



**HAL**  
open science

# Neural Network Reinforcement Learning for Audio-Visual Gaze Control in Human-Robot Interaction

Stéphane Lathuilière, Benoît Massé, Pablo Mesejo, Radu Horaud

► **To cite this version:**

Stéphane Lathuilière, Benoît Massé, Pablo Mesejo, Radu Horaud. Neural Network Reinforcement Learning for Audio-Visual Gaze Control in Human-Robot Interaction. 2017. hal-01643775v1

**HAL Id: hal-01643775**

**<https://inria.hal.science/hal-01643775v1>**

Preprint submitted on 21 Nov 2017 (v1), last revised 25 Apr 2018 (v2)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Neural Network Reinforcement Learning for Audio-Visual Gaze Control in Human-Robot Interaction

Stéphane Lathuilière<sup>1,1</sup>, Benoit Massé<sup>1,1</sup>, Pablo Mesejo<sup>1,1</sup>, Radu Horaud<sup>1,1</sup>

<sup>a</sup>*Inria Grenoble Rhône-Alpes, 655 Avenue de l'Europe, Montbonnot-Saint-Martin, 38330, France*

<sup>b</sup>*Univ. Grenoble Alpes, 621 avenue Centrale, Saint-Martin-d'Hères, 38400, France*

---

## Abstract

This paper introduces a novel neural network-based reinforcement learning approach for robot gaze control. Our approach enables a robot to learn and adapt its gaze control strategy for human-robot interaction without the use of external sensors or human supervision. The robot learns to focus its attention on groups of people from its own audio-visual experiences, and independently of the number of people in the environment, their position and physical appearance. In particular, we use recurrent neural networks and Q-learning to find an optimal action-selection policy, and we pretrain on a synthetic environment that simulates sound sources and moving participants to avoid the need of interacting with people for hours. Our experimental evaluation suggests that the proposed method is robust in terms of parameters configuration (i.e. the selection of the parameter values has not a decisive impact on the performance). The best results are obtained when audio and video information are jointly used, and when a late fusion strategy is employed (i.e. when both sources of information are separately processed and then fused). Successful experiments on a real environment with the Nao robot indicate that our framework is a step forward towards the autonomous learning of a perceivable and socially acceptable gaze behavior.

*Keywords:* Reinforcement Learning, Human-Robot Interaction, Robot Gaze Control, Neural Networks, Transfer Learning, Multimodal Data Fusion

---

## 1. Introduction

In recent years, there has been a growing interest in human-robot interaction (HRI), a research field dedicated to designing, evaluating and understanding robotic systems able to communicate with people [?]. The robotic agent must perceive humans and perform actions that, in turn, will have an impact on the interaction. For instance, it is known that the robot's verbal and gaze behavior has a strong effect on the turn-taking conduct of the participants [?]. Traditionally, HRI has been focused on the interaction between a single person with a robot. However, robots are increasingly part of groups and teams, e.g. performing delivery tasks in hospitals [?] or working closely alongside people on manufacturing floors [?]. In the case

of the gaze control problem in a multi-person scenario, the fact of focusing on only one person would lead to omit important information and, therefore, to make wrong decisions. Indeed, the robot needs to follow a strategy to maximize useful information, and such a strategy is difficult to design for two main reasons. First, the number of possible configurations grows exponentially with the sequence length, making difficult to obtain an optimal solution for long time sequences. Second, the robot needs to be able to adapt its strategy to currently available data, as provided by its sensors, cameras and microphones in our case. For instance, if a companion robot enters a room with very bad acoustic conditions, the strategy needs to be adapted by decreasing the importance given to audio information.

In this paper, we consider the general problem of gaze control, with the specific goal of finding good policies to control the orientation of a robot head during informal group gatherings. In particular, we propose a methodology for a robotic system to be able to autonomously learn to focus its attention on groups of people using audio-visual information. This is a very important topic of research since perception requires not only making inferences from observations, but also making decisions about where to look next. More specifically, we want a robot to learn to find people in the environment, hence maximize the number of people present in its field of view, and favor people who speak. We believe this could be useful in many real scenarios, such as a conversation between a humanoid robot and a group of humans, where the robot needs to learn to look at people, in order to behave properly. The reason for using multiple sources of information can be found in recent HRI research suggesting that no single sensor can reliably serve to robust interaction [? ]. Importantly, when it comes to the employment of multiple sensors in complex social interactions, it becomes difficult to implement an optimal policy based on handcrafted rules that take into consideration all possible situations that may occur. On the contrary, we propose to follow a data-driven approach to face such complexity. In particular, we propose to tackle this problem using a reinforcement learning (RL) approach [? ]. RL is a machine learning paradigm in which agents learn by themselves by trial-and-error to achieve successful strategies. As opposed to supervised learning, there is no need for optimal decisions at training time, only a way to evaluate how good a decision is: a reward. This paradigm, inspired by behaviorist psychology, can allow a robot to autonomously learn a policy that maximizes accumulated reward. In our case, the agent, a Nao robot, autonomously moves its head depending on its knowledge about the environment. This knowledge is called the agent’s state, and it is defined as a sequence of audio-visual and motor observations, actions and rewards. The optimal policy for making decisions is learned from the reward computed using the detected faces and the localized sound sources. The use of annotated data is not required to learn the best policy as the agent learns autonomously by trial-and-error in an unsupervised manner. Moreover, using our approach, it is not even necessary to make any assumption about the number of people moving in the environment or their initial

locations.

The use of RL techniques presents several advantages. First, training using optimal decisions is not required since the model learns from the reward obtained for each decision taken. The reward can be considered as a feedback signal that indicates how well the robot is doing at a given time step. Second, the robot must continuously make judgments so as to select good actions over bad ones. In this sense, the model can keep training at test time and benefits from a higher adaptation ability. Finally, we avoid the need of an annotated training set or calibration data, as our approach is unsupervised. In our opinion, it seems entirely natural to use RL techniques to “educate” a robot, since recent neuroscientific studies have suggested that reinforcement affects how infants interact with their environment, including what they look at [? ], and that face looking is not innate but that environmental importance influences viewing behavior.

This paper is an extension of a recently submitted conference paper [? ], where we presented the initial version of a neural network-based RL approach to address the robot gaze control problem. In this paper, we extend the system to deal with more complex and realistic scenarios, we delve into the impact of the main parameters of the model, and we make a much more precise description of the whole approach. The overall proposal of this paper can be summarized as follows:

- We propose to use recurrent neural networks, in combination with a fully connected layer, to autonomously learn the robot gaze control strategy by means of a value-based RL approach from multimodal data.
- We extend the visual observations to use a full body pose detector instead of a simple face detector. While this makes both the proposed simulation-based training and online algorithms more complex, it guarantees that more realistic scenarios are considered.
- We introduce a new algorithm to simulate moving persons together with their respective poses. We employ this synthetic environment to avoid the tedious training protocols that use real data, and we use

transfer learning from the simulated environment to the real one.

- We perform an extensive comparison with other neural network-based temporal architectures and evaluate the impact of the main parameters involved.
- We describe a real-time implementation of the proposed algorithm using a companion robot. The data, the network weights, and the codes used in this paper will be released upon acceptance of the paper <sup>1</sup>.

