

# The Head Turning Modulation system: an active multimodal paradigm for intrinsically motivated exploration of unknown environments

Benjamin Cohen-Lhyver<sup>1</sup>, Sylvain Argentieri<sup>1,\*</sup> and Bruno Gas<sup>1</sup>

<sup>1</sup> Sorbonne Université, CNRS, Institut des Systèmes Intelligents et de Robotique, ISIR, F-75005 Paris, France

Correspondence\*:

Sorbonne Université, ISIR, 4 place Jussieu, 75005 Paris, France  
sylvain.argentieri@sorbonne-universite.fr

## 2 ABSTRACT

Over the last twenty years, a significant part of the research in exploratory robotics partially switches from looking for the most efficient way of exploring an unknown environment to finding what could motivate a robot to autonomously explore it. Moreover, a growing literature focuses not only on the topological description of a space (dimensions, obstacles, usable paths, etc.) but rather on more semantic components, such as multimodal objects present in it. In the search of designing robots that behave autonomously by embedding life-long learning abilities, the inclusion of mechanisms of attention is of importance. Indeed, be it endogenous or exogenous, attention constitutes a form of intrinsic motivation for it can trigger motor command towards specific stimuli, thus leading to an exploration of the space. The Head Turning Modulation model presented in this paper is composed of two modules providing a robot with two different forms of intrinsic motivations leading to triggering head movements towards audiovisual sources appearing in unknown environments. First, the Dynamic Weighting module implements a motivation by the concept of Congruence, a concept defined as an adaptive form of semantic saliency specific for each explored environment. Then, the Multimodal Fusion & Inference module implements a motivation by the reduction of Uncertainty through a self-supervised online learning algorithm that can autonomously determine local consistencies. One of the novelty of the proposed model is to solely rely on semantic inputs (namely audio and visual labels the sources belong to), in opposition to the traditional analysis of the low-level characteristics of the perceived data. Another contribution is found in the way the exploration is exploited to actively learn the relationship between the visual and auditory modalities. Importantly, the robot—endowed with binocular vision, binaural audition and a rotating head—does not have access to prior information about the different environments it will explore. Consequently, it will have to learn in real-time what audiovisual objects are of “importance” in order to rotate its head towards them. Results presented in this paper have been obtained in simulated environments as well as with a real robot in realistic experimental conditions.

Keywords: multimodal perception, attention, motivation, active learning, binaural audition

## 1 INTRODUCTION

One of the most critical and important task humans are able to do is to explore unknown environments, topologically or semantically, while being able to create internal representations of them for localization in it and interaction with it. Such cerebral representations, or maps as it is often referred to (O'Keefe and Nadel, 1978; Cuperlier et al., 2007), enable humans and animals in general to gather and organize perceptual cues (visual, acoustic, tactile, olfactory, proprioceptive...) in semantic components. In parallel, in the mobile robotics community, exploration of unknown environments has been one of the most important fields studied, back to the artificial turtles of Walter (1951) and later to the vehicles of Braitenberg (1986). Indeed, being able for a mobile robot to simultaneously (i) map the world it is exploring, (ii) locate itself in it, and (iii) trigger relevant motor actions for further exploration (i.e. the three key tasks to perform in an exploration scheme according to Makarenko et al. (2002)), has shown to be a hard, but critical for robots' existence, problem to solve. While many artificial systems have been implemented with the sole purpose of *exploring the most of an environment* with only efficiency as a goal (Smith et al., 1988; Henneberger et al., 1992; Montemerlo et al., 2002; Carrillo et al., 2015), some more recent algorithms emerged on the basis of the precursor works of Berlyne (1950, 1965), who stated that *Motivation* is a fundamental mechanism in spontaneous exploratory behaviors in humans. Following this principle, exploration would not be driven by a goal defined by an external agent (such as the human experimenter) but rather by internal goals defined by the robot itself, that is *intrinsic* motivations (Ryan and Deci, 2000; Oudeyer and Kaplan, 2008). Amongst them are the motivations by *Curiosity*, first mathematically modeled by Schmidhuber (1991), by *Uncertainty* (Huang and Weng, 2002), by *Information gain* (Roy et al., 2001), or by *Empowerment* (Capdepuy et al., 2007). Intrinsic motivation has extensively been used during the last twenty years in several powerful systems, in particular by Oudeyer et al. (2007) with the development of the Independent Adaptive Curiosity algorithm (IAC) and the later updated systems (R-IAC (Baranes and Oudeyer, 2009) and SAGG-RIAC (Baranes and Oudeyer, 2010)). Systems based on such motivations to explore/understand an environment incorporate in particular the notion of *reward*, a principle that is of high importance in learning in primates and humans (Rushworth et al., 2011). As such, these systems are particularly suited for adaptive life-long learning robots for they bring to them wider motivations to react to their environments: instead of compelling the robot to "*explore as quickly as possible every inch of the room*", it becomes closer to "*just be curious*". But beyond the topological characteristics of unknown environments, their content also provides valuable information for the robots internal representation of the world (object formation, their affordance, etc.). Then, while one of the most predominant issue in driving topological exploration is to decide what is the next point or area to explore, semantical exploration can be also introduced to determine what is the next component to discover. Such considerations are close to attentional behaviors, which have also been extensively studied (Hopfinger et al., 2000; Downar et al., 2000; Corbetta and Shulman, 2002; Corbetta et al., 2008; Petersen and Posner, 2012).

Among others, saliency is known to be a key feature in attention thanks to its sensitivity to discontinuity in perceived data. A significant literature can be found on saliency-driven exploration: eye saccades modelization (Itti et al., 1998; Oliva et al., 2003; Le Meur and Liu, 2015), detection of auditory salient events (Kayser et al., 2005; Duangudom and Anderson, 2007), or audiovisual objects exploration (Ruesch et al., 2008; Tsiami et al., 2016). However, most of these models propose either a solely off-line solution requiring prior training from large databases, or an immutable saliency characterization of events. Moreover, the fact that these models only deal with the low-level characteristics of the perceived data leads often to an absence of wider context inclusion, be it through a form of memory, or through the semantics of the events. In addition, saliency can somehow differ from importance, depending on the task to accomplish: attention can be driven by behaviorally important but not salient stimuli while, on the other hand, very salient

73 stimuli but showing no behavioral importance can be disregarded by the attentional networks (Corbetta and  
74 Shulman, 2002; Indovina and Macaluso, 2007). However, it is worth mentioning the interesting feature of  
75 the multimodal model of salience of (Ruesch et al., 2008) as the implementation of an additional inhibition  
76 map to the ones already used for saliency. Such map promotes the exploration of unknown parts of the  
77 environments and avoids deadlock situations caused by local minima. This has also to be brought close to  
78 the notion of motivations for exploration mentioned above since a form of Curiosity is here implemented.

79 In this paper is presented a computational system, The Head Turning Modulation system (HTM), which  
80 aims at giving a mobile robot endowed with binaural hearing, binocular vision and a rotating head, the  
81 ability to decide which audiovisual sources present in unknown environments are worth the robot's attention.  
82 The principle of attention mentioned in this paper is based on the prime definition originating from James  
83 (1890): "Everyone knows what attention is. It is the taking possession by the mind, in clear and vivid form,  
84 of one out of what seem several simultaneously possible objects or trains of thought". More particularly,  
85 the proposed HTM system is dedicated to the implementation of an overt and endogenous (Driver and  
86 Spence, 1998; Le Meur et al., 2006) attentional reaction: *the head turning*. This reaction, known to be  
87 one of the attentional behavior involved in the mechanism of *attention reorientation* to unpredictable  
88 stimuli (Thompson and Masterton, 1978; Corbetta et al., 2008; Corneil et al., 2008), aims at bringing the  
89 visual sensors in front of the sources of interest hence enabling the robot to gather and analyze additional  
90 data. In addition, the HTM system provides the robot with an adaptive enough online learning behavior so  
91 that it can endlessly integrates new useful information to its self-created audiovisual database. However,  
92 this learning relying intensively upon the triggering of head movements, it is also necessary for the robot to  
93 understand when this knowledge is robust and relevant enough, thus not requiring further motor reaction.  
94 The HTM is part of a much wider system, implemented as the TWO!EARS software<sup>1</sup>, which aims at  
95 providing a computational framework for modeling active exploratory listening that assigns meaning to  
96 auditory scenes. More precisely, it consists in perceiving and analyzing a multimodal world through a  
97 paradigm that combines a classical bottom-up signal-driven processing step together with a top-down  
98 cognitive feedback. In there, the HTM is in charge of building an internal semantic map of the explored  
99 environment, made of localized audiovisual objects coupled with their respective semantic importance, the  
100 so-called congruence.

101 In comparison with other works, the proposed system is described as a *real-time* (Huang et al., 2006)  
102 and *online* behavioral unit, which is always able to learn new situations while also taking advantage of  
103 its previous experience of the past environments. In terms of architecture, the proposed system receives  
104 data from several "experts" from the TWO!EARS software, i.e. computational elements specialized in very  
105 particular tasks, such as the identification of sounds or images. It means that the HTM system is placed  
106 right after these experts, and thus receives already highly interpreted data. Two main parts constitute the  
107 system: an attentional component, the *Dynamic Weighting* module (Walther and Cohen-Lhyver, 2014), and  
108 a learning component, the *Multimodal Fusion & Inference* module (Cohen-Lhyver et al., 2015). On the one  
109 hand, the DW module is dedicated to the analysis of perceived audiovisual objects through the concept of  
110 *Congruence*, defined as a semantic saliency and rooted in the principle of optimal incongruity (Hunt, 1965).  
111 The DW module implements a form of motivation by surprise for it favors unexpected audiovisual events.  
112 On the other hand, the MFI module learns the association between auditory and visual data in order to  
113 make the notion of *multimodal object* arise from potentially erroneous data of the aforementioned experts.  
114 The MFI module implements a form of motivation by reduction of uncertainty for it aims at consolidating  
115 as much as needed its knowledge about the audiovisual objects that the robot encounters. This learning

<sup>1</sup> <http://www.twoears.eu>

116 serves two purposes. First, it might improve the robustness and reliability of the classification (Droniou  
117 et al., 2015). Secondly, it allows the system to perform missing information inference (Bauer and Wermter,  
118 2013), as when an object is placed behind the robot thus having only access to the auditory information.

119 The paper is organized as follows. To begin with, the overall TWO!EARS framework, together with the  
120 notations used all along the paper, are introduced in a first section. On this basis, the overall HTM system is  
121 thoroughly presented in a second section: after a short insight into the HTM system architecture, the way the  
122 DW module and the MFI module operate is formalized. This section also presents their respective evaluation  
123 in simulated conditions. Then, the combination of the two modules is investigated and the evaluation of  
124 the approach in real experimental conditions, that is including a real robot in a real environment, is made.  
125 Finally, a conclusion ends the paper.

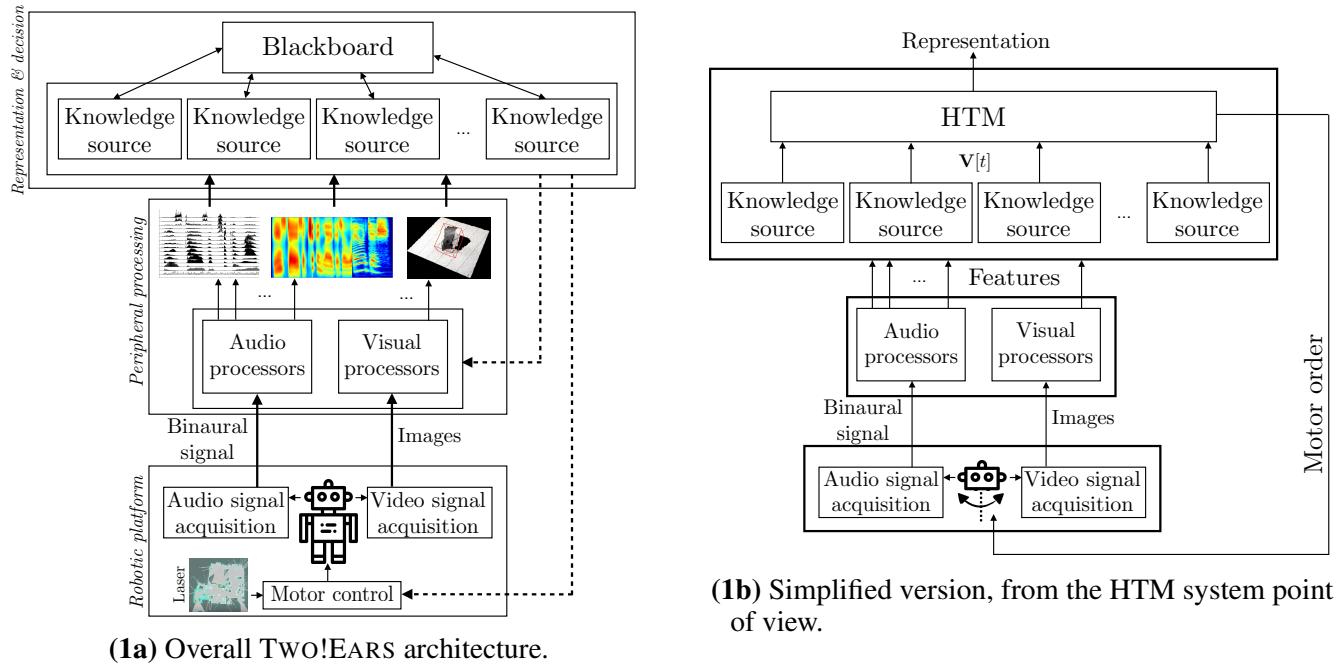
## 2 CONTEXT AND NOTATIONS

126 This section presents the context in which is rooted the proposed HTM system. All the forthcoming  
127 development has taken place inside a specific computational architecture aiming at modeling an integral,  
128 multimodal, intelligent and active auditory perception and experience. This model physically uses  
129 two human-like ears and visual inputs to make a mobile robot able to interactively explore unknown  
130 environments, see Two!Ears (2016b). Among other applications, this modular architecture targets evaluation  
131 of bottom-up audiovisual processing coupled with top-down cognitive processes. The proposed HTM  
132 system relies also on this top-down and bottom-up paradigm providing the robot with a reliable internal  
133 representation of its audiovisual environment. To begin with, a short overview of the overall architecture  
134 is proposed in a first subsection. A second subsection introduces the notations used all along the paper,  
135 together with all the notions required to understand the HTM system.

### 136 2.1 Global framework

137 All the forthcoming developments have been conducted inside the multilayer TWO!EARS architecture,  
138 see Figure 1a. This figure highlights two different pathways: first, a classical bottom-up processing  
139 way, where raw data coming from the sensors (microphones and cameras) are first analyzed (features  
140 extraction step), processed (through some specialized pattern recognition algorithms) and interpreted  
141 (representation and decisional layers). All of the above is computed by dedicated Knowledge Sources (KS).  
142 The main contribution of this architecture is that all these layers are highly and dynamically parameterizable:  
143 for instance, most of the feature extractions parameters (for audio data, one could cite the number of  
144 Gammatone filters used, their repartition on the frequency scale, etc.) can be changed on the fly. In general,  
145 the decision to change parameters comes from upper layers, resulting in a top-down pathway, also involving  
146 decisions concerning the movement of the robot itself. Such decisions concerning the robot actions are of  
147 particular importance, especially when dealing with attention reorientation and scene understanding for  
148 they add adaptability to new and unpredictable events.

