, the robots asked for repetition only once.

For the collaborative "Form a Farm" game, played by two children, it was observed that younger children faced difficulties with the rules of the game, even though primary school children are familiar with the animals of a farm. As a result, kids with ages six and seven played the game following the guidelines offered by one adult. The rest of the children played the game without any guidance. The average duration of the game was 8 minutes. In total, children assumed the role of the guesser for 24 rounds and found the correct answer using 2.4 guesses on average and 4 guesses maxi-
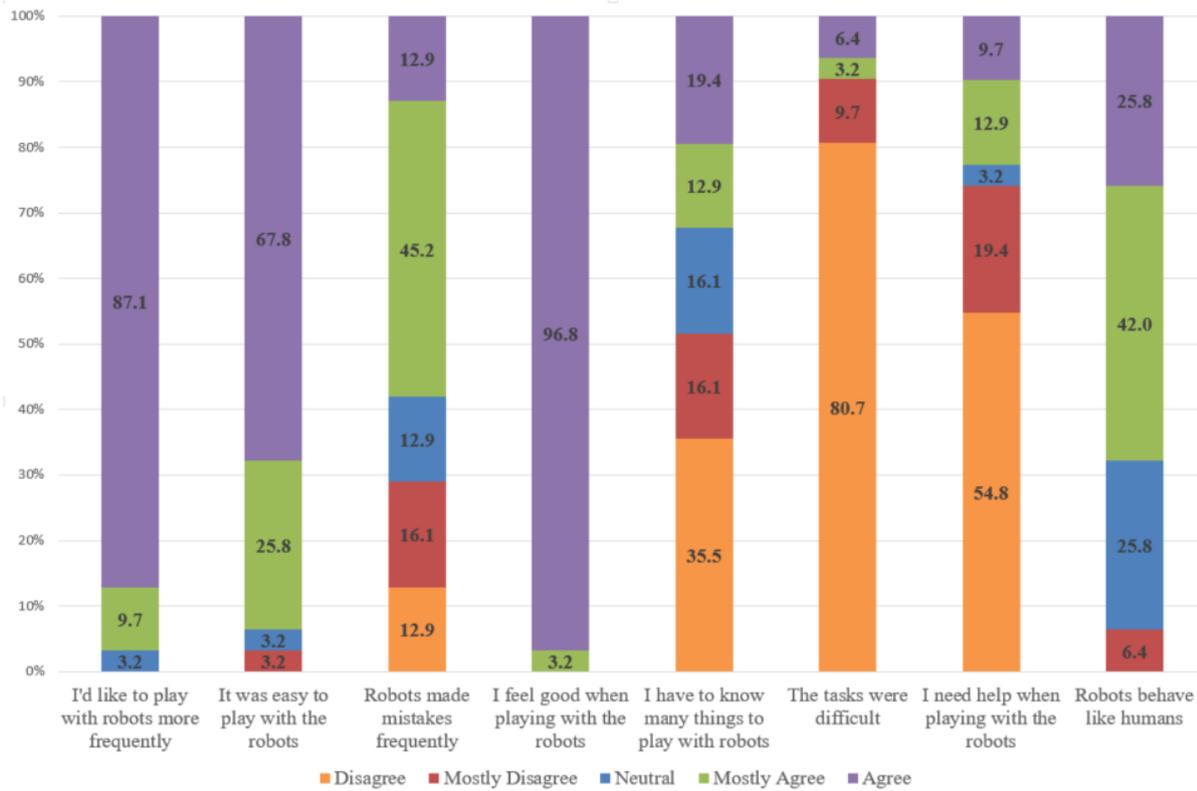
**Fig. 10** User experience of the entire ChildBot system. After each completed interaction, children were asked to fill a questionnaire with the shown questions on a Likert-type scale from 1 to 5 labeled as shown.

mum. On the other hand, the robot assumed the guesser role for 22 rounds and found the correct answer in 2.2 guesses on average, with a maximum of 6. Children did not manage to identify the picked animal in 4% of the guesses, while the robot in 32%. Generally, the children managed to guess the picked animal more easily, since the robot was programmed to always reveal more general animal characteristics in the beginning, and proceeding with more specific details.