## 2. Related Work

RL has been studied for decades [? ? ] and has been widely used in various topics, including robotics [? ]. Learning a policy is the main challenge in RL, and there are two main categories of methods to address it. First, policy-based methods define a space from the set of policies, and sample policies from it. The reward is then used, together with optimization techniques, *e.g.* gradient-based methods, to increase the quality of subsequent sampled policies [? ]. The other category, value-based methods, consists in estimating the expected reward for the set of possible actions. The actual policy uses this value function to decide the suitable action, *e.g.* choose the action that maximizes the value-function. In particular, popular value-based methods include Q-learning [? ] and its deep learning extension, Deep Q-Networks (or DQNs) [? ].

We now review some of the most relevant RL-based HRI methods. In [? ] an RL approach is employed to learn a robot to play a game with a human partner. The robot uses vision and force/torque feedback to choose the commands, and the uncertainty associated with human actions is modeled via Gaussian processes. Bayesian optimization selects an optimal action at each time step. In [? ] RL is employed to adjust motion speed, timing, interaction distances, and gaze in the context of HRI. The

reward is based on the amount of movement of the subject and the time spent gazing at the robot in one interaction. As external cameras are required, this cannot be easily applied in scenarios where the robot has to keep learning in a real environment. Moreover, the method is limited to the case of a single human participant. Another example of RL applied to HRI can be found in [? ] where a human-provided reward is used to teach a robot. This idea of interactive RL is exploited by [? ] in the context of a table-cleaning robot. Visual and audio recognition are used to get advice from a parent-like trainer to enable the robot to learn a good policy efficiently. An extrinsic reward is used in [? ] to learn how to point a camera towards the active speaker in a conversation. Audio information is used to determine where to point the camera, while the reward is provided by visual information: the active speaker raises a blue card that can be easily identified by the robot. The use of a multimodal DQN to learn human-like interactions is proposed both in [? ] and in [? ]. The robot must choose an action (waiting, looking to a person, hand waving and hand shaking) to perform a hand shake with a human. The reward is negative if the robot tries unsuccessfully to shake hands, positive if the hand shake is successful, and null otherwise. In practice, the reward is obtained from a sensor located in the hand of the robot and it takes fourteen days of training to learn this skill successfully. To the best of our knowledge, the closest work to ours is [? ] where an RL approach learns good policies to control the orientation of a mobile robot during social group conversations. The robot learns to turn its head towards the speaking person of the group. However, their model is learned on simulated data that are restricted to a few predefined scenarios with static people and a fixed spatial organization of the group.

In contrast to prior work, our approach allows a robot to autonomously learn an effective gaze control policy from audio and visual inputs in an unconstrained real-world environment. In particular, a simulated environment is used for pretraining, thus avoiding to spend several days of tedious real interactions with people, followed by transfer learning to map the learned strategies to real environments. Moreover, it requires neither external sensors nor human intervention to obtain a reward.

---

<sup>1</sup>A video showing additional offline and online experiments is already available at <https://team.inria.fr/perception/research/neural-reinforcement-learning-for-human-robot-interaction>

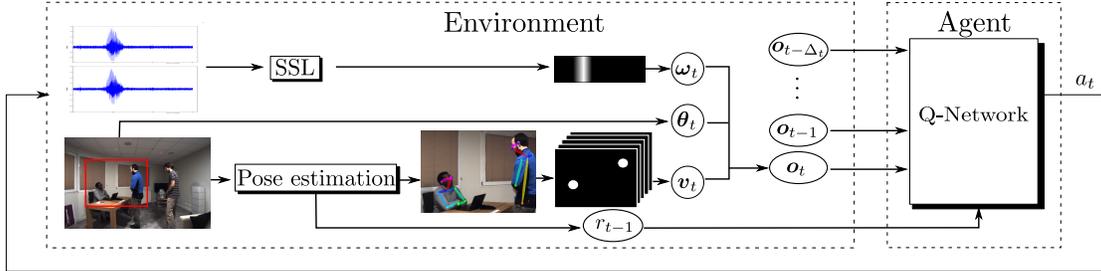


Figure 1: Outline of our neural network-based reinforcement learning (RL) approach. In this example, the reward,  $r_{t-1}$ , includes only visual information (*Face\_reward*).  $\omega_t$ ,  $\theta_t$ , and  $v_t$  denote the extracted audio (Sound Source Localization, or SSL), motor and video features that compose  $\mathbf{o}_t$ , a particular observation at time  $t$ . The observation at time  $t$  and the previous  $\Delta_t$  observations are used as input for the Q-Network that, as output, provides the action,  $a_t$ , with the largest expected return.

### 3. Reinforcement Learning for Gaze Control

#### 3.1. Problem Formulation

We consider the case of a robot in a social environment. The robot can move its head using its motors (2 degrees of freedom) and have access to video (stereo camera) and audio information (microphone array). We assume the existence of methods to perform body pose estimation and sound source localization (SSL) to process video and audio data, respectively. The goal is to “teach” the robot, by a trial-and-error learning procedure, to perform the suitable action employing those sources of information. In this case, the suitable action corresponds to move the head to maximize the number of people in the field of view, favoring the attention to people who speak. The main building blocks of our proposal are graphically displayed in Figure ???. Throughout the paper, random variables and their realizations are denoted by uppercase and lowercase letters, respectively. Vectors and matrices are represented using bold italic letters. The terms agent and robot are used indistinctly to refer to the autonomous entity that observes the environment and performs actions towards achieving a particular goal.

The variable  $\mathbf{O}_t$  gathers observations received at time step  $t$  by the agent, namely visual ( $\mathbf{V}_t$ ), audio ( $\mathbf{\Omega}_t$ ), and motor observations ( $\mathbf{\Theta}_t$ ). The agent uses  $\mathbf{O}_t$  to update its state from  $\mathcal{S}_{t-1}$  to  $\mathcal{S}_t$ , that represents the knowledge the agent possesses about the environment. The agent then selects an action,  $A_t$ , that consists on sending an input to the robot motors. Finally, the agent receives a reward,  $R_t$ ,

that determines the suitability of the action to carry out the specified task.

The current pitch and yaw angles of the robot head are denoted by  $\mathbf{\Theta}_t = (\Theta_t^p, \Theta_t^y)$ .  $\mathbf{V}_t \in \{0, 1\}^{N_v \times M_v \times J}$  denotes the visual observations. We consider an image of size  $N_v \times M_v$  and a multi-person 2D pose estimator returning the  $J$  joint locations of the  $N_t$  detected persons at each time  $t$ . The outputs of the pose estimator are  $p_n = (u_n^j, v_n^j, s_n^j)$ ,  $n \in [1..N_t]$ ,  $j \in [1..J]$  where  $(u_n^j, v_n^j) \in N_v \times M_v$  denotes the pixel coordinates of the  $j^{\text{th}}$  joints, and  $s_n^j$  is a binary value such that  $s_n^j = 1$  if the  $j^{\text{th}}$  joints is successfully detected and  $s_n^j = 0$  otherwise. We can now define the visual observations  $\mathbf{V}_t$  as follows:

$$\mathbf{V}_t(u, v, j) = \begin{cases} 1 & \text{if } \exists n \in [1..N_t], s_n^j = 1, (u_n^j, v_n^j) = (u, v) \\ 0 & \text{otherwise} \end{cases}$$