149 The HTM system inside the TWO!EARS architecture shown in Figure 1b is implemented as a Knowledge  
150 Source (KS). It gets data from other KSs available in the architecture through a blackboard (Schymura  
151 et al., 2014) (which can be seen, with a rough simplification, as a data structure), and provides as an output  
152 a proposition for a motor command, together with an interpreted representation of the robot's world. One  
153 originality of the approach is that the HTM system is placed behind other KS, thus not working directly  
154 with the features extracted from the raw audio and visual signals. All of the KSs the HTM relies upon  
155 contribute to the scene analysis and are fused by the HTM into a representation of the world that spans  
156 wider in time than the one provided by the individual KSs. This representation is made of all the unknown



(1a) Overall TWO!EARS architecture.

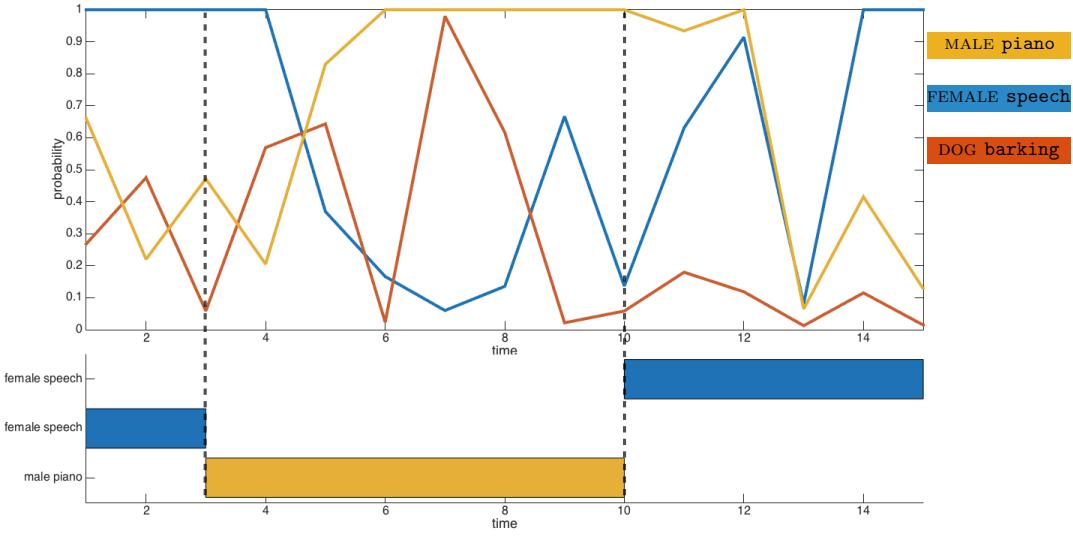
(1b) Simplified version, from the HTM system point of view.

**Figure 1.** Two!EARS architecture. (Left:) On the basis on audio and visual data, features are extracted to provide a compact description of the data. Several audio and visual experts (or Knowledge Sources, KS) exploit these features to analyze the signals. Each KS is specialized in one task: recognition of one type of sound, localization, separation, etc. All experts share their knowledge through a blackboard system, thus producing an internal representation of the world. On this basis, the overall system (but also individual KSs) can decide to modulate either the feature extraction step, or the action of the robot. The proposed HTM system, implemented as a KS, is—among others—responsible for this last modulation. (Right:) Focus on the implementation of the HTM system inside this architecture.

157 environments explored by the robot, each of them being characterized by the audiovisual objects observed  
 158 there in an allocentric representation, coupled with an additional semantic layer formalized through the  
 159 notion of Congruence. The data used by the HTM, together with their notations are described in the  
 160 following section.

## 161 2.2 Definitions and notations

162 The HTM system only relies upon KSs outputs to analyze the unknown environments the robot  
 163 explores. These KSs are classification experts specialized in the recognition of audio or visual  
 164 frames (Two!Ears, 2016a), classified in terms of audio classes  $c_i^a$ , with  $i = 1, \dots, N_a$  (such as,  
 165  $c_i^a \in \{\text{voice, barking, yelling, ...}\}$ ) or visual classes  $c_k^v$ , with  $k = 1, \dots, N_v$  (such as  $c_k^v \in$   
 166  $\{\text{DOG, BABY, MALE PERSON, ...}\}$ ) with  $N_a$  and  $N_v$  the number of audio and visual classes respectively.  
 167 All classifiers are mutually independent, each providing a probability  $p_i^a[t]$  and  $p_k^v[t]$  for the frame at  
 168 time  $t$  to belong to the class they represent. All these probabilities are regrouped by modality in the  
 169 two vectors  $\mathbf{P}^a[t] = (p_1^a[t], \dots, p_{N_a}^a[t])$  and  $\mathbf{P}^v[t] = (p_1^v[t], \dots, p_{N_v}^v[t])$ . In addition, the TWO!EARS  
 170 architecture provides  $N_\theta$  localization experts (May et al., 2011; Ma et al., 2015), aiming at localizing  
 171 audio and/or visual events in the horizontal plane with respect to the robot. Each of them outputs a  
 172 probability  $p_{\theta_u}^a[t]$  and  $p_{\theta_u}^v[t]$ , with  $u = 1, \dots, N_\theta$ , for an audio and/or visual event to originate from the  
 173 azimuth  $\theta_u^a$  or  $\theta_u^v$  (by convention,  $\theta = 0^\circ$  corresponds to an event placed in front of the robot). All these  
 174 probabilities are gathered into the audio and visual localization vectors  $\Theta^a[t] = (p_{\theta_1}^a[t], \dots, p_{\theta_{N_\theta}}^a[t])$  and  
 175  $\Theta^v[t] = (p_{\theta_1}^v[t], \dots, p_{\theta_{N_\theta}}^v[t])$ . In practice, all these classifiers outputs are regrouped into a single vector



**Figure 2.** Illustration of the audio classification experts on real perceived data. (Top panel:) Probabilities of the frames to belong to the corresponding audio classes. (Bottom panel:) Time description of the audiovisual objects appearance.

176  $\mathbf{V}[t]$  constituting the sole HTM system input, with

$$\mathbf{V}[t] = (\mathbf{P}[t], \Theta[t]), \text{ with } \mathbf{P}[t] = (\mathbf{P}^a[t], \mathbf{P}^v[t]) \text{ and } \Theta[t] = (\Theta^a[t], \Theta^v[t]). \quad (1)$$

177 From  $\mathbf{V}[t]$ , the HTM model attempts to build a stable and reliable internal representation of the world,  
178 environment by environment. Such a representation is obtained by transforming an audio and/or visual  
179 event  $\Psi_j$  objectively present in the environment at azimuth  $\theta(\Psi_j)$  and belonging to the ground truth  
180 audiovisual class  $c(\Psi_j) = \{c^a(\Psi_j), c^v(\Psi_j)\}$ , into an object  $o_j$  perceived by the robot, i.e.

$$\begin{aligned} \Psi_j = \{\theta(\Psi_j), c(\Psi_j)\} &\longrightarrow o_j = \{\hat{\theta}(o_j), \hat{c}(o_j)\}, \\ \text{with } \hat{\theta}(o_j) = \left\{ \begin{array}{ll} \theta_u^a, & \text{with } u = \arg \max_i (p_{\theta_i}^a), \text{ if } \theta_u^a \geq |\theta_{\text{HTM}} - \theta_{\text{FOV}}| \\ \theta_u^v, & \text{with } u = \arg \max_k (p_{\theta_k}^v) \text{ otherwise} \end{array} \right., \\ \text{and } \hat{c}(o_j) = \{\hat{c}^a(o_j), \hat{c}^v(o_j)\}, \end{aligned} \quad (2)$$

181 where  $\theta_{\text{HTM}}$  and  $\theta_{\text{FOV}}$  represent the current azimuthal head position and the field of view of the camera,  
182 respectively. Then, an object  $o_j$  is defined by its estimated angular position  $\hat{\theta}(o_j)$  and its estimated  
183 audiovisual class  $\hat{c}(o_j)$  made of the estimated audio class  $\hat{c}^a(o_j)$  and estimated visual class  $\hat{c}^v(o_j)$ .  
184 Equation (2) also indicates that the estimated angular position is obtained from the audio localization  
185 experts when the objects are out of the robot sight; otherwise, visual localization experts are exploited.  
186 Because of localization and/or classification errors, the object  $o_j$  might differ from the corresponding  
187  $\Psi_j$ . As an example, Figure 2 plots as a function of temporal frames experimental data from three audio  
188 classifiers outputs corresponding to the audio classes PIANO, SPEECH and BARKING. This figure shows  
189 first that potential classification errors can obviously occur: at time  $t = 7$ , the BARKING output probability  
190 reaches about 98% while a piano sound is perceived by the robot. Additionally, the data show the temporal  
191 dynamic audio experts can exhibit: while the piano starts playing at time  $t = 3$ , the corresponding audio  
192 expert becomes dominant a few frames later only. This delay observed experimentally will justify later  
193 technical implementation specifics.

194 At this point, the notion of object already constitutes more than just a structure of data. In particular, the  
 195 objects created by the HTM system embed a short-term temporal smoothing of the data  $\mathbf{P}(o_j)$  they are  
 196 associated with, as

$$\mathbf{P}(o_j)[t] = \frac{1}{N_t} \sum_{n=t-N_t}^{n=t} \mathbf{P}(o_j)[n], \quad (3)$$

197 where  $N_t \leq 10$  is the number of frames during which data  $\mathbf{P}$  have been associated to  $o_j$ . This temporal  
 198 smoothing enables the robot to take into account its past experience of the audiovisual data the robot  
 199 perceived and that have been associated with this object, but also to lower the impact of the early potential  
 200 erroneous outputs from the classification experts. Indeed, experiments have shown that most of them are  
 201 prone to making more errors during the very first frames of perceived events. Thus, it is one of the goal  
 202 of the HTM system to make the object identical to the event, even in the presence of classification errors.  
 203 In all the following, the internal representation  $e^{(l)}$  of the  $l$ -th environment being explored by the system  
 204 is defined as the collection of the  $N^{(l)}$  objects inside, i.e.  $e^{(l)} = \{o_1, \dots, o_{N^{(l)}}\}$ . Note that  $N^{(l)}$  evolves  
 205 along time all along the agent life on the basis of the perceived data. Importantly, this definition of an  
 206 environment—which will be augmented later on in Section 3.2.1)—aims at making the difference  
 207 between the topological and the semantic definition of an environment (see Section 1). While the robot,  
 208 through its navigation system, gets to know when a new topological environment is being explored, the  
 209 HTM analyzes its audiovisual content in order to assess whether this environment is really new or if it  
 210 similar to a previously explored one (as explained in Section 3.2.1). In that case, this audiovisual similarity  
 211 enables the robot to apply previously self-created behavioral rules making its reaction abilities way quicker.  
 212 Then, one can define audiovisual categories  $\mathcal{C}^{(l)}(c_i^a, c_k^v)$  of this  $l$ -th representation with

$$\mathcal{C}^{(l)}(c_i^a, c_k^v) = \{o_j \in e^{(l)}, \hat{c}^a(o_j) = c_i^a \text{ and } \hat{c}^v(o_j) = c_k^v\}. \quad (4)$$

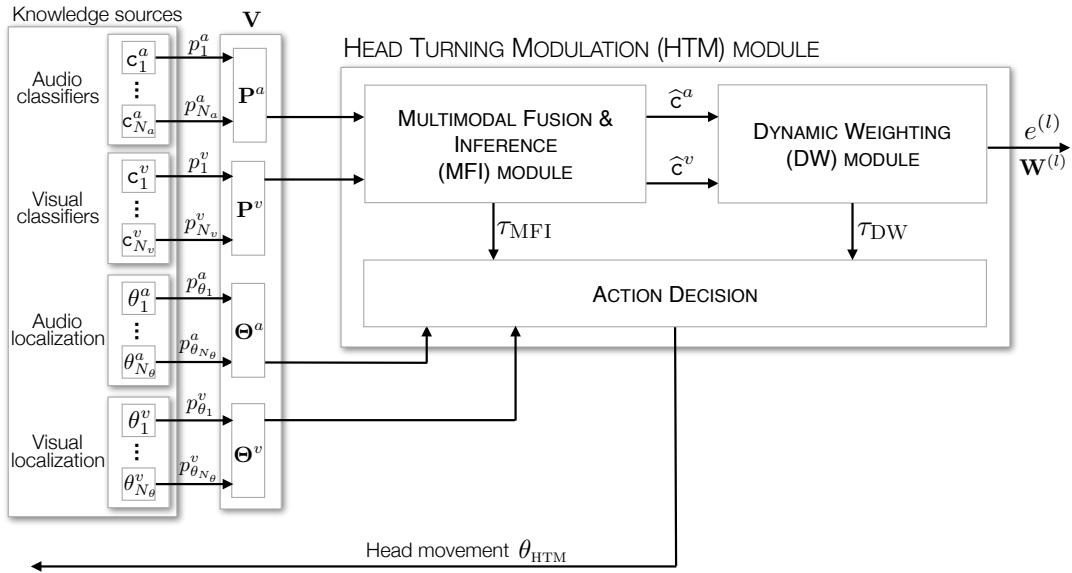
213 Once the events have been interpreted as objects within the internal representation  $e^{(l)}$  of the robot, the  
 214 HTM system analyses them through the notion of *Congruence*, described in the next section.

### 3 THE HEAD TURNING MODULATION SYSTEM

215 The *Head Turning Modulation* system is an attempt to provide a binaural and binocular humanoid robot  
 216 with the ability to learn by its own how to react to unpredictable events and to consequently trigger or inhibit  
 217 head movements toward them. Moreover, the system is endowed with a module that provides a multimodal  
 218 internal representation of the world through a real-time learning paradigm that has no access to any prior  
 219 knowledge about the environments to be explored. This system, partially introduced by the authors in  
 220 (Cohen-Lhyver et al., 2015, 2016; Cohen-Lhyver, 2017) is defined as a model of attention supported by  
 221 an object-based representation of the world. This section will thus present separately the two constitutive  
 222 modules of the HTM system. An evaluation of each of them will be presented in simulated conditions,  
 223 while the evaluation of the whole system, made in real conditions, will be presented in Section 4.

#### 224 3.1 Architecture of the proposed system

225 The overall architecture of the HTM system is depicted in Figure 3. It exhibits two modules inside,  
 226 each of them being dedicated to one specific task. As outlined in Section 2.2, the HTM inputs are made  
 227 of audio and visual classifiers outputs, which are used by the first module—the Multimodal Fusion &  
 228 Inference (MFI) module—to provide an estimation of the audiovisual class  $\hat{c}(o_j) = \{\hat{c}^a(o_j), \hat{c}^v(o_j)\}$  of  
 229 the currently analyzed frame. As will be shown later, such an estimation is made possible by a top-down



**Figure 3.** Architecture of the HTM system. It is made of two modules, dedicated to the estimation of the audiovisual class of an event and to the computation of its importance. An additional element is in charge with the respective motor orders integration and decision.