**User experience assessment** Regarding subjective evaluation of the experience, children were asked to complete a questionnaire containing the subjective statements that can be seen in Figure 10. Each statement was accompanied by a 5-point Likert type ordinal scale labeled from "disagree" to "agree", using smiley faces [28].

We also included two multiple choice questions asking children to justify which use-case was the most preferable and why, and which perception ability of the robots make them popular to the kids, as a means of getting more insight on their preferences. In general, the majority of children (12/31) stated that their favorite use-case was "Pantomime", due to the movements of the robot. As we can see in Figure 10, the majority of the children (27/31) stated that they like playing with the robots, while 22 of them enjoyed playing because robots understood both their movements and speech. Many of them (21/31) also found the interaction and

the use-cases easy to follow, without needing external help (19/31). Furthermore, children tended to agree (20 positive answers out of 31) that robots behave like humans. By analyzing the questionnaire responses, we noticed that older children stated that they did not need prior knowledge in order to play with the robots, while younger children stated that they did.

Similarly, in the assembly task that was evaluated in the primary school, 21 children were asked to express their opinion for the interaction. The questions are presented in Table 11, and the available responses were a 3-point Likert scale (Disagree - Neutral - Agree). The Table also presents the questionnaire results after being mapped to a scale of 0-2, with 0 being the most negative. Their answers indicate that children were pleased with the interaction (1.81 MOS on whether they would like to play again with the robot and the comfortableness of the interaction). However, clearly the robot supervision for the assembly task has room for improvements, since although the instructions of the robot were very clear (1.95 MOS), children were neutral on whether they were actually helpful (1.10 MOS), or wrong (0.95 MOS).

**Discussion** In general, the evaluation of user experience during interaction with our multi-robot, multi-tasking, and multi-sensor robotic system provided encouraging results,

proving that the system is technically capable of accommodating a complete CRI experience, with some adult intervention needed in certain cases and mainly for the collaborative task. Of course, there exists room for improvement, since many children stated that robots made mistakes frequently (18/31). In the future, we aim to also conduct a subjective evaluation focused on the pedagogical aspect of the system, based on the insights collected during this initial study.

| Question | Mean Opinion Score |
|---|---|
| Were you comfortable working with the robot? | 1.81 |
| Would you play with a robot again, sometime? | 1.81 |
| Was the robot helpful? | 1.10 |
| Did the robot make a lot of mistakes? | 0.95 |
| Were the robots instructions clear? | 1.95 |

**Table 11** Questions and results of the questionnaire presented to the children, following their "Assembly Game" with the robot.

# 6 Conclusion

In this paper we have presented ChildBot, a multimodal perception framework that is the culmination and the extension of several earlier works by the authors in multimodal perception and CRI. ChildBot constitutes a CRI framework with multiple robotic agents that can be successfully used for edutainment purposes, and its perception system includes several different modules: audio-visual active speaker localization, and 6-DoF object tracking, visual activity recognition, and distant speech recognition. The effectiveness and successful interconnection of the modules has been demonstrated via five indicative edutainment CRI use-cases, each using a different subset of the various perception modules.

In order to validate the performance and the capabilities of our system for CRI, we have carried out an extensive objective evaluation of the developed perception modules, as well as a user experience study that provides valuable initial insights for the interaction with ChildBot. The experiments took place in a specially designed area that was decorated to resemble a child's room, where we collected both development data necessary for training the individual system modules, as well as use-case related data that were essential for testing the system performance during actual CRI.

Our results have showed that the individual perception technologies successfully capture the environment surrounding the interaction with a high degree of accuracy, while the user experience study showed that children enjoyed playing with different robots.