Similarly,  $\mathbf{\Omega}_t \in [0, 1]^{N_a \times M_a}$  denotes the audio observations, where  $\mathbf{\Omega}_t$  is a heat map giving the probability that a source is emitting a sound at each location of a grid of size  $N_a \times M_a$ . It is important to notice that in practice the audio grid is wider than the visual grid as a sound source can be detected even if the source is outside the field of view of the cameras. The observation variable is formally defined as  $\mathbf{O}_t = \{\mathbf{\Theta}_t, \mathbf{V}_t, \mathbf{\Omega}_t\}$ , and the state variable is defined as  $\mathcal{S}_t = \{\mathbf{O}_1 \dots \mathbf{O}_t\} \in \mathcal{S}$ , the set of all possible states. The robot can perform the action  $A_t \in \mathcal{A} = \{\emptyset, \leftarrow, \uparrow, \rightarrow, \downarrow\}$ , namely staying in the same position or moving its head a fixed angle in one of the four cardinal directions. The reward  $R_t$  is defined after taking action  $A_t$ , either as the number of faces detected in the field of view (termed *Face\_reward*

in section ??, and displayed in Figure ??) or as the number of faces detected plus one if sound is also present in the field of view (termed *Speaker\_reward* in section ??). We consider interesting to compare a purely visual reward with a multimodal one including audio information.

In RL, the model is learned on sequences of  $T$  states, actions and rewards called episodes. At each time-step  $t$ , the action  $a_t$  should not be chosen aiming at maximizing only the immediate reward  $R_t$  but also the future rewards ( $R_{t+1} \dots R_T$ ). To do so, we make the standard assumption that future rewards are discounted by a factor of  $\gamma$ . The parameter  $\gamma$  defines how much we favor rewards returned in the next coming time-steps over longer term rewards. We then define the discounted future return  $\bar{R}_t$  as the discounted sum of future rewards  $\bar{R}_t = \sum_{\tau=0}^{T-t} \gamma^\tau R_{t+\tau}$ . We now aim at maximizing  $\bar{R}_t$  at each time step  $t$ . In other words, the goal is to learn a policy,  $\pi(a_t, s_t) = P(A_t = a_t | S_t = s_t)$ ,  $(a_t, s_t) \in \mathcal{A} \times \mathcal{S}$ , such that if the agent choose its actions according to  $\pi$ , the expected  $\bar{R}_t$  is maximal. The Q-function (also called action-value function) is defined as the expected future return from state  $S_t$  taking action  $A_t$  and then following any policy  $\pi$ :

$$Q_\pi(s_t, a_t) = \mathbb{E}_\pi(\bar{R}_t | S_t = s_t, A_t = a_t) \quad (1)$$

Learning the best policy corresponds to the following optimization problem  $Q^*(s, a) = \max_\pi [Q_\pi(s_t = s, a_t = a)]$ . The optimal Q-function obeys the identity known as the Bellman equation:

$$Q^*(s_t, a_t) = \mathbb{E}_{S_{t+1}, R_t} \left[ R_t + \gamma \max_a (Q^*(S_{t+1}, a)) \middle| S_t = s_t, A_t = a_t \right] \quad (2)$$

This equation corresponds to the following intuition: if we have an estimator  $Q^*(s_t, a_t)$  for  $\bar{R}_t$ , the optimal action  $a_t$  is the one that leads to the largest expected  $\bar{R}_t$ . The recursive application of this policy leads to equation (??). A straightforward approach would consist in updating  $Q$  at each training step  $i$  with:

$$Q^i(s_t, a_t) = \mathbb{E}_{S_{t+1}, R_t} \left[ R_t + \gamma \max_a (Q^{i-1}(S_{t+1}, a)) \middle| S_t = s_t, A_t = a_t \right] \quad (3)$$

However, in practice, we employ a network  $Q(s, a, \theta)$  parametrized by weights  $\theta$  to estimate the Q-function  $Q(s, a, \theta) \approx Q^*(s, a)$  and we minimize the following loss:

$$\mathcal{L}(\theta_i) = \mathbb{E}_{S_t, A_t, R_t, S_{t+1}} \left[ (Y_{i-1} - Q(S_t, A_t, \theta_i))^2 \right] \quad (4)$$

with  $Y_{i-1} = R_t + \gamma \max_a (Q(S_{t+1}, a, \theta_{i-1}))$ . It can be seen as minimizing the mean square distance between the approximations of the right and left hand sides of (??). When the robot is training, we obtain a quadruplet  $(S_t, A_t, R_t, S_{t+1})$  for each time-step allowing us to compute (??). However, instead of sampling only according to the policy implied by  $Q(s, a, \theta_i)$ , random actions  $a_t$  are taken in  $\epsilon$  percents of the time steps in order to explore new strategies. This approach is known as epsilon-greedy policy.  $\mathcal{L}$  is minimized over  $\theta_{i+1}$  by stochastic gradient descent. Refer to [?] for more technical details about the training algorithm.

### 3.2. Neural Network Architectures for Q-Learning

The Q-function is modeled by a neural network that takes as input part of the state variable  $S_t$ , that we define as  $S_t^{\Delta t} = \{\mathbf{O}_{t-\Delta t} \dots \mathbf{O}_t\}$ . The output is a vector of size  $\#\mathcal{A}$  that corresponds to each  $Q_\pi(s_t^{\Delta t}, a_t)$ ,  $a_t \in \mathcal{A}$ , where  $Q_\pi(s_t^{\Delta t}, a_t)$  is built analogously to Equation ???. Following [?], the output layer is a Fully-Connected Layer (FCL) with linear activations. We propose to use the Long Short-Term Memory (LSTM) [?] recurrent neural network to model the Q-function. Batch normalization is applied to the output of the LSTM. We argue that LSTM is well-suited for our task as it is capable of learning temporal dependencies better than other recurrent neural networks and hidden Markov models. In fact, our model needs to memorize the position and the motion of the people when it turns its head. When a person is not detected anymore, the network should be able to use previous detections back in time in order to predict the direction towards it should move. The  $J$  channels of  $V_t$  are flattened before the LSTM layers.

Four different network architectures are described in this section and evaluated in the experimental section. In order to evaluate when the two streams of information (audio and video) need to be fused, we propose to compare two strategies: early fusion and late fusion. In

early fusion, the unimodal features are combined into a single representation before modeling time dependencies (see Figure ??, called *EFNet*). Conversely, in late fusion, audio-visual features are modeled separately before fusing them (see Figure ??, called *LFNet*). In order to measure the impact of each modality, we propose two more networks using either only video ( $V$ ) or only audio ( $\Omega$ ) information. Figure ?? displays *AudNet*, the network using only audio information, while Figure ?? shows *VisNet*, that employs only visual information. Figure ?? employs a compact representation where time  $t$  is not explicitly included, while Figure ?? depicts the unfolded representation of *EFNet* where each node is associated with one particular time instance. Both figures follow the graphical representation used in [? ].