230 motor feedback that allows the system to gather additional audio and visual data. On the basis on the  
 231 frame estimated classification, a second module—the Dynamic Weighting (DW) module—is in charge  
 232 of deciding if the currently emitting object is of interest through the computation of its congruence to  
 233 the current environment. As a result, this module also exploits the motor feedback to modulate the robot  
 234 attention. Since both modules require motor actions for their operations, a supplemental element is in  
 235 charge with prioritizing them, depending on their respective motor activities  $\tau_{DW}$  and  $\tau_{MFI}$ , see Figure 3.  
 236 Motor decisions are taken by using the localization experts providing an estimated angle of the processed  
 237 event. Finally, the overall HTM system outputs a list of interpreted objects, i.e. an internal representation  
 238 of the explored environment, which can be used by other KS in the TWO!EARS architecture for other tasks  
 239 (modulating the exploration depending on the objects in the environment, deciding which object is of  
 240 particular interest in the current scenario on the basis on the DW module module conclusions, exploiting the  
 241 top-down architecture to refine the peripheral processing steps, etc.). All of these modules are introduced  
 242 in the next subsections together with some intermediate illustrations and evaluations of their functioning.

### 243 3.2 The Dynamic Weighting Module

244 The *Dynamic Weighting* module (DW module) is the attentional part of the HTM system aiming at giving  
 245 the robot an hypothesis about a possible relevant audiovisual object that would present an interest to it,  
 246 in the scope of the exploration of unknown environments. As already stated, this interest is formalized  
 247 through the new notion of *Congruence*, thereafter detailed.

#### 248 3.2.1 Congruence: definition and formalization

249 Congruence is a notion that defines the relationship between an audiovisual event to the environment  
 250 it is occurring in. It has to be brought next to the well-known and studied notion of Saliency (Treisman  
 251 and Gelade, 1980; Nothdurft, 2006; Duangudom and Anderson, 2007) that describes how the perceived  
 252 characteristics of a stimulus exhibit continuity, or not, with its direct surrounding. Whereas saliency is  
 253 based on low-level characteristics of the signals (such as intensity, frequencies, pitch, color, contrast, etc.),

254 Congruence relies on a higher representation of the audio and visual signals, namely the audiovisual class  
 255 they belong to (see Section 2.2). Congruence is thus defined as a *semantic saliency* for it relies on an  
 256 already interpreted representation of the perceived data. Since the robot does not have any prior knowledge  
 257 about the possible likelihood of an audiovisual event to occur in an environment, the DW module will only  
 258 base its analysis on a posteriori probabilities, that is computing statistics only on what has been observed  
 259 so far by the system, environment by environment. This probability of an object  $o_j$  to belong to a certain  
 260 audiovisual category  $\mathcal{C}^{(l)}(c_i^a, c_k^v)$  is thus defined as

$$p(o_j \in \mathcal{C}^{(l)}(c_i^a, c_k^v)) = p(\mathcal{C}^{(l)}(c_i^a, c_k^v)) = \frac{|\mathcal{C}^{(l)}(c_i^a, c_k^v)|}{N^{(l)}}, \quad (5)$$

261 where  $|\mathcal{C}^{(l)}(c_i^a, c_k^v)|$  depicts the number of objects that have already been associated to the audiovisual  
 262 category  $\mathcal{C}^{(l)}(c_i^a, c_k^v)$  (as a reminder,  $N^{(l)}$  corresponds to the number of objects detected so far in the  
 263 l-th environment). Still following the fact that no prior knowledge is available for the robot, the system  
 264 will compare this a posteriori probability to a threshold  $K^{(l)} = 1/N_{\mathcal{C}^{(l)}}$  defined as the equiprobability  
 265 of an object to belong to any of the categories detected so far, where  $N_{\mathcal{C}^{(l)}}$  is the number of different  
 266 audiovisual categories detected in the l-th environment. Such criterion has been chosen so that minimal bias  
 267 is introduced in order not to promote any audiovisual category. The criterion  $K^{(l)}$  evolves through time:  
 268 the more audiovisual classes observed, the lower the criterion. Finally, the congruence decision follows:

$$o_j \in \mathcal{C}^{(l)}(c_i^a, c_k^v) \text{ is incongruent} \Leftrightarrow p(\mathcal{C}^{(l)}(c_i^a, c_k^v)) \leq K^{(l)}. \quad (6)$$

269 All the “status of congruence”, that is whether they are congruent or not, of the audiovisual  
 270 categories detected by the system in a given environment are then gathered into a binary vector  
 271  $\mathbf{W}^{(l)} = \{p(\mathcal{C}^{(l)}(c_i^a, c_k^v)) \leq_{0,1} K^{(l)}\}, \forall(i, k)$ , with  $\leq_{0,1}$  a binary comparison operator. This vector of size  
 272  $|\mathcal{C}^{(l)}|$  completes the definition of environments as they become collections of objects  $e^{(l)}$  coupled with  
 273 their congruence status  $\mathbf{W}^{(l)}$ . In consequence, an audiovisual class can be incongruent in an environment,  
 274 but congruent in another. Since the robot would explore unknown environments during its whole life, the  
 275 knowledge gained from previous explorations has to be reusable for it might speed up the exploration  
 276 of new ones. Following a rule of strict inclusion of the sets of categories observed in every environment  
 277 explored so far by the robot, if the set of categories detected during the exploration of an environment  
 278  $e^{(i)}$  has already been observed in a previous environment  $e^{(j)}$ , then  $\mathbf{W}^{(i)} = \mathbf{W}^{(j)}$ . This redefinition of an  
 279 environment implies that there is one instantiation of the DW module per environment. In addition, even  
 280 in the case where there has been a reuse of information, the rules of Congruence are still computed as if  
 281 the current environment was a completely new one. Consequently, if  $e^{(j)}$  gets to differ at a point in time  
 282 from  $e^{(i)}$  and that there is no other correspondence with other environments, the  $\mathbf{W}^{(j)}$  vector computed in  
 283 parallel from the beginning of the exploration of  $e^{(j)}$  will be from now on applied.

### 284 3.2.2 Motor orders

285 Based on the congruence of all the objects, an active behavior is defined: if an object  $o_j$  is incongruent  
 286 according to Equation (6), then it is worth focusing on it. A head movement can consequently be triggered  
 287 in the direction of this object. At the opposite, if  $p(\mathcal{C}^{(l)}(c_i^a, c_k^v)) > K^{(l)}$  the robot would inhibit this  
 288 movement. But such a binary motor decision has several drawbacks, as demonstrated in Cohen-Lhyver  
 289 et al. (2015). Among others, it presents a high sensitivity to classification errors, leading to erroneous motor  
 290 decisions. Introducing a temporal weighting  $w_{o_j}$  of each object  $o_j$ , inspired by the temporal dynamic of the  
 291 Mismatch Negativity phenomenon (Näätänen et al., 1978), filters out efficiently most of these errors. These

292 weights are computed thanks to two different functions, depending upon the probability  $p(\mathcal{C}^{(l)}(c_i^a, c_k^v))$ ,  
 293 along

$$w_{o_j}[n] = \begin{cases} f_\omega^\bullet[n] = 1/(1 + 100 e^{-2n}) & \text{if } p(\mathcal{C}^{(l)}(c_i^a, c_k^v)) \leq K^{(l)}, \\ f_\omega^o[n] = (1/1 + 0.01 e^{2n}) - 1 & \text{else,} \end{cases} \quad (7)$$

294 where  $f_\omega^\bullet[n]$  and  $f_\omega^o[n]$  are increasing positive and decreasing negative functions dedicated to the weighting  
 295 of incongruent and congruent objects respectively, and  $n$  a time index. Note that  $n$  is systematically reset to  
 296 0 whenever the congruence status of the object  $o_j$  switches. From these weights, it is possible to decide  
 297 which object has to be focused on. Such a decision is implemented through an adaptation of the GPR  
 298 model (Gurney et al., 2001a,b) of the basal ganglia-thalamus-cortex loop involved in the motor order  
 299 decision in humans. According to this model, all possible motor actions are expressed as channels of  
 300 information which are by default inhibited by several afferent external connections. Depending on the  
 301 goal or on the perceived stimuli, one of the channels is excited, thus promoting the motor action it is  
 302 representing. Inspired by this functioning, all the objects perceived by the robot are similarly represented  
 303 as information channels having a dedicated activity  $\tau_{\text{DW}}(o_j)$ . The vector of canal activities  $\boldsymbol{\tau}_{\text{DW}}$  can be then  
 304 defined as

$$\boldsymbol{\tau}_{\text{DW}} = (\tau_{\text{DW}}(o_1), \dots, \tau_{\text{DW}}(o_{N_l})) , \text{ with } \tau_{\text{DW}}(o_j) = -\frac{p(\mathcal{C}^{(l)}(c_i^a, c_k^v))}{K^{(l)}}. \quad (8)$$

305 Thus, the higher the weight  $w_{o_j}$  of an object, the lowest the activity of its corresponding canal. The angle  
 306  $\hat{\theta}(o_j)$  estimated by the audio localization expert corresponding to the canal with the lowest activity will  
 307 then be selected as the winning motor order  $\theta_{\text{DW}}$ , i.e.

$$\theta_{\text{DW}} = \hat{\theta}(o_j), \text{ with } j = \arg \min_l (\tau_{\text{DW}}(o_l)). \quad (9)$$

308 If two different objects  $o_j$  and  $o_l$  have the same weight  $w_{o_j} = w_{o_l}$ , then their corresponding channels  
 309  $\tau_{\text{DW}}(o_j)$  and  $\tau_{\text{DW}}(o_l)$  have the same value. In such a case, the most recent object in the representation is  
 310 promoted, thus introducing a motivation by novelty (Huang and Weng, 2002, 2004) (see also Walther et al.  
 311 (2005)). Then, Equation (8) is slightly modified by introducing a weight which is minimized for recently  
 312 appeared objects, i.e.

$$\tau_{\text{DW}}(o_j) = -\frac{p(\mathcal{C}^{(l)}(c_i^a, c_k^v))}{K^{(l)}} \times \frac{1}{\Delta_t(o_j)}, \text{ with } \Delta_t(o_j) = t - t_{\text{emit}}(o_j), \quad (10)$$

313 where  $\Delta_t(o_j)$  represents the time elapsed between the object appearance  $t_{\text{emit}}(o_j)$  (reset to  $t$  when the  
 314 object starts emitting again after having stopped previously) and current frame  $t$ . Note that the temporal  
 315 smoothing introduced by Equation (3) does not influence the global reactivity to unexpected events, for the  
 316 dynamics of the smoothing has the same order of magnitude to the dynamic of the weighting function in  
 317 Equation (7).

### 318 3.2.3 Simulations and evaluation of the DW module

319 The DW module aims at controlling the head movements of an exploratory robot through the notion  
 320 of Congruence of perceived audiovisual objects. Thus, what is expected from the DW module is to  
 321 either trigger movements towards important audiovisual sources, and to also be able to inhibit them  
 322 when necessary. To illustrate this, simulations have been conducted on the basis of the TWO!EARS  
 323 architecture. Importantly, twelve audio classifiers and ten visual classifiers are actually implemented inside

324 the software (Two!Ears, 2016a), making evaluation scenarios quite limited. Thus, instead of simulating  
 325 raw (audio and visual) data used by real classifiers, their outputs  $p_i^a[t]$  and  $p_k^v[t]$  are rather simulated.  
 326 Nevertheless, real conditions will be used later to evaluate the overall HTM system in Section 4. Note that  
 327 the forthcoming simulated localization experts are designed to provide the exact object audio and visual  
 328 localization, the focus being put here on the congruence analysis performed by the DW module.

329 **3.2.3.1 Simulations**

330 Multiple evaluation scenarios are proposed, each of them being described by the number  $n_S$  of different  
 331 sources in the simulated environment, the description of their azimuthal localization, their temporal  
 332 appearance and disappearance, and their ground truth audiovisual classes  $c(\Psi)$ —obviously, the HTM  
 333 system does not have access to any of these. The scenarios are also defined by the maximal number  
 334 of simultaneously emitting sound sources  $n_{sim}^{max}$ . While this number never exceeds five in real extreme  
 335 experimental conditions, the simulations allow to incorporate up to ten audiovisual sources. At every time  
 336 step  $t$  of a simulation, a vector  $\mathbf{P}[t] = (\mathbf{P}^a[t], \mathbf{P}^v[t])$ , from Equation (1) is sent to the HTM system. In  
 337 the scope of the sole DW module evaluation, the estimated audio and visual classes of an event is directly  
 338 obtained from  $\mathbf{P}[t]$ , i.e. on the KS outputs, according to a maximum a posteriori (MAP) estimation, with

$$\widehat{c}_{MAP}^a = c_i^a, i = \arg \max_l (p_l^a) \text{ and } \widehat{c}_{MAP}^v = c_k^v, k = \arg \max_l (p_l^v). \quad (11)$$

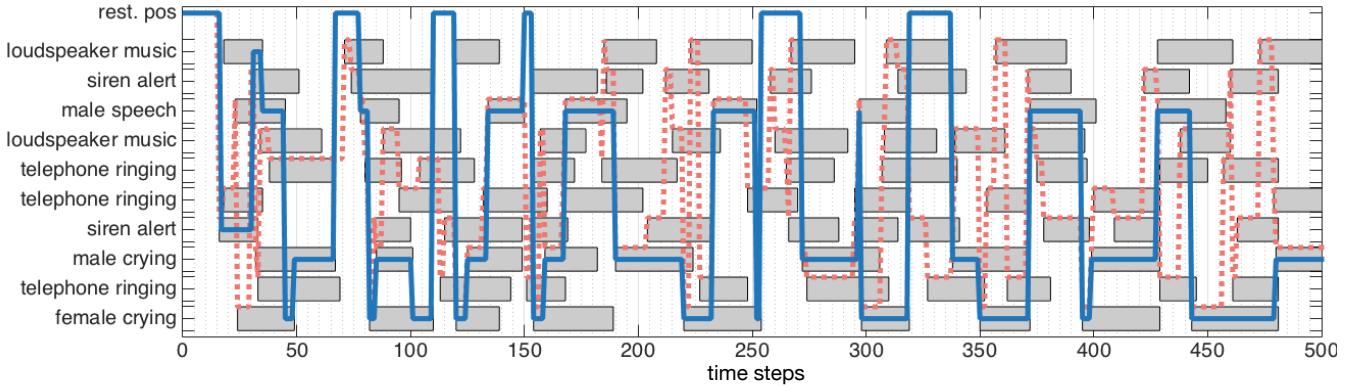
339 Note that this audiovisual class estimation will be later provided by the MFI module introduced in §3.3, as  
 340 shown in Figure 3. However, because of the inevitable presence of classification errors, the corresponding  
 341 audio and/or visual classes can be wrong (see Figure 2). It has been simulated through the implementation  
 342 of an error rate  $\varepsilon_P \in [0, 100]\%$ . At time  $t$ , a ground truth probability vector corresponding to the simulated  
 343 event is generated. With respect to  $\varepsilon_P$ , a “wrong” classification expert index is randomly selected by  
 344 drawing its value from a uniform pseudorandom number generator. Then, its associated probability is set to  
 345 be the maximal value of the whole vectors  $\mathbf{P}[t]$ . In the end, this will allow to judge the robustness of the  
 346 approach to such classification errors.