For future work, we would like to also extend Child-Bot for other applications, and it would be interesting to see how some of the novel perception methods we presented can generalize to other fields, such as rehabilitation or assistive applications for ASD children. Further, we aim to conduct a more thorough subjective evaluation on the pedagogical aspect of the system.

In conclusion, our work shows that through the integration of multiple robots, sensors, and modalities, we can achieve a high level of unconstrained and autonomous CRI, opening up new prospects for the educational and entertainment social robotics of the future.

# References

1. ANIMATAS project. http://www.animatas.eu/.
2. BabyRobot project. http://babyrobot.eu.
3. Furhat Robotics. http://furhatrobotics.com.
4. NAO, Softbank Robotics. https://www.softbankrobotics.com/.
5. Robokind. Advanced Social Robots. http://robokind.com/.
6. X. Anguera, S. Bozonnet, N. Evans, C. Fredouille, G. Friedland, and O. Vinyals. Speaker diarization: A review of recent research. *IEEE Trans. Audio, Speech, and Language Processing*, 20(2):356–370, 2012.
7. R. Arandjelovic and A. Zisserman. All about VLAD. In *Proc. CVPR*, 2013.
8. T. Belpaeme, P.E. Baxter, R. Read, R. Wood, H. Cuayáhuitl, B. Kiefer, S. Racioppa, I. Kruijff-Korbayová, G. Athanasopoulos, V. Enescu, et al. Multimodal child-robot interaction: Building social bonds. *Journal of Human-Robot Interaction*, 1(2):33–53, 2012.
9. T. Belpaeme, J. Kennedy, P. Baxter, P. Vogt, E. E. Krahmer, S. Kopp, K. Bergmann, P. Leseman, A. C. Küntay, T. Göksun, et al. L2TOR-second language tutoring using social robots. In *Proc. of the ICSR 2015 WONDER Workshop*, 2015.
10. A. Brutti, M. Omologo, P. Svaizer, and C. Zieger. Classification of acoustic maps to determine speaker position and orientation from a distributed microphone network. In *Proc. ICASSP*, 2007.
11. J. Cech, R. Mittal, A. Deleforge, J. Sanchez-Riera, X. Alameda-Pineda, and R. Horaud. Active-speaker detection and localization with microphones and cameras embedded into a robotic head. In *Proc. Humanoid Robots*, 2013.
12. C.C. Chang and C. J. Lin. LIBSVM: a library for support vector machines. *ACM Trans. Intelligent Systems and Technology*, 2(3):27, 2011.
13. V. Digalakis, D. Oikonomidis, D. Pratsolis, N. Tsourakis, C. Vosnidis, N. Chatzichrisafis, and V. Diakoloukas. Large vocabulary continuous speech recognition in Greek: Corpus and an automatic dictation system. In *Proc. Interspeech*, 2003.
14. H. Do, H.F. Silverman, and Y. Yu. A real-time SRP-PHAT source location implementation using stochastic region contraction (SRC) on a large-aperture microphone array. In *Proc. ICASSP*, 2007.
15. N. Efthymiou, P. Koutras, P. P. Filntisis, G. Potamianos, and P. Maragos. Multi- view fusion for action recognition in child-robot interaction. In *Proc. ICIP*, 2018.
16. P. G. Esteban, P. Baxter, T. Belpaeme, E. Billing, H. Cai, H. L. Cao, M. Coeckelbergh, C. Costescu, D. David, A. De Beir, et al. How to build a supervised autonomous system for robot-enhanced therapy for children with autism spectrum disorder. *Paladyn, Journal of Behavioral Robotics*, 8(1):18–38, 2017.
17. C. Evers, Y. Dorfan, S. Gannot, and P.A. Naylor. Source tracking using moving microphone arrays for robot audition. In *Proc. ICASSP*, 2017.
18. G. Farnebäck. Two-frame motion estimation based on polynomial expansion. In *Scandinavian conference on Image analysis*, 2003.
19. D. Fischinger, P. Einramhof, K. Papoutsakis, W. Wohlkinger, P. Mayer, P. Panek, S. Hofmann, T. Koertner, A. Weiss, A. Argyros, and M. Vincze. Hobbit, a care robot supporting independent living at home: First prototype and lessons learned. *Robotics and Autonomous Systems*, 75:60–78, 2016.
20. S. Frennert and B. Östlund. Review: Seven matters of concern of social robots and older people. *International Journal of Social Robotics*, 6:299–310, 2014.
21. G. Garau, A. Dielmann, and H. Bourlard. Audio-visual synchronisation for speaker diarisation. In *Proc. Interspeech*, 2010.