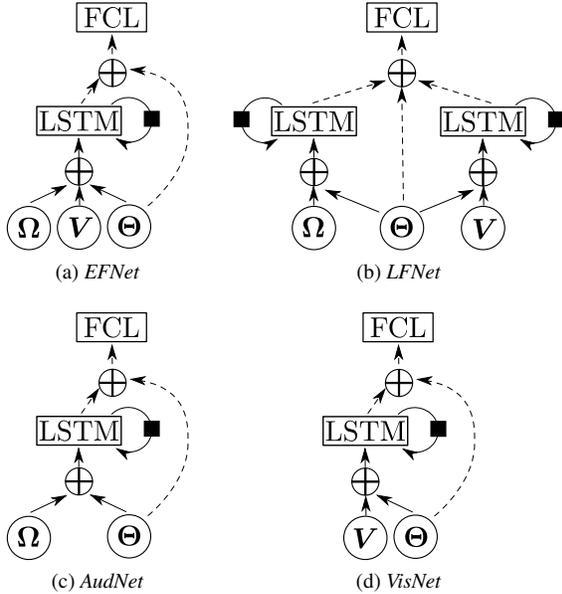


Figure 2: Proposed architectures to model the Q-function. Dashed lines indicate connections only used in the last time step. Black squares represent a delay of a single time step. Encircled crosses depict the concatenation of inputs.

### 3.3. Pretraining on Synthetic Environment

Training from scratch a DQN model can take a long time (in our case  $\sim 150000$  time steps to converge), and training directly on a robot would not be convenient for

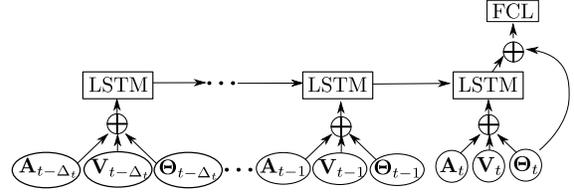


Figure 3: Unfolded representation of *EFNet* to better capture the sequential nature of the recurrent model. Encircled crosses depict the concatenation of inputs.

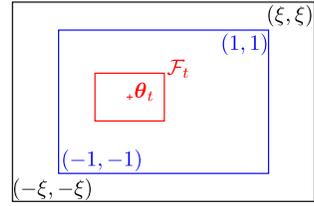


Figure 4: Diagram showing all fields used in the proposed synthetic environment. The robot's field of view (in red) can move within the reachable field (in blue), whereas the participants can freely move within a larger field (in black).

two reasons. First, it would entail a long period of training, since each physical action by the robot takes an amount of time that cannot be reduced neither by code optimization nor by increasing our computational capabilities. Second, in the case of HRI, participants would need to move in front of the robot for several hours or days (like in [? ]). For these two reasons, we propose to use a transfer learning approach. The Q-function is first learned on a synthetic environment, where we simulate people moving and talking, and it is then used to initialize the network employed by the robot. Importantly, the network learned from this synthetic environment can be successfully used in the robot without the need of fine-tuning in real data. In this synthetic environment, we do not need to generate images and sound signals, but only the observations and rewards the Q-Network receives as input.

We consider that the robot can cover the field  $[-1, 1]^2$  by moving its head, but can only visually observe the people within a small rectangular region  $\mathcal{F}_t \subset [-1, 1]^2$  centered in position vector  $\Theta_t$ . The audio observations cover the whole reachable region  $[-1, 1]^2$ . However, the actual robot we use is only able to locate the yaw angle of the sound sources, therefore we decided to solely provide

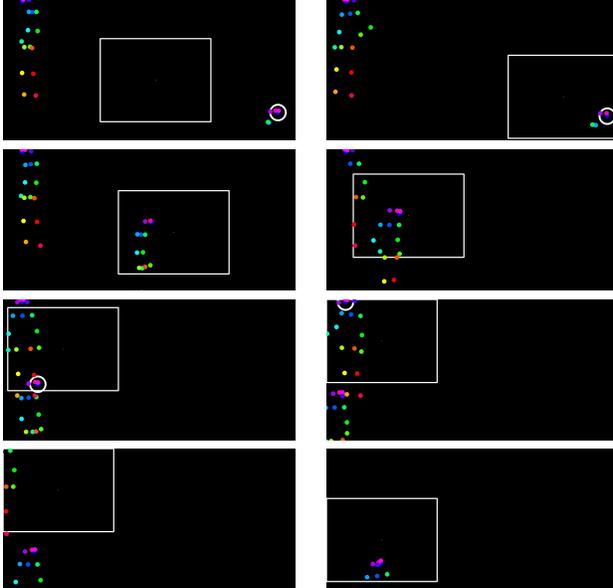


Figure 5: Illustrative sequence taken from the synthetic environment and employed to pretrain our neural network-based RL approach. The moving square represents the camera field of view  $\mathcal{F}_t$  of the robot. The colored circles represent the joints of a person in the environment. The large white circle represents a person speaking and, therefore, producing speech that can be detected by the SSL system. Frames are displayed from top to bottom and left to right.

sound observations on the horizontal axis  $[-1, 1]$ . On each episode, we simulate one or two persons moving with random speeds and accelerations within a field  $[-\xi, \xi]^2$  where  $\xi > 1$ . In other words, people can go to regions that are unreachable for the robot. For each simulated person in the current episode, we consider the position and velocity of their head at time  $t$ ,  $\mathbf{h}_t = (u_t^h, v_t^h) \in [-\xi, \xi]^2$  and  $\dot{\mathbf{h}} = (\dot{u}_t^h, \dot{v}_t^h) \in \mathbb{R}^2$ , respectively. At each frame, the person can keep moving, stay without moving, or choose another random direction. The details of the simulated environment generator are given in Algorithm ???. In a real scenario, people can leave the scene so, in order to simulate this phenomenon, we consider two equally probable cases when a person is going out horizontally of the field ( $v_t^h \notin [-\xi, \xi]$ ). In the first case, the person is deleted and instantly recreated on the other side of the field ( $v_{t+1}^h = -v_t^h$ ) keeping the same velocity ( $\dot{v}_{t+1}^h = \dot{v}_t^h$ ). In the second case, the person is going back towards the center ( $v_{t+1}^h = v_t^h$  and ( $\dot{v}_{t+1}^h = -\dot{v}_t^h$ )). A similar approach is

used when a person is going out vertically except that we do not create new persons on top of the field because that would imply the unrealistic sudden appearance of new legs within the field. Figure ?? displays a visual representation of the different fields (or areas) defined in our synthetic environment, and Figure ?? shows an example of a sequence of frames taken from the synthetic environment and used during training.

Moreover, in order to favor tracking abilities, we bias the person motion probabilities such that a person that is faraway from the robot head orientation has a low probability to move, and a person within the field of view has a high probability to move. Thus, when there is nobody in the field of view, the robot cannot simply wait for a person to come in. On the contrary, the robot needs to track the persons that are visible. More precisely, we consider 4 different cases. First, when a person has never been seen by the robot, the person does not move. Second, when a person is in the robot field of view ( $\mathbf{h}_t \in \mathcal{F}_t$ ), they move with a probability of 95%. Third, when the person is further than a threshold  $\tau \in \mathbb{R}$  from the field of view ( $\|\mathbf{h}_t - \mathbf{O}_t\|_2 > \tau$ ), the probability of moving is only 25%. Finally, when the person is not visible but close to the field of view ( $\|\mathbf{h}_t - \mathbf{O}_t\|_2 < \tau$  and  $\mathbf{h}_t \notin \mathcal{F}_t$ ), or when the person is unreachable ( $\mathbf{h}_t \in [-\xi, \xi] \setminus [-1, 1]$ ), this probability is 85%. Regarding the simulation of missing detections, we randomly ignore some faces when computing the face features. Concerning the sound modality, we randomly choose between the following cases: 1 person speaking, 2 persons speaking, and nobody speaking. We use a Markov model to enforce continuity in the speaking status of the persons, and we also simulate wrong SSL observations.