347 Like proposed in Girard et al. (2002), the performance of the system is partially evaluated in comparison  
 348 with a virtual “naive robot” noted  $\mathfrak{R}_n$ . In particular,  $\mathfrak{R}_n$  will systematically turn its head towards any  
 349 audiovisual source occurring in the environment, independently of its importance. For now, the simulations  
 350 are made with an important restriction (explained and justified later): all the sources are in the field of view  
 351 of the robot, i.e. the robot always has access to visual data.

352 **3.2.3.2 Evaluation 1: head movements modulation by the DW module**

353 A rather complex environment is used in the following to illustrate the functioning of the DW module:  
 354  $n_S = 10$  audiovisual sources are present with a maximum of  $n_{sim}^{max} = 7$  simultaneously emitting sources.  
 355 At first, let’s focus on the ability of the DW module to modulate head movements by selecting only the  
 356 sources of importance through the congruence analysis. Here will only be assessed the behavioral role  
 357 conferred by the DW module to the robot; in consequence the simulated classification experts will be set as  
 358 outputting perfect data, that is  $\varepsilon_P = 0$  (evaluations with higher error rates are made later in the paper).

359 Figure 4 exhibits one simulated environment, made of sources (represented as gray boxes) emitting  
 360 sound along time (horizontal axis). Each source belongs to an audiovisual category represented on the  
 361 left axis. Some sources might have the same audiovisual category: for instance, in this simulated scenario,  
 362 the environment is made of three different telephones ringing. In addition to this “objects along time”  
 363 description of the scene, Figure 4 shows two different lines: both “pass” through objects, indicating



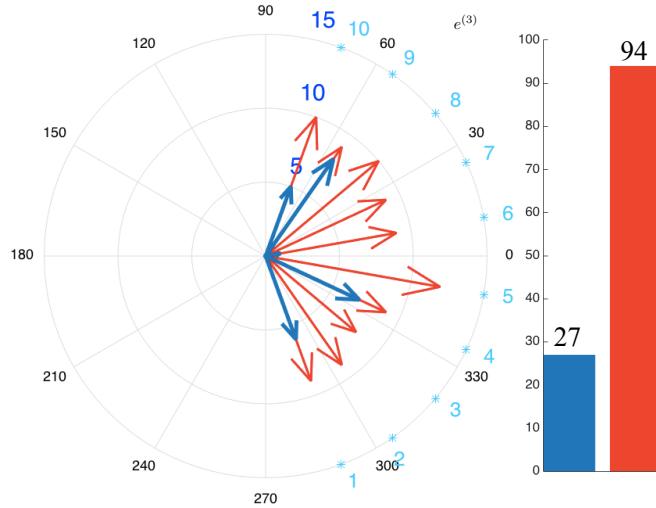
**Figure 4.** Audiovisual sources towards which a head movement has been triggered (blue line) by the DW module, (dotted red line) by the naive robot. (Gray boxes) audiovisual sources emitting sound. A source is focused on when the lines crosses the corresponding box.

364 that the robot has decided to focus on them. The blue line corresponds to the decision taken by the DW  
 365 module, while the red dashed one corresponds to  $\mathfrak{R}_n$ . Simulations show that the DW module considers the  
 366 audiovisual classes (RINGING, telephone) and (MUSIC, loudspeaker) as congruent in less than 100  
 367 time steps. This is because of their distribution with respect to the other categories: in the beginning of  
 368 the simulation, objects belonging to these two categories are often present, making them less important.  
 369 Consequently, the robot will not turn its head towards those sources: there is no (motor) attentional reaction  
 370 anymore. On the other hand, the categories (ALERT, siren), (SPEECH, male), (CRYING, female), and  
 371 (CRYING, male) are considered as incongruent, thus requiring the robot to focus on them. Importantly, the  
 372 actual meaning of those sources is not used here to decide of a reaction: one could have trade the congruent  
 373 categories with the incongruent ones without any change in the global reaction. Only the frequency of  
 374 apparition defined in Equation (5) is taken into account to decide the importance of a source.

375 In comparison, the naive robot  $\mathfrak{R}_n$  turns its head every time a source starts to emit sound: it is particularly  
 376 noticeable between  $t = 200$  and  $t = 250$  where a lot of movements can be observed. The comparison  
 377 between the two behaviors is highlighted in Figure 5, where is depicted the total number of head movements  
 378 triggered to the audiovisual sources in the environment for the DW module (blue) and naive robot (red). It  
 379 appears that a drastic modulation of the exploratory behavior is obtained: using the DW module conducts to  
 380 a reduction of 71.3% of the number of head movements in comparison with the naive robot. Furthermore,  
 381 the DW module only triggers movements toward five sources, instead of ten for the naive robot, thus  
 382 showing how Congruence—even with its simple and intuitive definition—can provide an efficient filter  
 383 for the attentional behavior of the robot. Importantly, such a modulation allows the robot to use head  
 384 movements, and more generally its exploratory actions, for other unrelated tasks. As long as no incongruent  
 385 source is detected, head movements are free to be used for anything else. But as soon as an incongruent  
 386 source pops up in the environment, the DW module will drive the head towards this source: the robot then  
 387 puts its attention on it. In the end, this simple illustration shows how important it is to be able to inhibit or  
 388 trigger head movements.

### 389 3.2.4 Conclusion and limitations

390 The DW module is a crucial part of the HTM system in charge with providing a semantic understanding  
 391 of the unknown environments the robot is supposed to explore. One of the cornerstone of this module  
 392 is to be able to work without prior knowledge about the potential distribution of the audiovisual sources

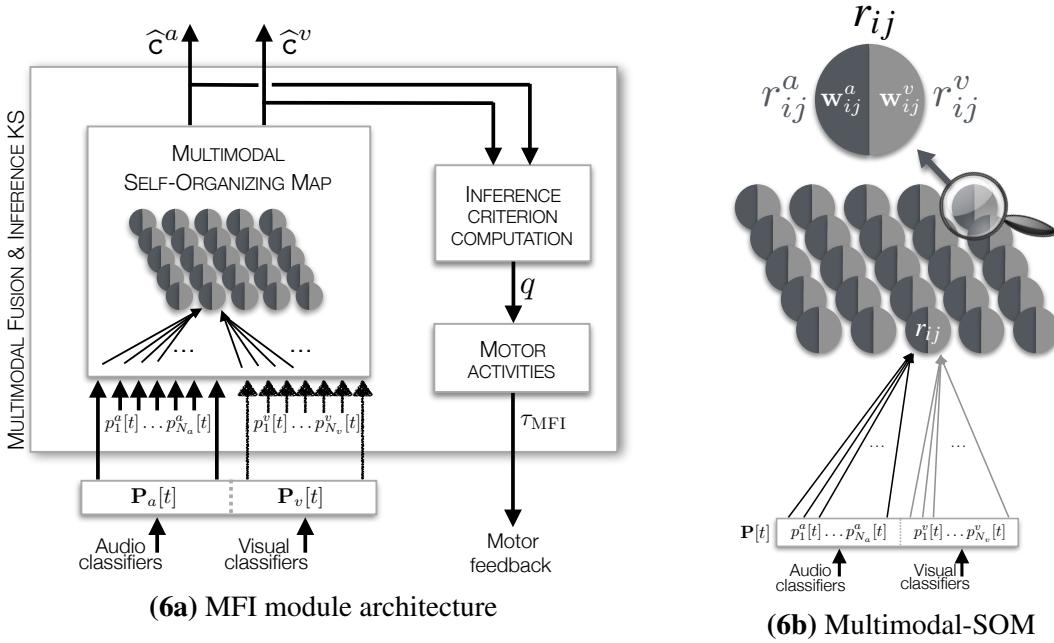


**Figure 5.** Head movements triggered by (blue) the DW module, (red) the virtual naive robot. Every arrow points towards the position of a source and their lengths depict the number of movements towards every pointed source. (Bars:) number of movements triggered by the the MFI module (blue) and the virtual naive robot (red). The light blue numbers correspond to the position of the audiovisual sources.

393 occurring in these environments. Thus, the DW module has to create congruence rules on the sole basis  
 394 of what the robot sees and hears, that is the audio and visual labels the classification experts output. The  
 395 behavior rules created are, firstly, adaptive enough to always take into account new information, since the  
 396 congruence status of all the objects are computed every time a new object is detected in the environment;  
 397 and secondly, broad enough to limit any bias possible in the interpretation of the perceived information: an  
 398 audiovisual class can be incongruent in an environment but congruent in another one, as will be illustrated  
 399 in Section 4. Moreover, by not creating any prior behavioral rules (such as *if-else* statements) and by  
 400 letting the system continuously being sensitive to new information, the DW module provides the robot  
 401 with a life-long learning of the environments composing the world it is living in. However, one important  
 402 limitation appears here: the DW module needs to have access to a *complete* audiovisual information in order  
 403 to compute the congruence of any object appearing in the scene. Indeed, in the situation where a source is  
 404 placed behind the robot, it would have to first turn its head towards it in order to get the full audiovisual  
 405 data, to then be able to take a decision on whether or not a head movement is necessary... which is what  
 406 can be called a *deadlock* situation. This is why the previous illustration of the DW module has used a setup  
 407 where all the visual data were always perceivable to the robot. Obviously, this is not a realistic context at  
 408 all. This is where the second module of the HTM system comes into play.

### 409 3.3 The Multimodal Fusion and Inference module

410 The *Multimodal Fusion & Inference* module (MFI module) is in charge of providing the DW module  
 411 with a complete information about the audiovisual sources, even when they are placed behind the robot.  
 412 Moreover, the MFI module is able to cope with classification errors, i.e. to provide a stable and reliable  
 413 estimation of the audiovisual classes of an object. This module is based on an online self-supervised active  
 414 learning paradigm that enables the overall system to create knowledge about the audiovisual classes that are  
 415 present in the environments the robot is exploring. Basically, the idea is to exploit head movements to learn  
 416 the relationship between the audio and visual classes of the sources, making the robot becoming afterwards  
 417 able to infer a missing modality. To begin with, the learning paradigm of the MFI module is described in a  
 418 first subsection. Then, the way motor orders are triggered to learn the association between audio and visual



**Figure 6.** Illustration of (a) the Multimodal Fusion & Inference module, (b) the Multimodal Self-Organizing Map. The M-SOM embeds one SOM per modality used for the definition of an object (audio and vision in our case). The representation here depicts the two subnetworks as a map containing neurons split in two parts defined by their own weights vectors, one part being dedicated to the mapping of audio data, the other to visual data.

419 classes is presented. An illustration of the MFI module functioning together with new details concerning  
420 the simulations, are then provided. A short discussion ends this MFI module presentation.

### 421 3.3.1 The Multimodal Self-Organizing Map

422 The MFI module is based on a Self-Organizing Map (SOM) Kohonen (1982) which is a learning algorithm  
423 relying upon a low dimensional map on which is performed a vector quantization of a high dimensional  
424 input matrix of data, while allowing its categorization. The input data are here made of classification  
425 experts outputs gathered in the vector  $P[t]$ , see Figure 6a and Equation (1). However, the traditional SOM  
426 algorithm shows one important limitation: it is unable to cope with missing data. In the case where an  
427 event originates from behind the robot, visual classifier outputs will not be relevant: the visual modality is  
428 missing. Then, two options can be chosen: (i) remove the corresponding visual components of  $P[t]$ , or (ii)  
429 set the corresponding components to the same arbitrary value. In the former case, this would imply a change  
430 in the data dimensionality. In the latter case, this would create arbitrary meaningful data which would be  
431 misinterpreted by the SOM. Then, these two options do not offer any solution to missing data inference.  
432 This is why it is proposed to transform a classical SOM into a *Multimodal-SOM* in order to keep what makes  
433 it powerful and usable with the constraints listed before. Interestingly, Papliński and Gustafsson (2005) have  
434 developed a bio-inspired system of interconnected SOMs allowing the learning of complex multimodal data  
435 for classification purpose. But while this system possesses interesting multimodal classification properties,  
436 it lacks the essential capability of inferring missing information. More recently, Bauer and Wermter (2013)  
437 and Schillaci et al. (2014) have proposed original models based on the SOM paradigm. But while they allow  
438 the multimodal learning of perceptual data in an unsupervised way, their major drawbacks reside either in  
439 their need of significant amount of data or in the time required to converge to a stable representation of the  
440 processed data.

441 **3.3.1.1 The subnetworks**

442 Lets recall that a SOM is a map composed of  $I \times J$  interconnected  $r_{ij}$  nodes, or neurons. The proposed  
 443 modification of the original SOM consists in creating one SOM per modality, as shown in Figure 6b.  
 444 Thus, the M-SOM is made of two (interdependent) maps, also composed by  $I \times J$  interconnected  
 445  $r_{ij}^{a/v}$  nodes, of size  $\lceil \sqrt{N_a \times N_v} \rceil \times \lceil \sqrt{N_a \times N_v} \rceil$  (where  $a/v$  stands for *audio or visual* in a compact  
 446 notation). This size has been selected to ensure that there will be at least one node available per possible  
 447 audiovisual class combination, given that no prior information is available about the plausible audiovisual  
 448 classes the robot will perceive during its life-long exploration. To each node is associated (i) a weights  
 449 vector  $\mathbf{w}_{ij}^a = (w_{ij}^a(1), \dots, w_{ij}^a(N_a))$  of size  $N_a$  for the audio subnetwork, and a weights vector  $\mathbf{w}_{ij}^v =$   
 450  $(w_{ij}^v(1), \dots, w_{ij}^v(N_v))$  of size  $N_v$  for the visual one, (ii) a  $(i, j)$  position in the map, and (iii) connections  
 451  $\chi_{(ij) \rightarrow (kl)}$  between the  $r_{ij}^{a/v}$  nodes and their neighbors in the same map, where  $[i, k] \in [1, I]$  and  $[j, l] \in$   
 452  $[1, J]$  (with an exception for the nodes located at the edges of the map where the connectivity is reduced).  
 453 The weights vectors  $\mathbf{w}_{ij}^{a/v}$  associated to all the  $r_{ij}^{a/v}$  nodes will become, through the iterative learning phase,  
 454 the representatives of the different kinds of vectors constituting the input matrix, and thus, of the different  
 455 audiovisual classes the input data capture.