22. E. Gat, R. P. Bonnasso, R. Murphy, et al. On three-layer architectures. *Artificial Intelligence and Mobile Robots*, 195:210, 1998.
23. I.D. Gebru, C. Evers, P.A. Naylor, and R. Horaud. Audio-visual tracking by density approximation in a sequential Bayesian filtering framework. In *Proc. HSCMA*, 2017.
24. M. Gombolay, X. J. Yang, B. Hayes, N. Seo, Z. Liu, S. Wadhwania, T. Yu, N. Shah, T. Golen, and J. Shah. Robotic assistance in the coordination of patient care. *The International Journal of Robotics Research*, 37(10):1300–1316, 2018.
25. M.A. Goodrich and A. Schultz. Human-robot interaction: a survey. *Foundations and Trends in Human-Computer Interaction*, 1(3):203–275, 2007.
26. G. Gordon, C. Breazeal, and S. Engel. Can children catch curiosity from a social robot? In *Proc. HRI*, 2015.
27. J. Hadfield, P. Koutras, N. Efthymiou, G. Potamianos, C. S. Tzafestas, and P. Maragos. Object assembly guidance in child-robot interaction using RGB-D based 3D tracking. In *Proc. IROS*, 2018.
28. L. Hall, C. Hume, and S. Tazzyman. Five degrees of happiness: Effective smiley face Likert scales for evaluating with children. In *Proc. 15th International Conference on Interaction Design and Children*, 2016.
29. D. Harel. Statecharts: A visual formalism for complex systems. *Science of computer programming*, 8(3):231–274, 1987.
30. W. Huo, S. Mohammed, J. C. Moreno, and Y. Amirat. Lower limb wearable robots for assistance and rehabilitation: A state of the art. *IEEE Systems Journal*, 10(3):1068–1081, 2016.
31. C.T. Ishi, S. Matsuda, T. Kanda, T. Jitsuhiro, H. Ishiguro, S. Nakamura, and N. Hagita. A robust speech recognition system for communication robots in noisy environments. *IEEE Trans. Robotics*, 24(3):759–763, 2008.
32. H. Jégou, M. Douze, C. Schmid, and P. Pérez. Aggregating local descriptors into a compact image representation. In *Proc. CVPR*, 2010.
33. N. Kardaris, I. Rodomagoulakis, V. Pitsikalis, A. Arvanitakis, and P. Maragos. A platform for building new human-computer interface systems that support online automatic recognition of audiogestural commands. In *Proc. ACMMM*, 2016.
34. J. Kennedy, P. Baxter, E. Senft, and T. Belpaeme. Higher nonverbal immediacy leads to greater learning gains in child-robot tutoring interactions. In *Proc. ICSR*, 2015.
35. I. Laptev, M. Marszalek, C. Schmid, and B. Rozenfeld. Learning realistic human actions from movies. In *Proc. CVPR*, 2008.
36. L. Lucignano, F. Cutugno, S. Rossi, and A. Finzi. A dialogue system for multimodal human-robot interaction. In *Proc. ICMI*, 2013.
37. E. Marinoiu, M. Zanfir, V. Olaru, and C. Sminchisescu. 3d human sensing, action and emotion recognition in robot assisted therapy of children with autism. In *Proc. CVPR*, 2018.
38. P. Mayer, C. Beck, and P. Panek. Examples of multimodal user interfaces for socially assistive robots in ambient assisted living environments. In *Proc. CogInfoCom*, 2012.
39. F. S. Melo, A. Sardinha, D. Belo, M. Couto, M. Faria, A. Farias, H. Gamba, C. Jesus, M. Kinarullathil, P. Lima, L. Luz, A. Mateus, I. Melo, P. Moreno, D. Osrio, A. Paiva, J. Pimentel, J. Rodrigues, P. Sequeira, R. Solera-Urea, M. Vasco, M. Veloso, and R. Ventura. Project INSIDE: towards autonomous semi-unstructured human-robot social interaction in autism therapy. *Artificial Intelligence in Medicine*, 96:198–216, 2019.
40. V.P. Minotto, C.R. Jung, and B. Lee. Multimodal multi-channel on-line speaker diarization using sensor fusion through svm. *IEEE Trans. Multimedia*, 17(10):1694–1705, 2015.
41. K. Murphy and S. Russell. Rao-blackwellised particle filtering for dynamic bayesian networks. In *Sequential Monte Carlo methods in practice*, pages 499–515. Springer, 2001.
42. M. Nani, P. Caleb-Solly, S. Dogramadzi, T. Fear, and H. van den Heuvel. MOBISERV: an integrated intelligent home environment for the provision of health, nutrition and mobility services to the elderly. In *Proc. 4th Companion Robotics Workshop*, 2010.