From, the head position, we need to generate the position of all body joints. To do so, we propose to collect a set  $\mathcal{P}$  of poses from an external dataset (the AVDIAR dataset [? ]). We use a multiple person pose estimator on this dataset and use the detected poses for our synthetic environment. This task is not trivial since we need to simulate a realistic and consistent sequence of poses. Applying tracking to the AVDIAR videos could provide good pose sequences, but we would suffer from three major drawbacks. First, we would have a tracking error that could affect the quality of the generated sequences. Second, each sequence would have a different

and constant size, whereas we would like to simulate sequences without size constraints. Finally, the number of sequences would be relatively limited. In order to tackle these three concerns, we first standardize the output coordinates obtained on AVDIAR. Considering the pose  $p_t^n$  of the  $n^{\text{th}}$  person, we sample a subset  $\mathcal{P}_t^M \subset \mathcal{P}$  of  $M$  poses. Then, we select the closest pose to the current pose:  $p_{t+1}^n = \underset{p \in \Pi}{\operatorname{argmin}} d(p, p_t^n)$  where

$$d\left(\begin{pmatrix} u_1 \\ v_1 \\ s_1 \end{pmatrix}, \begin{pmatrix} u_2 \\ v_2 \\ s_2 \end{pmatrix}\right) = \frac{1}{\sum_{j=1}^J s_1^j s_2^j} \sum_{j=1}^J (s_1^j s_2^j) \sqrt{(u_1^j - u_2^j)^2 + (v_1^j - v_2^j)^2} \quad (5)$$

This distance is designed to face poses with different number of detected joints. It can be interpreted as an  $L_2$  distance weighted by the number of visible joints in common. The intuition behind this sampling process is that when the size  $M$  of  $\mathcal{P}_t^M$  increases, the probability of obtaining a pose closer to  $p_t^n$  increases. Consequently, the motion variability can be adjusted with the parameter  $M$  in order to obtain a natural motion. With this method we can obtain diverse sequences of any size.

## 4. Experiments

This section begins with the description of the quantitative evaluation performed on AVDIAR and synthetic datasets. After this offline evaluation, it continues with the description of the experiments in real time with the Nao robot, performed to qualitatively evaluate our approach in a real environment. Finally, the section ends with implementation details, and the results obtained and their analysis.

### 4.1. Offline Evaluation on the AVDIAR Dataset

The evaluation of HRI systems is not an easy task. First, the definition of a metric to measure a correct, socially acceptable behavior is far from trivial [? ]. In our particular case, since gaze behavior is an important nonverbal communication cue in human-human social encounters [? ? ], we evaluate our approach according to the robot capability of finding and tracking faces. Second, in

**Data:**  $\mathcal{P}$ : a set of poses,  $\delta$ : time-step  
 $\sigma$ : velocity variance,  $M$ : pose continuity parameter

```

Randomly chose  $N$  in [1..3].
for  $n \in [1..N]$  do
  Initialize
   $(\mathbf{h}_0^n, \dot{\mathbf{h}}_0^n) \sim \mathcal{U}([-1, 1])^2 \times \mathcal{U}([-1, -0.5] \cup [0.5, 1])^2$ .
  Randomly chose  $p_0^n$  in  $\mathcal{P}$ .
end
for  $t \in [1..T - 1]$  do
  for  $n \in [1..N]$  do
    Randomly chose  $motion \in \{Stay, Move\}$ 
    if  $motion = Move$  then
      if  $\mathbf{h}_t^n \notin [-\xi, \xi]^2$  then
        The person is leaving the scene.
        See section ???.
      else
         $\mathbf{h}_{t+1}^n \leftarrow \mathbf{h}_t^n + \delta(\dot{\mathbf{h}}_t^n + \mathcal{N}((0, 0), \sigma))$ .
         $\dot{\mathbf{h}}_{t+1}^n \leftarrow \frac{1}{\delta}(\mathbf{h}_{t+1}^n - \mathbf{h}_t^n)$ 
      end
    else
       $\mathbf{h}_{t+1}^n \leftarrow \mathbf{h}_t^n$ 
       $\dot{\mathbf{h}}_{t+1}^n \sim \mathcal{U}([-1, -0.5] \cup [0.5, 1])^2$ 
    end
    Draw  $\mathcal{P}_t^M$ , a random set of  $M$  elements of  $\mathcal{P}$ 
     $p_{t+1}^n \leftarrow \underset{p \in \mathcal{P}_t^M}{\operatorname{argmin}} d(p, p_t^n)$ 
  end
end

```

**Algorithm 1:** Generation of simulated moving poses for our synthetic environment.

order to fairly compare different models, we need to train and test the different models on the exact same data. In the context of RL in HRI, this is problematic because the data (e.g. what the robot sees and hears) depends on the action the robot has taken. Thus, we propose to first evaluate our proposal on an offline dataset. To mimic the real behavior of a robot, we use the audio-visual AVDIAR dataset [? ]. This dataset has been recorded with 4 microphones and high-resolution binocular cameras ( $1920 \times 1080$ ), of which we use only one. These images, due to their wide field of view, are suitable to simulate the whole field the robot can cover by moving its head. In practical terms, only a small part of the full image is considered as seen

by the robot.

#### 4.2. Real Time Experiments on a Nao Robot

In order to carry out an online evaluation of our proposal, we perform experiments on a Nao robot, developed by Aldebaran Robotics. Nao provides a camera of  $640 \times 480$  pixels and four microphones. This robot is particularly well suited for HRI applications because of its design, hardware specifications and affordable cost. Nao can detect and identify people, localize sounds, understand some spoken words, synthesize speech and engage itself in simple and goal-directed dialogs. Our gaze control system is implemented on top of the NAOLab middleware [?] that synchronizes proprioceptive and perceptive information. The reason why we use a middleware is three-fold. First, the implementation is platform-independent and, thus, easily portable. Platform-independence is crucial since we employ a transfer learning approach to transfer the knowledge gathered on our proposed synthetic environment to the Nao robot. Second, the use of external computational resources is transparent. This is also a crucial matter in our case, since the full-body pose estimator requires GPU computation for fast inference. Third, the use of middleware makes prototyping much faster. For all these reasons, we employed the remote and modular layer-based middleware architecture named NAOLab. NAOLab consists of 4 layers: drivers, shared memory, synchronization engine and application programming interface (API). Each layer is divided into 3 modules devoted to vision, audio and proprioception, respectively. The last layer of NAOLab provides a general programming interface in C++ to handle the robot’s sensor data and manage its actuators. NAOLab provides, at each time step, the camera images and the yaw angle of the detected sound sources using [? ?].

It is important to highlight that we pretrain on the synthetic environment before running experiments on the Nao robot. The synthetic environment is flexible and allows us to be closer to the conditions Nao would face in reality (field of view range, uniform location of the people of the field, etc.). For instance, in AVDIAR, as the camera is fixed, heads are almost always at the same height. As a consequence, the learned model would not be sufficiently general and flexible to perform well in real scenar-

ios. Figure ?? shows a synthetic sequence employed for pretraining our neural network-based RL system.