456 **3.3.1.2 Weights update**

457 Traditionally, at every iteration  $n_{it}$  of the original SOM algorithm (the total number of iterations classically  
 458 going from thirty to thousands, given the complexity of the data to be processed), the input matrix is  
 459 parsed randomly until every vector has been processed once (Kohonen, 2013). For every vector explored  
 460 the algorithm looks then for the closest weights vector  $\mathbf{w}_{ij}$  associated to the node  $r_{ij}$  to the current input  
 461 vector, in terms of their Euclidean distance. The winning neuron, that is the one presenting the closest  
 462 distance to the input vector, is called the *Best Matching Unit* (BMU). It will be the location in the map  
 463 where the propagation of the resemblance between the input vector and the weights vector  $\mathbf{w}_{BMU}$  will start.  
 464 This propagation follows a Gaussian neighborhood function  $h_{ij}[n_{it}]$  (see Equation (14)) of variance  $\sigma[n_{it}]$   
 465 that defines the spread of the propagation. The neighborhood function is modulated by a factor  $\alpha[n_{it}]$ , the  
 466 learning rate, making the learning powerful in the first iterations but almost non-existent in the last ones.  
 467 Spreading the resemblance to the BMU's neighbors has two effects: (i) lowering the distance between the  
 468 BMU and the input vector so that this neuron becomes more and more the representative of the information  
 469 coded by this vector, and (ii) partially shaping the map around the BMU so that the closest to the BMU  
 470 in terms of distance, the closest also in terms of information coded by the input vector. This leads to an  
 471 self-organized map where regions have emerged, regions that code for similar categories. Once every vector  
 472 of the matrix has been explored, a new iteration of learning starts. At every iteration  $n_{it}$  is incremented  
 473 making  $\alpha[n_{it}]$  and  $\sigma[n_{it}]$  both decrease. Such decreasing leads to the following behavior of the learning  
 474 process: at start, the propagation spreads largely in the SOM and the learning rate is at its highest; at the  
 475 end of the learning, the propagation barely spreads around the BMU and the learning rate is at its lowest.

476 Within the M-SOM however, several changes of the traditional algorithm have been performed, changes  
 477 that impact the way weights are updated. First, an audiovisual BMU  $r_{BMU}^{av}$  is now computed as the  
 478 combination of the two (audio and visual) subnetworks, according to

$$r_{BMU}^{av} = r_{IJ}, \text{ with } (I, J) = \arg \min_{i,j} \left( \|\mathbf{P}^a - \mathbf{w}_{ij}^a\| \times \|\mathbf{P}^v - \mathbf{w}_{ij}^v\| \right), \quad (12)$$

479 where  $\|\cdot\|$  depicts the Euclidean distance between the vectors. This combined audiovisual BMU is associated  
 480 to the combined weights vector  $\mathbf{w}_{BMU}^{av} = (\mathbf{w}_{BMU}^a, \mathbf{w}_{BMU}^v)$ .

481 Secondly, the HTM does not have access to the whole matrix of data: the robot gets one vector at a time,  
 482 every time a frame is analyzed by the set of KSs in the architecture. Thus, the iterative process has been  
 483 revisited accordingly to this online paradigm. At every time step, the M-SOM will perform only 1 iteration  
 484 of learning with the current vector (that is, there is no infinite memory of the past perceived data). However,  
 485 the key principle of augmenting the resemblance between the BMU and the current vector, together with  
 486 its spread, must be kept in order to reach an organized map. Taking also into account the fact that the  
 487 audio classification experts from TWO!EARS get more and more precise the longer they gather data from  
 488 a same sound source, the Traditionally, at every iteration  $n_{it}$  of the original SOM algorithm (the total  
 489 number of iterations classically going from thirty to thousands, given the complexity of the data to be  
 490 processed), the input matrix is parsed randomly until every vector has been processed once (Kohonen,  
 491 2013). For every vector explored the algorithm looks then for the closest weights vector  $w_{ij}$  associated  
 492 to the node  $r_{ij}$  to the current input vector, in terms of their Euclidean distance. The winning neuron, that  
 493 is the one presenting the closest distance to the input vector, is called the *Best Matching Unit* (BMU). It  
 494 will be the location in the map where the propagation of the resemblance between the input vector and the  
 495 weights vector  $w_{BMU}$  will start. This propagation follows a Gaussian neighborhood function  $h_{ij}[n_{it}]$  (see  
 496 Equation (14)) of variance  $\sigma[n_{it}]$  that defines the spread of the propagation. The neighborhood function  
 497 is modulated by a factor  $\alpha[n_{it}]$ , the learning rate, making the learning powerful in the first iterations but  
 498 almost non-existent in the last ones. Spreading the resemblance to the BMU's neighbors has two effects: (i)  
 499 lowering the distance between the BMU and the input vector so that this neuron becomes more and more  
 500 the representative of the information coded by this vector, and (ii) partially shaping the map around the  
 501 BMU so that the closest to the BMU in terms of distance, the closest also in terms of information coded by  
 502 the input vector. This leads to an self-organized map where regions have emerged, regions that code for  
 503 similar categories. Once every vector of the matrix has been explored, a new iteration of learning starts.  
 504 At every iteration  $n_{it}$  is incremented making  $\alpha[n_{it}]$  and  $\sigma[n_{it}]$  both decrease. Such decreasing leads to the  
 505 following behavior of the learning process: at start, the propagation spreads largely in the SOM and the  
 506 learning rate is at its highest; at the end of the learning, the propagation barely spreads around the BMU  
 507 and the learning rate is at its lowest. Thirdly, still from the fact that the system does not have access to  
 508 the whole data to be processed, it is necessary to adapt how the algorithm converges. Since the robot will  
 509 always get to explore new environments during its life, there is no priorly known solutions to this learning  
 510 problem. Consequently, instead of trying to reach a global convergence of the overall M-SOM, the MFI  
 511 implements a *local consistency* (Chapelle et al., 2002; Zhou et al., 2004) at the audiovisual-class level (see  
 512 also Section 3.3.1.4). This local consistency enables the M-SOM to judge by itself whenever the learning  
 513 of a particular class can be stopped or has to be continued. Thus, the value of the iteration  $n_{it}$ , that will  
 514 have an impact on the values of  $\alpha$  and  $\sigma$ , will be computed *object by object*: every object has its own  
 515 iteration value corresponding to a certain degree in the learning process of the audiovisual class it belongs  
 516 to. The choice of implementing an iteration index object by object instead of class by class, which would  
 517 seem more logical, comes also from the potentially erroneous behavior of the classification experts during  
 518 the first perceived audio or visual frames associated to the objects (see 2.2). Indeed, relying directly on  
 519 these outputs could promote, on the mid- to long-term, the learning of false audiovisual classes that could  
 520 hamper the learning of the correct ones. The learning iteration  $n_{it}$  is now defined by

$$n_{it}[t] = \max((N_{it} - n_{it}^{o_j}[t]) + 1, 1) \text{ with } n_{it}^{o_j}[t] = t_{\text{init}}(o_j) + (t - t_{\text{init}}(o_j)), \quad (13)$$

521 where  $t_{\text{init}}(o_j)$  is the temporal index corresponding to the initial time the object emitted sound in the current  
 522 environment, and  $N_{it} = 10$  corresponds to the maximal number of iterations. The value of  $N_{it} = 10$   
 523 time steps has been defined experimentally with respect to two factors: (i) a too low value would put too

524 much importance on the very first frames detected by the classifiers for a given object, and (ii) a too high  
 525 value would significantly delay the local convergence of the learning for it would also delay the moment  $\alpha$   
 526 and  $\sigma$  would be high enough to make the learning actually efficient.

527 Once  $r_{\text{BMU}}^{av}$  is found, all the weights vectors associated with every node are then updated, as described  
 528 above, and according to

$$\mathbf{w}_{ij}^{a/v}[t+1] = \mathbf{w}_{ij}^{a/v}[t] + \alpha[n_{it}] h_{ij}[t, n_{it}] \|\mathbf{P}^{a/v}[t] - \mathbf{w}_{ij}^{a/v}[t]\|,$$

with  $h_{ij}[t, n_{it}] = \exp\left(-\frac{\|r_{\text{BMU}}^{av}[t] - r_{ij}\|^2}{2\sigma[n_{it}]^2}\right)$ , (14)

529 where  $\alpha \in [0.02, 0.9]$  represents the increasing learning rate (first and last values from (Kohonen, 1990),  
 530 and  $h_{i,j}[t, n_{it}] \rightarrow \mathbb{R}$  is the Gaussian neighborhood function of variance  $\sigma[n_{it}]$ .

### 531 3.3.1.3 Estimation of the audio and/or visual classes

532 Every time data  $\mathbf{P}[t]$  are available from the KS, the M-SOM proposes a corresponding estimated audio  
 533 and visual classes  $\hat{c}^a$  and  $\hat{c}^v$  respectively. In the case where all the data are available, then the corresponding  
 534 classes can be estimated along

$$\hat{c}^a = c_i^a, i = \arg \max_l w_{\text{BMU}}^a(l), \text{ and } \hat{c}^v = c_k^v, k = \arg \max_l w_{\text{BMU}}^v(l). \quad (15)$$

535 Thus, the audiovisual class  $\hat{c}^{\text{all}}$ , estimated when all the modalities are available, is given by  $\hat{c}^{\text{all}} = \{\hat{c}^a, \hat{c}^v\}$ .  
 536 All the interest of the M-SOM is its ability to provide both audio and visual classes, even if a part the KS  
 537 outputs are not available. Of course, no learning is then performed, but it is the step where the network is  
 538 actually exploited for inference. In the case where, for instance, the visual data are missing (i.e. the event is  
 539 out of the field of view of the robot), then:

- 540 1. audio only is exploited to determine the winning (audio) node  $r_{\text{BMU}}^a$  in the audio map, whose associated  
 541 weight vector  $\mathbf{w}_{\text{BMU}}^a$  can be used to estimate the audio class  $\hat{c}^a = c_i^a$ , with  $i = \arg \max_l w_{\text{BMU}}^a(l)$  like  
 542 in Equation (15);
- 543 2. the winning (visual) node is deduced from audio by  $r_{\text{BMU}}^v = r_{\text{BMU}}^a$ : this is exactly where the learned  
 544 interlink between audio and visual data is exploited. The corresponding visual class  $\hat{c}^v$  can then be  
 545 deduced from the associated weight vector  $\mathbf{w}_{\text{BMU}}^v$  along  $\hat{c}^v = c_k^v$ , with  $k = \arg \max_l w_{\text{BMU}}^v(l)$ .

546 In the end, the audiovisual class  $\hat{c}^{\text{miss}}$ , estimated when one modality is missing, is then given by  
 547  $\hat{c}^{\text{miss}} = \{\hat{c}^a, \hat{c}^v\}$ . Of course all the reasoning is identical when the other modality is missing: the available  
 548 data drive the missing modality for inference.

### 549 3.3.1.4 Convergence and the inference criterion

550 A key principle in learning algorithms is their ability to converge to one of the acceptable solutions of the  
 551 problem to be solved. However in the proposed context, different environments made of possibly different  
 552 audiovisual sources might be explored during the robot life. Then, it is clearly impossible to define one  
 553 global good solution to the problem. Nevertheless, the proposed M-SOM possesses a characteristic of  
 554 *local consistency* (see Section 3.3.1.2). Within the classical SOM algorithm, convergence means that the  
 555 whole map is organized such that the different nodes are grouped in meaningful entities that code part  
 556 of the input data. In the proposed M-SOM, it is proposed that the algorithm always keeps a free part in  
 557 the map, i.e. nodes not coding for any audio or visual classes. This would allow the network to include

558 new audiovisual classes, discovered all along the interaction with new environments during the robot life.  
 559 Looking for local consistencies, rather than reaching for global convergence, is implemented through the  
 560 definition of a criterion for each audiovisual category already created, indicating how much this category  
 561 has been learned so far and if its learning can be stopped. The multimodal learning performed by the MFI  
 562 module is supported by head rotations to the sources to be learned. It allows to bring the visual sensors in  
 563 front of them in order to learn the association between the corresponding audio and visual classes. But  
 564 these head movements are no longer useful once the M-SOM has enough knowledge about the audiovisual  
 565 classes, thus justifying the need to (i) inhibit these head movements, and (ii) being able to judge when this  
 566 amount of knowledge is sufficient. Then, an *inference ratio*  $q(\mathcal{C}^{(l)}(c_i^a, c_k^v))$  for the audiovisual category  
 567  $\mathcal{C}^{(l)}(c_i^a, c_k^v)$  is defined as

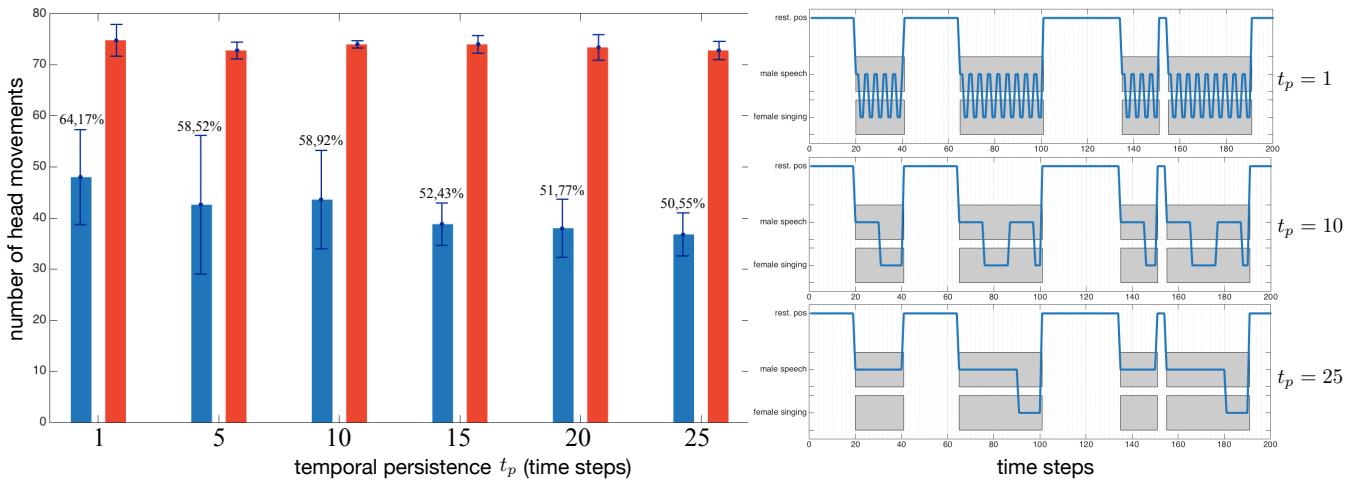
$$q\left(\mathcal{C}^{(l)}(c_i^a, c_k^v)\right) = \frac{\sum_{n=1}^{n=t} \delta_{i,k}^{\text{miss}}[n-1] \delta_{i,k}^{\text{all}}[n]}{\sum_{n=1}^{n=t} \delta_{i,k}^{\text{miss}}[n]}, \text{ with } \delta_{i,k}^{\text{all/miss}} = \begin{cases} 1 & \text{if } \widehat{c}_{i,k}^{\text{all/miss}}(o_j) = \{c_i^a, c_k^v\}, \\ 0 & \text{else.} \end{cases} \quad (16)$$

568 This inference ratio is computed by comparing the number of times the audiovisual category  $\mathcal{C}^{(l)}(c_i^a, c_k^v)$   
 569 has been obtained (or inferred) with one missing modality ( $\delta_{i,k}^{(\text{miss})} = 1$ ) at time  $n-1$  and confirmed at  
 570 time  $n$  ( $\delta_{i,k}^{(\text{all})} = 1$ ) by a head movement with all modalities available, with the total number of inference.  
 571 Thus,  $q(\mathcal{C}^{(l)}(c_i^a, c_k^v))$  captures the ability of the MFI module to infer correctly a missing modality, category  
 572 by category. The inference ratio always lies between 0 and 1, where 1 means that the category has always  
 573 been perfectly inferred. On this basis,  $q(\mathcal{C}^{(l)}(c_i^a, c_k^v))$  is compared to a criterion  $K_q \in \mathcal{R}^+ = [0, 1]$ : if  
 574 a modality is missing, the MFI module will attempt to infer it, and as long as the inference ratio of the  
 575 corresponding audiovisual category is lower than  $K_q$ , a head movement will be triggered towards the  
 576 corresponding source. Thus, the system grabs the missing information and feeds the M-SOM, which can  
 577 then learn the audiovisual association. Of course, once the full audiovisual data is obtained, a comparison  
 578 with the previous inference is made and the inference ratio is updated accordingly. If the inference ratio  
 579 gets higher than the criterion  $K_q$ , the learning is considered as being good enough to trust the inference  
 580 made by the MFI module, and inhibit consequent head movements towards the sources belonging to the  
 581 corresponding audiovisual category. Remark that the criterion  $K_q$  has an influence on the behavior of  
 582 the MFI module (Cohen-Lhyver, 2017). A low threshold allows a quick confidence in the inference, thus  
 583 freeing head movements for other tasks, whereas a high  $K_q$  value pushes the system to be very careful  
 584 about its inferences.