43. T. Pachidis, E. Vrochidou, V.G. Kaburlasos, S. Kostova, M. Bonković, and V. Papić. Social robotics in education: State-of-the-art and directions. In *Proc. International Conference on Robotics in Alpe-Adria Danube Region*, 2018.

44. X. Peng, L. Wang, X. Wang, and Y. Qiao. Bag of visual words and fusion methods for action recognition: Comprehensive study and good practice. *Computer Vision and Image Understanding*, 150:109–125, 2016.

45. F. Perronnin, J. Sánchez, and T. Mensink. Improving the Fisher kernel for large-scale image classification. In *Proc. ECCV*, 2010.

46. J. C. Pulido, J. C. González, C. Suárez-Mejías, A. Bandera, P. Bustos, and F. Fernández. Evaluating the child–robot interaction of the NAOTherapist platform in pediatric rehabilitation. *International Journal of Social Robotics*, 9(3):343–358, 2017.

47. Z. Qian and Z. Bi. Recent development of rehabilitation robots. *Advances in Mechanical Engineering*, 7(2), 2015.

48. C. Y. Ren, V. Prisacariu, O. Kaehler, I. Reid, and D. Murray. 3D tracking of multiple objects with identical appearance using RGB-D input. In *Proc. International Conference on 3D Vision*, 2014.

49. B. Robins, K. Dautenhahn, R. Te Boekhorst, and A. Billard. Robotic assistants in therapy and education of children with autism: can a small humanoid robot help encourage social interaction skills? *Universal Access in the Information Society*, 4(2):105–120, 2005.

50. H. Robinson, B. MacDonald, N. Kerse, and E. Broadbent. The psychosocial effects of a companion robot: A randomized controlled trial. *Journal of the American Medical Directors Association*, 14:661 – 667, 2013.

51. I. Rodomagoulakis, A. Katsamanis, G. Potamianos, P. Giannoulis, A. Tsiami, and P. Maragos. Room-localized spoken command recognition in multi-room, multi-microphone environments. *Computer Speech & Language*, 46:419–443, 2017.

52. M. Saerbeck, T. Schut, C. Bartneck, and M.D. Janse. Expressive robots in education: varying the degree of social supportive behavior of a robotic tutor. In *Proc. CHI*, 2010.

53. S. Serholt, W. Barendregt, T. Ribeiro, G. Castellano, A. Paiva, A. Kappas, R. Aylett, and F. Nabais. EMOTE: Emboided-perceptive tutors for empathy-based learning in game environment. In *ECGBL*, 2013.