#### 4.3. Implementation Details

In all experiments we employ the full-body pose estimator described in [? ], considering the nose as the landmark that represents the face. On the Nao robot, we manage to obtain the pose in less than 100ms by selecting carefully the research scale and downsampling the images. Considering that the Nao cameras provide images with 10 fps, this pose estimator method can be considered as fast enough for our scenario. Moreover, [? ] has the particularity of following a bottom-up approach: each body joint is first detected in the image, and then connected by solving a graph matching problem. In our case, as we use a joint heatmap, we do not need to perform this association step in order to save computation time.

In all scenarios we set  $\Delta_T = 4$  such that each decision is based on the last 5 observations. Different values for  $\Delta_T$  were tested, see Table ??, and we kept the value that provided the best possible results without increasing the computational complexity (in fact, values for  $\Delta_T$  larger than 1 provided an almost equivalent final reward). The output size of the LSTMs is set to 30 (since larger sizes do not provide an improvement in performance, see Table ??), and the output size of the FCLs is set to 5 (one per action). We use a discount factor ( $\gamma$ ) of 0.90 (that yields a good performance in both the AVDIAR test and the synthetic environment, see Table ??). Concerning the training phases, we employed the Adam optimizer [? ] and a batch size of 128. In order to help the model to explore the policy space, we use an  $\epsilon$ -greedy algorithm: while training, a random action is chosen in  $\epsilon\%$  of the cases; we decrease linearly the  $\epsilon$  value from  $\epsilon = 90\%$  to  $\epsilon = 10\%$  after 150000 iterations. Concerning the observations, we employed visual and SSL heatmaps of sizes  $7 \times 5$  for the three environments used in our experiments. The models were trained in approximately 45 minutes on both AVDIAR and the synthetic environment. It is interesting to notice that we obtain this training time without using GPU, because a GPU is only required to compute the full-body pose (in our case, a Nvidia GTX 1070 GPU).

Concerning the details related specifically to the AVDIAR dataset, we employed 16 videos for training. The

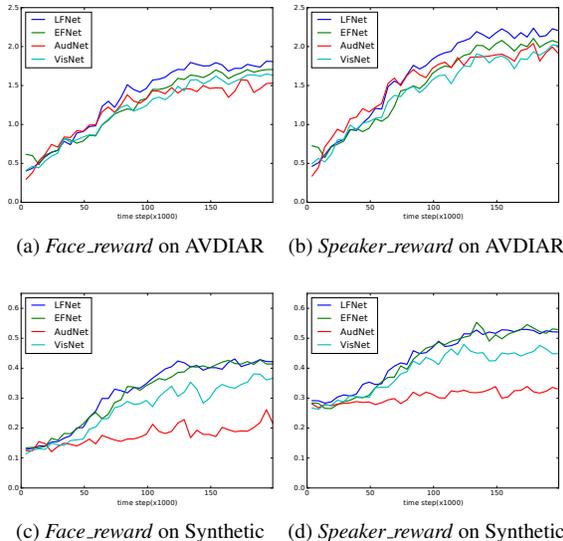


Figure 6: Evolution of the reward obtained while training with the two proposed rewards on the AVDIAR dataset and on the synthetic environment. We average over a 5000 time-step window for a cleaner visualization.

amount of training data is doubled by flipping the video and the SSL outputs. In order to save computation time, the original videos are down-sampled to  $1024 \times 640$  pixels. The size of the field of view where faces can be detected is set to  $300 \times 200$  pixels using motion steps of 36 pixels each. At the beginning of each episode, the position of the field of view is randomly sampled such that no face can be seen. We noticed that this initialization procedure favors the exploration abilities of the agent. To avoid a bias due to the initialization procedure, we used the same seed for all our experiments and iterated 3 times over the 10 test videos (20 when counting the flipped sequences). An action is taken every 5 frames (0.2 seconds) and the SSL is obtained using [?]. In the synthetic environment, the size of field in which the people can move is set to  $\xi = 1.4$ . In the case of Nao, the delay between two successive observations is  $\sim 0.3$  seconds. The head is free to move in a field corresponding to 180 degrees. The motion of a single action corresponds to 0.15 radians ( $\sim 9^\circ$ ) and 0.10 radians ( $\sim 6^\circ$ ) for horizontal and vertical moves, respectively.

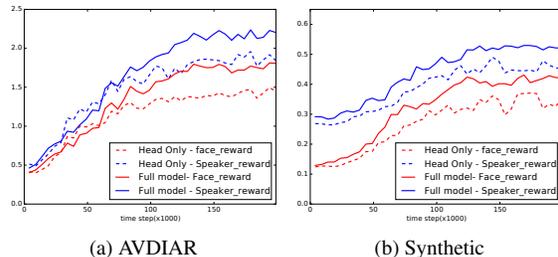


Figure 7: Evolution of the training reward obtained when using as visual observation the result of either the full-body pose estimation or the face detection.

#### 4.4. Results and Discussion

In all our experiments, we run five times each model and display the mean of five runs to lower the impact of the stochastic training procedure. On AVDIAR, the results on both training and test sets are reported in the tables. As described previously, the synthetic environment is randomly generated in real time, so there is no need for a separated test set. Consequently, the mean reward over the last 10000 time steps is reported.

Table 1: Comparison of the final reward obtained using different window lengths ( $\Delta_T$ ). The mean and standard deviation over 5 runs are reported. The best average results obtained are displayed in bold.

$\Delta_T + 1$	AVDIAR		Synthetic
	Training	Test	
1	$1.75 \pm 0.04$	$1.55 \pm 0.08$	$0.26 \pm 0.04$
2	$1.80 \pm 0.02$	$1.59 \pm 0.03$	$0.36 \pm 0.04$
3	<b><math>1.82 \pm 0.03</math></b>	<b><math>1.61 \pm 0.01</math></b>	$0.42 \pm 0.02$
5	$1.81 \pm 0.01$	$1.57 \pm 0.03$	<b><math>0.43 \pm 0.01</math></b>
10	$1.80 \pm 0.01$	$1.60 \pm 0.03$	$0.40 \pm 0.02$
20	$1.80 \pm 0.03$	$1.55 \pm 0.05$	$0.42 \pm 0.02$

First, we describe the experiments devoted to evaluate the impact of some of the principal parameters involved. Different window sizes (i.e. the number of past observations necessary to make a decision) are compared in Table ???. We can conclude that the worst results are obtained when only the current observation is used (window size of 1). We also observe that, on AVDIAR, the model performs well even with short window lengths (2 and 3). In turn, with a more complex environment, as the proposed

synthetic environment, a longer window length tends to perform better. We interpret that using a larger window size helps the network to ignore the noisy observations and to remember the position of people that left the field of view. In Table ??, different discount factors are compared. We notice that, on AVDIAR, high discount factors are prone to overfit as the difference in performance between training and test is higher. On the synthetic environment, low discount values perform worse because we think that, as the environment is more complex, the model may need several actions to reach a face. Consequently, a model that is able to take into account the future benefit of each action performs better. Finally, in Table ??, we compare different LSTM sizes. We observe that increasing the size does not lead to better results; an interesting conclusion since, from a practical point of view, smaller LSTMs are faster to train.