### 585 3.3.2 Motor orders

586 As for the DW module, the MFI module is able to trigger head movements towards sources of *interest*.  
 587 This interest is now formalized by the lack of confidence in the knowledge of the audiovisual category a  
 588 source might belong to. As previously explained, turning the head towards a source might enable the visual  
 589 sensors to get the missing visual data, thus giving to the MFI module the opportunity to learn the interlink  
 590 between the audio and visual modalities, but also to eventually confirm/refute an inference. Like for the DW  
 591 module, the head movements modulation is inspired by the GPR model (see Section 3.2.2), but through a  
 592 different expression of the activities  $\tau_{\text{MFI}}(o_j)$  for the object  $o_j$  with audiovisual category  $\mathcal{C}^{(l)}(c_i^a, c_k^v)$ , now  
 593 given by

$$\tau_{\text{MFI}}(o_j) = \frac{q(\mathcal{C}^{(l)}(c_i^a, c_k^v))}{K_q} \times \delta^{(i,k)}(n), \text{ with } \delta^{(i,k)}(n) = \begin{cases} -1 & \text{if } n < (t_p = 10), \\ 1 & \text{else,} \end{cases} \quad (17)$$



**Figure 7.** Impact of the temporal persistence, introduced in Equation (17), (left) on the number of triggered head movements in a complex environment, and (right) the behavior of the robot in a simplified case for illustration purposes. (Blue bars:) robot driven by the MFI module, (red bars) naive robot. Percentages depict the ratio between the naive robot and the MFI module.

594 where  $n = t - t_{\text{foc}}(o_j)$ , with  $t_{\text{foc}}(o_j)$  the first time the object has been focused on by the MFI module. Then,  
 595 the angle  $\hat{\theta}(o_l)$  estimated by the localization expert and corresponding to the canal with the lowest activity  
 596 is selected as the winning motor order  $\theta_{\text{MFI}}$ , i.e.

$$\theta_{\text{MFI}} = \hat{\theta}(o_l), \text{ with } l = \arg \min_j (\tau_{\text{MFI}}(o_j)). \quad (18)$$

597 The term  $\delta^{(i,k)}(n)$  in (17) introduces a form of temporal persistence through a positive feedback loop,  
 598 as observed in the thalamus by Redgrave et al. (1999); Gurney et al. (2001a); Meyer et al. (2005). The  
 599 value of  $t_p = 10$  has been set experimentally after several comparisons and evaluations. The impact of  
 600 this persistence in a complex environment (eight sources with five simultaneously emitting) is illustrated  
 601 in the left panel of Figure 7, where the blue bars depict the number of head movements triggered by  
 602 the MFI module, while the red bars, by the naive robot  $\mathfrak{R}_n$  (these numbers are obviously not affected by  
 603 the temporal persistence applied to the MFI module). The main point is that the temporal persistence  $t_p$   
 604 constitutes only a small part of the head movements control: 13.6% less head movements between  $t_p = 1$   
 605 and  $t_p = 25$ . The real benefits of temporal persistence is shown in Figure 7 (right): with  $t_p = 1$ , the robot  
 606 exhibits oscillations between two sources, potentially damaging the internal representation of the world  
 607 (confusions in binaural cues computations, speed of the movement...). With  $t_p = 25$ , a pervert effect of a  
 608 too long persistence is also shown: the system often ghosts completely the (SINGING, female) source,  
 609 preventing itself from learning it.

### 610 3.3.3 Evaluation 2: classification rates of the MFI module

611 The MFI module aims providing a corrected audiovisual information from the classification experts. In  
 612 order to assess the contribution brought by this module, a good audiovisual classification rate  $\Gamma(o_j)[t]$  is  
 613 defined by comparing the audio and visual classes associated to all the objects detected by the system with  
 614 the ground truth, according to

$$\Gamma(o_j)[t] = \mathbf{a} \times \sum_{k=t_i}^t \gamma(o_j)[k] \text{ with } \gamma(o_j)[k] = \begin{cases} 1 & \text{if } \hat{\mathbf{c}}(o_j)[k] = \mathbf{c}(\Psi_j)[k], \\ 0 & \text{else,} \end{cases} \quad (19)$$

**Table 1.** Simulation setup for Evaluation 2. Each setup is repeated 5 times for a total of 100 simulations.

Evaluation 2						
$e^{(l)}$	$n_{\mathcal{S}}$	$n_{sim}^{max}$	$T$	$ \mathcal{C}^{(l)} $	$K_q$	$\varepsilon_{\mathcal{P}}$
1 to 5	3	3	1000	2	0.8	0.1, 0.3, 0.5, 0.7, 0.9
6 to 10	5	5	1000	3	0.8	0.1, 0.3, 0.5, 0.7, 0.9
11 to 15	7	7	1000	4	0.8	0.1, 0.3, 0.5, 0.7, 0.9
16 to 20	10	10	1000	6	0.8	0.1, 0.3, 0.5, 0.7, 0.9

615 with  $c(\Psi_j)[k]$  being the ground truth audiovisual class of the event  $\Psi_j$  captured as the object  $o_j$  at time  $k$   
616 in the internal representation, and  $a = 1/[1, \dots, (t - t_i) + 1]$  corresponding to the elapsed time between  
617 the first time step  $t_i$  when the MFI module provided a classification of the object  $o_j$ , and the current time  $t$ .  
618 The overall good classification rate is given by applying a sliding window on all  $\Gamma(o_j)$  computed from the  
619 beginning of the exploration, along

$$\bar{\Gamma}_{MFI}[t] = \frac{1}{N_{obj}^c[t]} \sum_{j=1}^{N_{obj}^c[t]} \Gamma(o_j)[t] \quad (20)$$

620 with  $N_{obj}^c[t]$  the number of processed objects by the MFI module at time  $t$  (this number could be inferior or  
621 equal to the total number of objects present and emitting, noted  $N_{obj}$ ). In parallel, the same process is made  
622 for the naive robot  $\mathfrak{R}_n$  (knowing that this one performs the fusion of the classification experts themselves  
623 through a maximum a posteriori approach, along Equation (11)), according to

$$\bar{\Gamma}_{\mathfrak{R}_n}[t] = \frac{1}{N_{obj}[t]} \sum_{j=1}^{N_{obj}[t]} \Gamma(o_j)[t]. \quad (21)$$

624 In addition, a measure of the classification performance of an omniscient (thus unrealistic) robot is also  
625 computed, noted  $\bar{\Gamma}'_{\mathfrak{R}_n}[t]$ . This robot has full access to every auditory and visual information, even when  
626 the objects are out of the sight of the robot. The simulation setup is presented in Table 1. Twenty different  
627 multisource environments are simulated, each of them in possibly different conditions (number of sources,  
628 number of simultaneously emitting sources, error rates, etc.). The resulting good audiovisual classification  
629 rates are regrouped in Table 2, mainly organized by increasing error rates  $\varepsilon_{\mathcal{P}} = \{0.1, 0.3, 0.5, 0.7, 0.9\}$ .

630 At first, let us consider the naive omniscient robot  $\mathfrak{R}'_n$ . As expected, it presents a mean good audiovisual  
631 classification rate  $\bar{\Gamma}_{\mathfrak{R}'_n}$  almost equal to  $1 - \varepsilon_{\mathcal{P}}$  for all tested conditions. In contrast, the realistic naive robot  
632  $\mathfrak{R}_n$  (having only access to the data it is able to perceive) systematically exhibits lower rates  $\bar{\Gamma}_{\mathfrak{R}_n}$ . Clearly,  
633 the main flaw of this robot is its incapacity to perform any inference, which turns to be a critical capability  
634 in multisource environments. In comparison, the proposed MFI module outperforms both naive robots, for  
635 almost any error rates and number of sources (except for only one case:  $\varepsilon_{\mathcal{P}} = 0.1$  and  $n_{\mathcal{S}} = 10$ ). The last  
636 column in Table 2 exhibits the ratio between the best naive robot  $\mathfrak{R}_n$  (given by  $\bar{\Gamma}'_{\mathfrak{R}_n}[t = T]$ ) and the MFI  
637 module: the greater  $\varepsilon_{\mathcal{P}}$ , the higher the ratio, except with  $\varepsilon_{\mathcal{P}} = 0.9$ . In this case, the error rate is anyway so  
638 high that the interest in exploiting such corrupted data is almost null. However, even in very challenging  
639 conditions involving a very high  $\varepsilon_{\mathcal{P}} = 0.7$  in a multisource context, the MFI module provides on average a  
640 2.4 times better good audiovisual classification rate than with the classifier outputs.

**Table 2.** Good classification rates for different error rates and different numbers of sources. Every results is an average of 5 repetitions of every conditions with standard deviation in parentheses, for a total of 100 simulations.  $\mathfrak{R}_n$  corresponds to the naive robot, and  $\mathfrak{R}'_n$  to the unrealistic omniscient robot. Values are rounded up to the third decimal.

Evaluation 2: Results					
$\varepsilon_{\mathcal{P}}$	$n_{\mathcal{S}} - n_{sim}^{max}$	$\bar{\Gamma}_{MFI}[t = T]$	$\bar{\Gamma}'_{\mathfrak{R}_n}[t = T]$	$\bar{\Gamma}_{\mathfrak{R}_n}[t = T]$	ratio: $\frac{\bar{\Gamma}_{MFI}[t = T]}{\bar{\Gamma}'_{\mathfrak{R}_n}[t = T]}$
0.1	3 — 3	0.982 (0.027)	0.894 (0.021)	0.503 (0.073)	1.098
	5 — 5	0.988 (0.025)	0.899 (0.012)	0.339 (0.039)	1.099
	7 — 7	0.960 (0.023)	0.893 (0.016)	0.264 (0.021)	1.075
	10 — 10	0.866 (0.047)	0.887 (0.018)	0.182 (0.014)	0.976
	<b>mean</b>	<b>0.949</b>	<b>0.893</b>	<b>0.322</b>	<b>1.063</b>
0.3	3 — 3	0.992 (0.020)	0.703 (0.042)	0.414 (0.055)	1.411
	5 — 5	0.987 (0.022)	0.692 (0.017)	0.265 (0.014)	1.426
	7 — 7	0.942 (0.028)	0.691 (0.014)	0.198 (0.017)	1.363
	10 — 10	0.883 (0.041)	0.689 (0.011)	0.145 (0.014)	1.281
	<b>mean</b>	<b>0.951</b>	<b>0.693</b>	<b>0.255</b>	<b>1.372</b>
0.5	3 — 3	0.973 (0.026)	0.493 (0.020)	0.280 (0.031)	1.973
	5 — 5	0.965 (0.043)	0.496 (0.021)	0.189 (0.034)	1.945
	7 — 7	0.899 (0.048)	0.492 (0.018)	0.145 (0.019)	1.827
	10 — 10	0.836 (0.042)	0.492 (0.018)	0.103 (0.010)	1.699
	<b>mean</b>	<b>0.918</b>	<b>0.493</b>	<b>0.179</b>	<b>1.862</b>
0.7	3 — 3	0.774 (0.087)	0.282 (0.030)	0.165 (0.028)	2.744
	5 — 5	0.737 (0.105)	0.294 (0.014)	0.120 (0.023)	2.506
	7 — 7	0.683 (0.133)	0.296 (0.016)	0.081 (0.012)	2.307
	10 — 10	0.550 (0.117)	0.293 (0.016)	0.064 (0.011)	1.877
	<b>mean</b>	<b>0.686</b>	<b>0.291</b>	<b>0.107</b>	<b>2.357</b>
0.9	3 — 3	0.213 (0.060)	0.092 (0.019)	0.054 (0.019)	2.315
	5 — 5	0.152 (0.064)	0.102 (0.012)	0.039 (0.007)	1.490
	7 — 7	0.174 (0.075)	0.100 (0.009)	0.031 (0.005)	1.740
	10 — 10	0.140 (0.066)	0.100 (0.009)	0.019 (0.006)	1.400
	<b>mean</b>	<b>0.169</b>	<b>0.098</b>	<b>0.035</b>	<b>1.724</b>

### 641 3.3.4 Discussion

642 The proposed MFI module, mainly based on the M-SOM, provides an online self-supervised active  
 643 learning paradigm to be able to process erroneous and/or missing data in the particular context of the  
 644 exploration of unknown environments. The overall goal of the MFI module is thus to feed the DW module  
 645 with correct audiovisual classes the perceived objects belong to, with respect to a very short learning  
 646 time constraint (down to a few seconds only). The *active* capabilities of the MFI module is of very much  
 647 importance here, for it enables the intensive use of head movements to gather, whenever it is necessary,  
 648 and in real-time, additional data to refine the knowledge the module has of the world under exploration. A  
 649 fundamental question arises with the problem of audio and visual classes association when considering  
 650 one-to-one audiovisual pairs, i.e. that each audio label is associated with only one visual label, and vice  
 651 and versa. In the evaluations presented in this section, such pairing limitation was not used: an audio label  
 652 could have several visual correspondences, such as SPEAKING, male, SPEAKING, female, or SPEAKING,  
 653 child. However, given these audiovisual labels examples, it is not possible for the MFI module to create

654 an information that does not exist: from the audio label SPEAKING, it is impossible to determine whether  
 655 the corresponding visual label is male, female, or child. The MFI module still outputs an hypothesis  
 656 corresponding, given how the M-SOM learning algorithm works, to the most observed so far audiovisual  
 657 pair. Such limitation of the MFI module only comes from the limits of the classification experts themselves:  
 658 if the classifiers cannot distinguish a female voice from a male one, nor would the MFI module. Such a case  
 659 will be shown and also discussed in Section 4.4, when evaluating the whole system in real environments.