54. M. Shishehgar, D. Kerr, and J. Blake. A systematic review of research into how robotic technology can help older people. *Smart Health*, 7:1–18, 2018.

55. G. Skantze and S. Al Moubayed. IrisTK: a statechart-based toolkit for multi-party face-to-face interaction. In *Proc. ICMI*, 2012.

56. Rainer Stiefelhagen, C Fugen, R Gieselmann, Hartwig Holzapfel, Kai Nickel, and Alex Waibel. Natural human-robot interaction using speech, head pose and gestures. In *2004 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)(IEEE Cat. No. 04CH37566)*, volume 3, pages 2422–2427. IEEE, 2004.

57. A. Sullivan and M. U. Bers. Dancing robots: integrating art, music, and robotics in singapores early childhood centers. *International Journal of Technology and Design Education*, 28(2):325–346, 2018.

58. A. Tsiami, P. P. Filntisis, N. Efthymiou, P. Koutras, G. Potamianos, and P. Maragos. Far-field audio-visual scene perception of multi-party human-robot interaction for children and adults. In *Proc. ICASSP*, 2018.

59. A. Tsiami, P. Koutras, N. Efthymiou, P. P. Filntisis, G. Potamianos, and P. Maragos. Multi3: Multi-sensory perception system for multi-modal child interaction with multiple robots. In *Proc. ICRA*, 2018.

60. M. V. Soler, L. Agüera-Ortiz, J. O. Rodríguez, C. M. Rebolledo, and A. P. Muñoz et al. Social robots in advanced dementia. *Frontiers in Aging Neuroscience*, 7:133, 2015.

61. S. Valipour, C. Perez, and M. Jagersand. Incremental learning for robot perception through HRI. In *Proc. IROS*, 2017.

62. V. Vouloutsi, M. Blancas, R. Zucca, P. Omedas, D. Reidsma, D. Davison, V. Charisi, F. Wijnen, J. van der Meij, V. Evers, et al. Towards a synthetic tutor assistant: the EASEL project and its architecture. In *Conference on Biomimetic and Biohybrid Systems*, 2016.

63. H. Wang, A. Klaser, C. Schmid, and C.L. Liu. Action recognition by dense trajectories. In *Proc. CVPR*, 2011.

64. Heng Wang, Muhammad Muneeb Ullah, Alexander Kläser, Ivan Laptev, and Cordelia Schmid. Evaluation of local spatio-temporal features for action recognition. In *BMVC*, 2009.

65. M. Wölfel and J. McDonough. *Distant Speech Recognition*. John Wiley & Sons, 2009.

66. Y. H. Wu, C. Fassert, and A. S. Rigaud. Designing robots for the elderly: appearance issue and beyond. *Archives of Gerontology and Geriatrics*, 54:121–126, 2012.

67. M. Wüthrich, P. Pastor, M. Kalakrishnan, J. Bohg, and S. Schaal. Probabilistic object tracking using a range camera. In *Proc. IROS*, 2013.

68. H. Yan, M. H. Ang, and A. N. Poo. A survey on perception methods for human–robot interaction in social robots. *International Journal of Social Robotics*, 6(1):85–119, 2014.

69. S. J. Young, D. Kershaw, J. Odell, D. Ollason, V. Valtchev, and P. Woodland. *The HTK Book Version 3.4*. Cambridge University Press, 2006.

70. A. Zaraki, M. Pieroni, D. De Rossi, D. Mazzei, R. Garofalo, L. Cominelli, and M. B. Dehkordi. Design and evaluation of a unique social perception system for human–robot interaction. *IEEE Transactions on Cognitive and Developmental Systems*, 9(4):341–355, 2017.

71. A. Zlatintsi, I. Rodomagoulakis, V. Pitsikalis, P. Koutras, N. Kardaris, X. Papageorgiou, C. Tzafestas, and P. Maragos. Social human-robot interaction for the elderly: two real-life use cases. In *Proc. HRI*, 2017.