Table 2: Comparison of the final reward obtained using different discounted factors ( $\gamma$ ). The mean and standard deviation over 5 runs are reported. The best average results obtained are displayed in bold.

$\gamma$	AVDIAR		<i>Synthetic</i>
	Training	Test	
0.25	1.80 $\pm$ 0.47	<b>1.60 <math>\pm</math> 0.01</b>	0.33 $\pm$ 0.09
0.50	1.80 $\pm$ 0.48	<b>1.60 <math>\pm</math> 0.01</b>	0.35 $\pm$ 0.08
0.75	1.80 $\pm$ 0.49	1.54 $\pm$ 0.08	<b>0.43 <math>\pm</math> 0.11</b>
0.90	1.80 $\pm$ 0.48	1.58 $\pm$ 0.03	0.42 $\pm$ 0.12
0.99	<b>1.81 <math>\pm</math> 0.49</b>	1.55 $\pm$ 0.04	0.42 $\pm$ 0.12

Table 3: Comparison of the final reward obtained using different LSTM sizes. The mean and standard deviation over 5 runs are reported. The best average results obtained are displayed in bold.

LSTM size	AVDIAR		<i>Synthetic</i>
	Training	Test	
30	<b>1.81 <math>\pm</math> 0.47</b>	1.56 $\pm$ 0.02	0.42 $\pm$ 0.11
60	<b>1.81 <math>\pm</math> 0.50</b>	<b>1.60 <math>\pm</math> 0.03</b>	<b>0.43 <math>\pm</math> 0.12</b>
120	1.79 $\pm$ 0.46	1.56 $\pm$ 0.02	0.41 $\pm$ 0.10

In Figure ??, we compare the evolution of the reward obtained while training on the AVDIAR dataset and on our synthetic environment with the two proposed rewards (*Face\_reward* and *Speaker\_reward*). Four different networks are tested: *EFNet*, *LFNet*, *VisNet*, and *AudNet*. The y-axis of Figure ?? shows the average reward per episode, with a clear growing trend as the training time

passes (specially in the experiments with the AVDIAR dataset), meaning that the agent is learning (improving performance) from experience. The best results are generally provided by the late fusion strategy (*LFNet*) and the *Speaker\_reward*. We observe that the rewards we obtain on AVDIAR are generally higher than those obtained on the synthetic environment. We suggest two possible reasons. First, the synthetic environment, as described in section ??, has been specifically designed to enforce exploration and tracking abilities. Consequently, it poses a more difficult problem to solve. Second, the number of people in AVDIAR is higher (about 4 in average), thus finding a first person to track would be easier.

Figure ?? displays the reward obtained when using only faces as visual observation (dashed lines) in contrast to using the full-body pose estimation (continuous lines). The former represents the results obtained by our previous proposal [? ]. We observe that for both datasets, the rewards are significantly higher when using the full-body pose estimator. This figure intends to respond empirically to the legitimate question of why a full-body pose estimator is used instead of a simple face detector. From a qualitative point of view, the answer can be found in the type of situations that can solve one and the other. Let’s imagine that the robot looks at the legs of a user; in case of using only a face detector, there is no clue that could help the robot to move up its head in order to see a face; however, if a human full-body pose detector is used, the detection of legs implies that there is a torso over them, and a head over the torso. Figure ?? shows a short sequence of the AVDIAR environment, displaying the whole field covered by the AVDIAR videos as well as the smaller field of view captured by the robot (the red rectangle in the figure).

Finally, Table ?? shows the mean reward on the test set for all architectures and rewards, using both AVDIAR and synthetic data. We can notice that, on the AVDIAR dataset using the *Face\_reward*, we obtain a mean reward greater than 1, meaning that, on average, our model can see more than one face per frame. Similarly to Figure ??, higher rewards are obtained in the AVDIAR dataset, and the best results are yielded when both modalities are taken into account with *LFNet*. That led us to select the *LFNet* model to perform experiments on Nao. We observe that *AudNet* is the worst performing approach. However, it performs quite well on AVDIAR compared to the syn-

Table 4: Comparison of the reward obtained with different architectures. The best results obtained are displayed in bold.

Network	AVDIAR				Synthetic	
	Face		Speaker		Face	Speaker
	Training	Test	Training	Test		
<i>AudNet</i>	1.53 ± 0.02	1.47 ± 0.02	1.91 ± 0.03	1.84 ± 0.02	0.21 ± 0.01	0.33 ± 0.01
<i>VisNet</i>	1.63 ± 0.03	1.54 ± 0.04	2.01 ± 0.02	1.87 ± 0.05	0.37 ± 0.04	0.45 ± 0.06
<i>EFNet</i>	1.71 ± 0.04	1.53 ± 0.06	2.05 ± 0.02	1.85 ± 0.09	0.41 ± 0.03	<b>0.53 ± 0.03</b>
<i>LFNet</i>	<b>1.81 ± 0.01</b>	<b>1.56 ± 0.04</b>	<b>2.20 ± 0.04</b>	<b>1.96 ± 0.05</b>	<b>0.42 ± 0.01</b>	0.52 ± 0.03

thetic environment. This behavior can be explained by the fact that, on *AVDIAR*, the SSL algorithm returns a 2D heatmap whereas only the yaw angle is used in the synthetic environment.

Concerning the experiments performed on Nao, Figure ?? shows an example of a two-person scenario using the *LFNet* architecture. We managed to transfer the exploration and tracking abilities learned using the synthetic environment. In our experiments, we see that our model behaves well independently of the number of participants, and the main failure cases are related to quick movements of the participants.

## 5. Conclusions

In this paper we have presented a neural network-based reinforcement learning approach to solve the gaze robot control problem. In particular, our agent is able to autonomously learn how to find people in the environment by maximizing the number of people present in its field of view (and favoring people who speak). A synthetic environment is used for pretraining in order to perform transfer learning to the real environment. Neither external sensors nor human intervention are necessary to compute the reward. Several architectures and rewards are compared on three different environments: two offline (a real and a synthetic datasets) and one online (real time experiments using the Nao robot). Our results suggest that the late fusion of audio and visual information represents the best performing alternative, as well as that pretraining on synthetic data can even make unnecessary to train on real data.

## Acknowledgments

Funding from the EU through the ERC Advanced Grant VHIA #340113 is greatly acknowledged.

## References

- [ ] Arcaro, M.J., Schade, P.F., Vincent, J.L., Ponce, C.R., Livingstone, M.S., 2017. Seeing faces is necessary for face-domain formation. *Nature Neuroscience* 20, 1404–1412.
- [ ] Argyle, M., 1975. *Bodily communication*. 1st ed., Routledge.
- [ ] Badeig, F., Pelorson, Q., Arias, S., Drouard, V., Gebu, I., Li, X., Evangelidis, G., Horaud, R., 2015. A distributed architecture for interacting with nao, in: *ACM International Conference on Multimodal Interaction*, pp. 385–386.
- [ ] Cao, Z., Simon, T., Wei, S.E., Sheikh, Y., 2017. Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields, in: *IEEE Conference on Computer Vision and Pattern Recognition*.
- [ ] Cruz, F., Parisi, G.I., Twiefel, J., Wermter, S., 2016. Multi-modal integration of dynamic audiovisual patterns for an interactive reinforcement learning scenario, in: *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 759–766.
- [ ] Gebu, I., Ba, S., Li, X., Horaud, R., 2017. Audio-visual speaker diarization based on spatiotemporal bayesian fusion. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.