### 660 3.4 Conclusion

661 The Head Turning Modulation system is composed with two modules: the Dynamic Weighting module  
 662 (DW) and the Multimodal Fusion & Inference module (MFI), each of them having been described in  
 663 this section. The DW module is an attentional component, working on the sole basis of observed data in  
 664 unknown environments, from which it enables the robot to turn its head to audiovisual sources considered  
 665 as “of importance”. Coupled to it is the MFI module that learns the relationship between the modalities  
 666 that are used to define an object (audition and vision in this case). Based on a Multimodal Self-Organizing  
 667 Map (M-SOM), the MFI module is able to create the knowledge required by the DW module to work  
 668 properly. This knowledge consists in the fusion of multimodal data into a corrected database of audiovisual  
 669 categories, knowledge that is created through online active self-supervised exploration of the audiovisual  
 670 sources appearing in the unknown environments. Both modules can trigger head movements independently,  
 671 and their combination necessitates an adaptation of the motor orders expressions of the modules.

672 The next section will present the results obtained in real environments with the real robot embedding the  
 673 whole TWO!EARS software (including the integration of the HTM system), and processing real audio and  
 674 visual data.

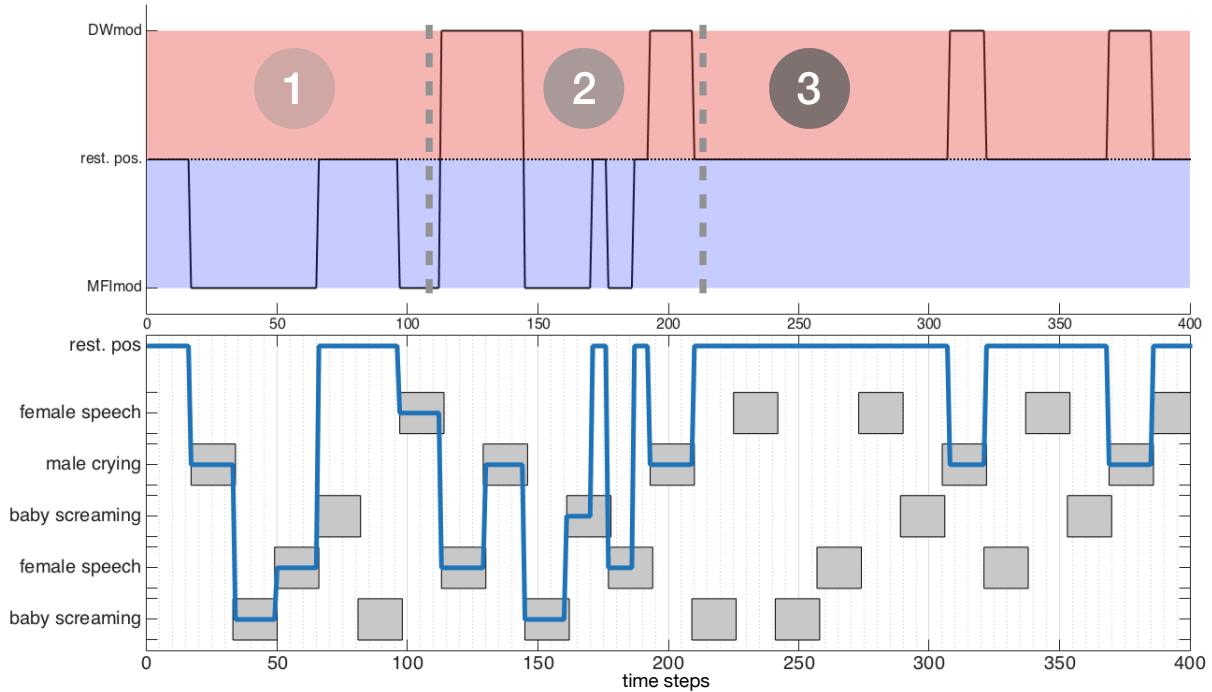
## 4 COMBINATION OF THE TWO MODULES

675 The previous section was dedicated to the individual presentation of each module constituting the HTM  
 676 system, while providing limited evaluations in simulated conditions. This section is now concerned with the  
 677 combination of the DW module and the MFI module together, with their evaluation in realistic conditions,  
 678 i.e. on a real robot and with real audio and visual data. At first, one have to deal with the fact that theses two  
 679 modules are both able to generate competitive head movements. The way they are prioritized is described  
 680 in a first subsection. Next, the experimental setup is carefully described in a second subsection. Then,  
 681 experimental results are provided in a third subsection, aiming at demonstrating the benefits of the overall  
 682 system in the audiovisual scene understanding.

### 683 4.1 Combined motor orders: Evaluation 3

684 It has been shown in Section 3 that the DW module and the MFI module both exploit head movements to  
 685 better their respective operations. Trying to make them able to work together then requires a prioritization  
 686 of them. On the one hand, the DW module provides the robot with potential sources to be focused on, on  
 687 the basis of their computed congruence; on the other hand, the MFI module aims at estimating audiovisual  
 688 classes of objects inside the environment, even with potential classification errors and lack of data. It seems  
 689 then obvious to set the priority to the MFI module: having a reliable audiovisual classes estimation system is  
 690 required for the attentional module to take relevant decisions. This prioritization introduces a new activity  
 691  $\tau'_{\text{DW}}$  for the DW module which is now defined, for an object  $o_j$ , by

$$\tau'_{\text{DW}}(o_j) = \tau_{\text{MFI}}(o_j) - \tau_{\text{DW}}(o_j) \times \delta(\tau_{\text{MFI}}(o_j)), \text{ with } \delta(x) = \begin{cases} 1, & \text{if } x \geq 1, \\ 0, & \text{otherwise.} \end{cases} \quad (22)$$



**Figure 8.** Behavior of the combined modules in three phases. The testing environment is composed of five audiovisual sources and is willingly simple for illustration purposes. (Top panel:) module ordering the head movement. (Bottom panel:) temporal course of the exploration of the environment; gray boxes depict the temporal course of emitting sources.

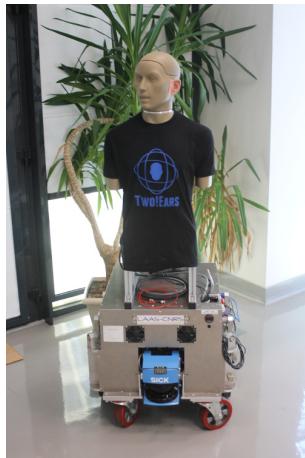
692 On this basis, the motor order  $\theta_{\text{HTM}}$  selected to drive the head is computed along

$$\theta_{\text{HTM}} = \widehat{\theta}(o_l), \text{ with } l = \arg_{j_1, j_2} \min(\tau'_{\text{DW}}(o_{j_1}), \tau_{\text{MFI}}(o_{j_2})) \text{ where } \begin{cases} j_1 = \arg \min_l (\tau'_{\text{DW}}(o_l)), \\ j_2 = \arg \min_k (\tau_{\text{MFI}}(o_k)), \end{cases} \quad (23)$$

693 i.e. the object with the lowest DW module or MFI module activity is selected. Such a modification of the  
 694 motor activity expression enables the MFI module to take over the lead on the DW module. The evaluation  
 695 of such a modification in the motor commands decision system can be performed again in simulation,  
 696 along the same procedure as in the previous simulations, see Figure 8. Let us consider an environment  
 697 made of five objects, belonging to three different audiovisual categories. Each of these objects emit sounds  
 698 along time, according to the time plot shown in Figure 8 (bottom). Figure 8 (top) exhibits the three-phase  
 699 behavior of the motor decision algorithm. At the very beginning, only the MFI module is responsible for  
 700 the head movements: the system is learning the association between audiovisual classes. Little by little, the  
 701 inference provided by the MFI module does not need motor confirmation for some of the classes: the DW  
 702 module can now compute congruence of the corresponding objects and potentially trigger head movements.  
 703 In the end, all the audiovisual classes are correctly learned by the MFI module, letting the sole DW module  
 704 in charge with head rotations. Of course, the head movements triggered by the DW module are also used to  
 705 feed the M-SOM.

## 706 4.2 Experimental setup and data generation

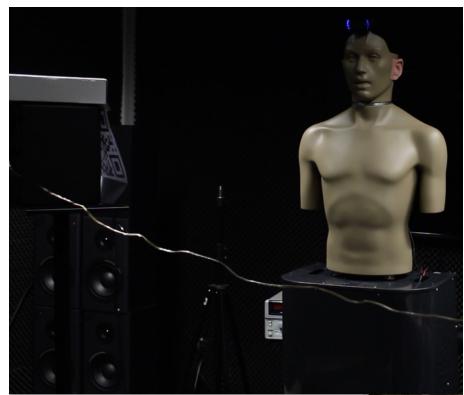
707 The overall system has been evaluated in a realistic environment by using a real robot integrating the  
 708 whole Two!EARS software and evolving in a real room. In practice, two different robots have been  
 709 actually used: one mobile platform from LAAS-CNRS (Toulouse, France) named JIDO, the other one



(9a) The Jido platform.



(9b) The ODI platform.



(9c) ODI, facing a loudspeaker with a QR code attached on it.

**Figure 9.** The two robotic platforms used in the project, both supporting a motorized KEMAR HATS. ODI has been used in this paper for the HTM evaluation.

from ISIR (Paris, France) named ODI, see pictures in Figure 9. Both platforms support a KEMAR HATS (Head And Torso Simulator), whose necks have been motorized to control their head movements in azimuth (Bustamante et al., 2016). A HATS is a manikin endowed with two microphones placed inside two pinnae which mimics the acoustic effect of the head (and torso) on the left and right ear signals, thus producing a realistic binaural information, close to what a human could actually hears. The servo control of the head is ensured by a set including a motor, its gear head, an encoder, and an Harmonica electronic controller from ELMO, mounted inside the HATS. A ROS node dedicated to the head control is in charge of controlling this motorization, allowing real-time servoing of the head movements by using possibly different feedback control options like position or velocity setpoints. In this paper, the positions deduced from Equation (23) are directly sent to the ROS node to control the head in position. These two robots are very much alike, except for vision: the one used at ISIR for the experiments used in this paper is only endowed with monocular vision. However, as already argued, the HTM system is not dependent on the way each modality works, but only on the identification experts, be they dedicated to monocular or binocular vision for instance.

Everything related to the platform and data acquisition is handled by the ROS middleware, running directly on the robot: navigation, obstacle avoidance, image and audio captures, etc. Note that a dedicated ROS binaural processing node has been developed during the project, so that most of the audio cues required for sound localization, recognition and separation are directly computed in real-time on the robot. State-of-the-art ROS nodes dedicated to vision (acquisition and processing) have also been used. All the data computed on the robot are then transmitted to another computer running the Two!EARS framework thanks to a MATLAB-to-ROS bridge. This bridge has been entirely designed to deal with the proposed bottom-up and top-down approach of the project, so that all the ROS nodes can be easily parameterized on the fly and in real time. Then, all the steps required for the “cognitive” analysis (i.e. object localization, recognition, fusion, etc.) runs under MATLAB.

Experiments used in this paper have been conducted in a pseudo-anechoic room populated with loudspeakers over which QR codes have been attached to, see Figure 9c. These are used by a ROS node to extract the visual labels of each object directly and with a recognition rate similar to the one obtained through the binocular vision of JIDO with the Line-Mod algorithm Hinterstoisser et al. (2012).

All the sounds emitted from the loudspeakers belong to a database constituted of sounds used to train the audio experts in recognition. In other terms, all the sounds can be recognized by at least one expert in the architecture. Then, the HTM has been evaluated in experimental conditions by two scenarios: the first emphasizes the global behavior of the system, while the second focuses on the fusion and classification abilities of the MFI module. Whatever the scenarios, they all work the following way: sounds are emitted from one or multiple loudspeakers, possibly at the same time. Depending on how the head of the KEMAR is turned, some QR codes can be manually changed from one loudspeaker to another to simulate an object movement in the environment. The HTM system then gathers classification and localization results coming from the audio and visual experts, and triggers some head movements accordingly. A scenario is entirely described by the number of different objects in the scene and by the time description of their localization, appearance and disappearance, exactly like in the previous simulations. Of course, ground truth audio and visual classes of each object are known, thus allowing a careful evaluation of the overall system performance. Note that the audio experts used in the following experiments have been set up by using data from a database recorded in a different acoustic environment. Since they all rely on a prior learning step exploiting these data, there will be a mismatch between their learning and testing phase. The main consequence is mainly a lower frame recognition rate, evaluated to about 37% for the four classifiers used here, and that have been chosen amongst the most performing ones (Two!Ears, 2016a). The same applies to the localization algorithm, with less consequences: experiments still show a good ability to localize sounds with a precision of about 7.7° (including front-back confusion). Finally, the visual recognition of QR codes works almost perfectly, while being quite sensitive to changes in illuminations. Of course, both phenomena are dealt with the HTM system, which has been entirely designed to cope with recognition errors and lack of data, as shown in the next subsections.