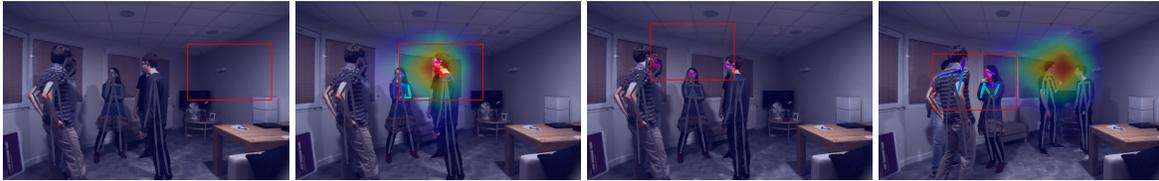


Figure 8: Example of a sequence from the AVDIAR dataset. The SSL heatmap is superposed on the frame, and visible joints/limbs are displayed using a colored skeleton. The agent’s field of view (in red) is randomly initialized (1st column), detects sound emitted by the a person and moves accordingly (2nd). The agent manages to get all the persons in the field of view (3rd), and it looks at the three people group when two persons move apart (4th).



Figure 9: Example of a sequence in a two-person scenario. First row shows an overview of the scene, including the participants and the robot. Second row shows the robot’s field of view. The robot’s head is first initialized in a position where no face is visible (1st column), and the model uses the available detections (elbow/wrist) to find the person on the right (2nd column). The robot finds the second person by looking around while keeping the first person in its field of view (3rd column), and tracks the two people walking together (4th column).

- [ ] Ghadirzadeh, A., Bütepage, J., Maki, A., Kragic, D., Björkman, M., 2016. A sensorimotor reinforcement learning framework for physical Human-Robot Interaction, in: IEEE/RSJ International Conference on Intelligent Robots and Systems, pp. 2682–2688.
- [ ] Goodfellow, I., Bengio, Y., Courville, A., 2016. Deep learning. MIT press.
- [ ] Goodrich, M.A., Schultz, A.C., 2007. Human-robot Interaction: A Survey. Foundations and Trends in Human-Computer Interaction 1, 203–275.
- [ ] Hochreiter, S., Schmidhuber, J., 1997. Long short-term memory. Neural Computation 9, 1735–1780.
- [ ] Kaelbling, L.P., Littman, M.L., Moore, A.W., 1996. Reinforcement learning: A survey. Journal of artificial intelligence research 4, 237–285.
- [ ] Kendon, A., 1967. Some functions of gaze-direction in social interaction. Acta Psychologica 26, 22 – 63.
- [ ] Kingma, D.P., Ba, J., 2014. Adam: A method for stochastic optimization, in: International Conference on Learning Representations.
- [ ] Kober, J., Bagnell, J.A., Peters, J., 2013. Reinforcement learning in robotics: A survey. The International Journal of Robotics Research 32, 1238–1274.
- [ ] Li, X., Girin, L., Badeig, F., Horaud, R., 2016. Reverberant sound localization with a robot head based on direct-path relative transfer function, in: IEEE/RSJ International Conference on Intelligent Robots and Systems, pp. 2819–2826.
- [ ] Li, X., Girin, L., Horaud, R., Gannot, S., 2017. Multiple-speaker localization based on direct-path

- features and likelihood maximization with spatial sparsity regularization. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 25, 1997–2012.
- [ ] Ljungblad, S., Kotrbova, J., Jacobsson, M., Cramer, H., Niechwiadowicz, K., 2012. Hospital Robot at Work: Something Alien or an Intelligent Colleague?, in: *ACM Conference on Computer Supported Cooperative Work*, pp. 177–186.
  - [ ] Massé, B., Lathuilière, S., Mesejo, P., Horaud, R., 2017. A reinforcement learning approach to sensorimotor control in human-robot interaction, in: *Submitted to IEEE International Conference on Robotics and Automation*.
  - [ ] Mitsunaga, N., Smith, C., Kanda, T., Ishiguro, H., Hagita, N., 2006. Robot behavior adaptation for human-robot interaction based on policy gradient reinforcement learning. *Journal of the Robotics Society of Japan* 24, 820–829.
  - [ ] Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., Riedmiller, M., 2013. Playing Atari With Deep Reinforcement Learning, in: *NIPS Deep Learning Workshop*.
  - [ ] Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A.A., Veness, J., Bellemare, M.G., Graves, A., Riedmiller, M., Fidjeland, A.K., Ostrovski, G., et al., 2015. Human-level control through deep reinforcement learning. *Nature* 518, 529–533.
  - [ ] Pourmehr, S., Thomas, J., Bruce, J., Wawerla, J., Vaughan, R., 2017. Robust sensor fusion for finding HRI partners in a crowd, in: *IEEE International Conference on Robotics and Automation*, pp. 3272–3278.
  - [ ] Qureshi, A.H., Nakamura, Y., Yoshikawa, Y., Ishiguro, H., 2016. Robot gains social intelligence through multimodal deep reinforcement learning, in: *IEEE International Conference on Humanoid Robots*, pp. 745–751.
  - [ ] Qureshi, A.H., Nakamura, Y., Yoshikawa, Y., Ishiguro, H., 2017. Show, attend and interact: Perceivable human-robot social interaction through neural attention Q-network, in: *IEEE International Conference on Robotics and Automation*, pp. 1639–1645.
  - [ ] Rothbucher, M., Denk, C., Diepold, K., 2012. Robotic gaze control using reinforcement learning, in: *IEEE International Workshop on Haptic Audio Visual Environments and Games*, pp. 83–88.
  - [ ] Sauppé, A., Mutlu, B., 2015. The Social Impact of a Robot Co-Worker in Industrial Settings, in: *ACM Conference on Human Factors in Computing Systems*, pp. 3613–3622.
  - [ ] Skantze, G., Hjalmarsson, A., Oertel, C., 2014. Turn-taking, feedback and joint attention in situated human-robot interaction. *Speech Communication* 65, 50–66.
  - [ ] Sutton, R.S., Barto, A.G., 1998. *Introduction to Reinforcement Learning*. 1st ed., MIT Press.
  - [ ] Thomaz, A.L., Hoffman, G., Breazeal, C., 2006. Reinforcement learning with human teachers: Understanding how people want to teach robots, in: *IEEE International Symposium on Robot and Human Interactive Communication*, pp. 352–357.
  - [ ] Vázquez, M., Steinfeld, A., Hudson, S.E., 2016. Maintaining awareness of the focus of attention of a conversation: A robot-centric reinforcement learning approach, in: *IEEE International Symposium on Robot and Human Interactive Communication*, pp. 36–43.
  - [ ] Watkins, C.J.C.H., Dayan, P., 1992. Q-learning. *Machine Learning* 8, 279–292.
  - [ ] Williams, R.J., 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning* .
  - [ ] Zarak, A., Mazzei, D., Giuliani, M., Rossi, D.D., 2014. Designing and Evaluating a Social Gaze-Control System for a Humanoid Robot. *IEEE Transactions on Human-Machine Systems* 44, 157–168.