#### 4.3 Evaluation 4: global behavior

This first evaluation aims at demonstrating how the two modules constituting the HTM system cooperate together in order to give the exploratory robot an additional understanding of the world. The evaluation consists in presenting to the system three successive environments made of three to four objects, as summarized in Table 3. The audiovisual sources of the environments are placed around the robot and emit sound intermittently, according to the time scenario shown in Figure 10 (bottom). Exactly like in simulations, the real robot is compared to its naive counterpart  $\mathfrak{R}_n$ , turning its head towards every audiovisual events regardless of their meaning. To begin with, the HTM builds a first representation  $e^{(1)}$ . As shown in Figure 10, the robot starts by turning its head towards the first two audiovisual sources (BARKING, dog and SPEECH, female), driven by the MFI module since these audiovisual classes are brand new to it. As already outlined in the previous subsection, the HTM tries to learn the audiovisual association between these two classes. This learning is done very quickly: one can observe at time index  $t = 28$  (corresponding to the “real” time 14s) that the robot turns its head to its resting state (blue line going at the top of the figure), meaning that neither the DW module nor the MFI module requires a head movement towards the sources BARKING, dog: these sources are not of interest anymore, and hearing the sound BARKING is sufficient to infer the visual class dog. Nevertheless, one can remark a glitch in the head movement decision at  $t = 30$ , as the last attempt of the MFI module to learn the BARKING, dog audiovisual association. At  $t = 41$  (20.5s), the robot turns its head again towards the source (SPEECH, female): with two BARKING, dog for one SPEECH, female in  $e^{(1)}$ , the probability for this last audiovisual category  $p(C^{(1)}(\text{SPEECH, female})) = 1/3$  falls below  $K^{(1)} = 1/2$ , thus making any object of this audiovisual category incongruent. Then, the robot explores a second environment. It is similar in terms of frequencies of apparition of each audiovisual categories, even if their meaning (at least, to us) is different: the two BARKING, dog are trade for two

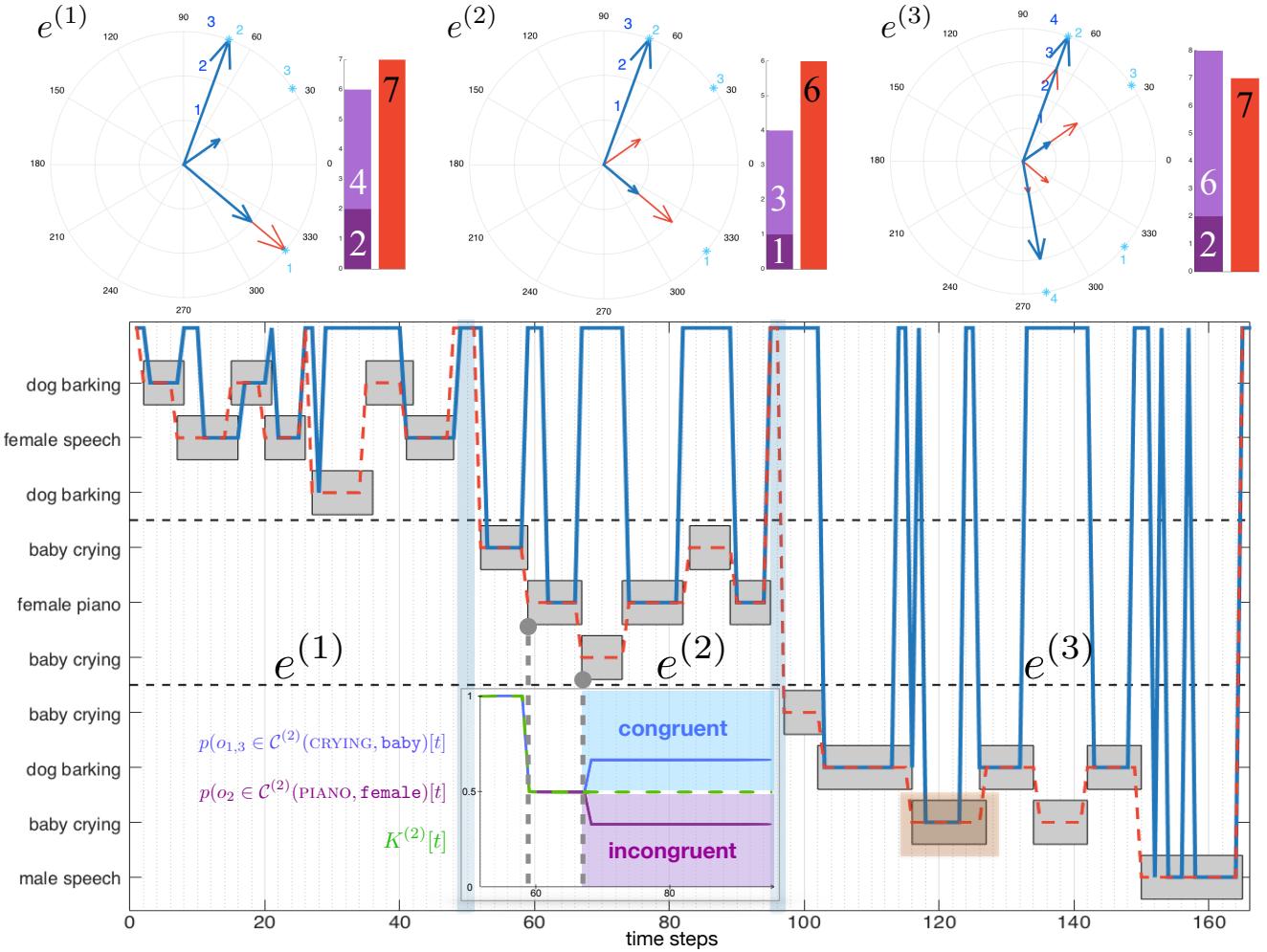
**Table 3.** Experimental setup for Evaluations 4 & 5.

Evaluation 4					
$e^{(i)}$	$n_{\mathcal{S}}$	$n_{sim}^{max}$	$c(\Psi_j)$	$\theta(\Psi_j)$	$K_q$
1	3	1	<i>dog barking n°1</i>	$320^\circ$	0.6
			<i>dog barking n°2</i>	$35^\circ$	
			<i>female speech</i>	$70^\circ$	
2	3	1	<i>baby crying n°1</i>	$70^\circ$	0.6
			<i>baby crying n°2</i>	$35^\circ$	
			<i>female piano</i>	$320^\circ$	
3	4	1	<i>baby crying n°1</i>	$70^\circ$	0.6
			<i>baby crying n°2</i>	$35^\circ$	
			<i>dog barking</i>	$320^\circ$	
			<i>male speech</i>	$280^\circ$	

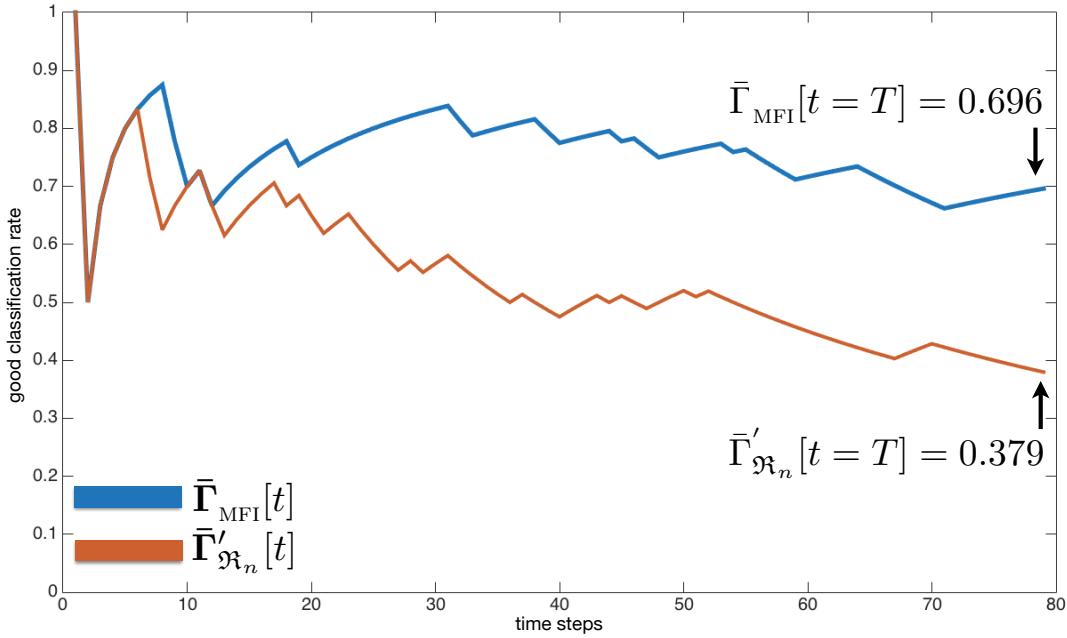
Evaluation 5					
1	5	1	<i>female speech</i> <i>female piano</i> <i>male speech</i> <i>dog barking</i> <i>baby screaming</i>	$320^\circ$ $30^\circ$ $60^\circ$ $90^\circ$ $280^\circ$	0.6

782 CRYING, baby, while the category SPEECH, female is replaced by PIANO, female. Logically, the  
 783 obtained behavior is similar: a quick learning of the audiovisual association allows then the head to be  
 784 controlled by the DW module on the basis on congruency computations. Interestingly, the understanding  
 785 of this second environment by the DW module could appear as counterintuitive in comparison with how  
 786 humans might have reacted by favoring the two objects CRYING, baby. This more social reaction could  
 787 nevertheless be handled by some additional KS from the TWO!EARS architecture which could modulate the  
 788 overall reaction of the robot w.r.t. the current task (Ferreira and Dias, 2014). Finally, a third environment is  
 789 explored. It will allow to demonstrate the benefits of reusing information between the representation of  
 790 environments, see Section3.2.1. Indeed, the scene begins with a CRYING, baby which does not trigger  
 791 any head movement: while being in a new environment, the HTM system considers at this point that this  
 792 third environment is very likely to be the same as  $e^{(2)}$  where this audiovisual class was considered as  
 793 congruent. Consequently, the congruence computations of each audiovisual categories in the previous  
 794 environment can still be used, and no head movements towards this now object is performed. However,  
 795 as soon as a new object eliminates the possibility to be in an environment similar to  $e^{(2)}$  pops up, a new  
 796 representation  $e^{(3)}$  is created. Thus, when the source BARKING, dog appears in the scene, a head movement  
 797 is immediately triggered towards it, since it is incongruent in  $e^{(3)}$ . Once again, a small glitch in the motor  
 798 decision appears in  $t = 116$ , caused by the experts outputs and the signal non-stationarity (a BARKING  
 799 sound includes indeed some silence). The movement triggered at  $t = 118$  is an error from the system  
 800 since the object CRYING, baby should have been considered as congruent. The audio data perceived at  
 801 this time is MALESPEECH, data from a never encountered audio class, thus enjoining the MFI module  
 802 to trigger a head movement. From  $t = 119$  and  $t = 122$ , the experts' data changed and became of class  
 803 CRYING, dog, an audiovisual pair the MFI module never encountered before, consequently still promoting  
 804 the focus on the object. However, at time  $t = 123$ , the correction of the MFI module has been applied  
 805 and the “correct” audiovisual class CRYING, baby is now output by the module. The DW module, in  
 806 response, analyses it and consider it as congruent in this environment, thus inhibiting the head movement.  
 807 Finally, the new source SPEECH, male appears in the environment at  $t = 150$  and the robot is focused on  
 808 it. Two (apparently) erroneous movements to the resting position can be observed, at  $t = 153$  and  $t = 157$ ,



**Figure 10.** (Top panel:) Number of head movements triggered during the exploration of each environment by (blue) the HTMKS, (red) the virtual naive robot. Each arrow points at a source and their length represent the number of movements. The light blue numbers correspond to the position of the sources. (Purple bars:) total number of movements triggered by (dark) the MFI module, (light) the DW module (black numbers are their sum). (Red bars:) number of movements triggered by the virtual naive robot. (Bottom panel:) movements triggered by (blue line) the HTMKS, and (dotted red line) the naive robot. (Grey boxes:) temporal course of the scenarios. The semi-transparent red box at  $t=116$  highlights the significant wrongful discrepancy that occurred between the actual audiovisual class of the object and the perception of the HTM (error that is corrected soon after, see text for more details). Additionally, the subfigure present in the delineated box at the bottom of  $e^{(2)}$  represents the evolution of  $K^{(2)}$  together with the posterior probabilities of the two audiovisual classes observed in  $e^{(2)}$  (in light blue for  $o_1$  and  $o_3$ , in purple for  $o_2$ ). The comparison of the all the  $p(o_j)$  and  $K^{(l)}$  justifies the potential triggering of head movements by the DW module, as observed at  $t = 74$  and  $t = 90$ .

due to the discontinuity of the sound signal: the audio experts did not detect any sound for these two frames (to give an idea:  $\arg \max(\mathbf{P}^a[t = 157]) = 0.176$ , whereas for the frame right before, at  $t = 156$ , five components out of thirteen are lying between  $p^a = 0.403$  and  $p^a = 1.00$ ). Going back to the resting position when an object stops emitting sound is part of the attempt of the overall HTM system to also inhibit the head movements in order to free the head for other potential purposes.

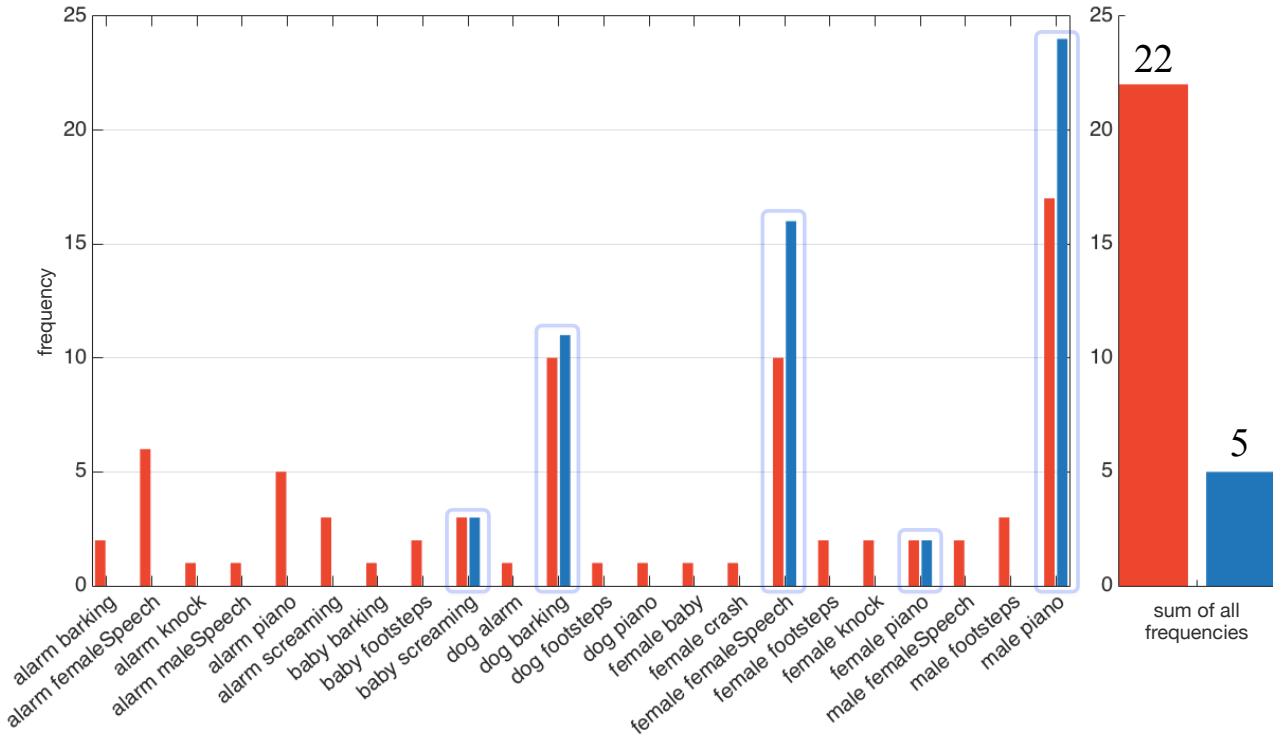


**Figure 11.** Results of the audiovisual classification (including the inference by the MFI module) obtained by (blue) the HTM system, (red) the naive robot. The two numbers on the right correspond to the value at the end of the exploration.

#### 814 4.4 Evaluation 5: Fusion & Classification

815 After having performed numerous evaluation in simulated conditions (see Table 2), this experiment  
 816 is focused on the evolution of the good audiovisual classification rate along the exploration of a real  
 817 environment. For that purpose, an environment is set up with five different sources, as presented in Table 3.  
 818 The three audio classes populating the environment have been selected because of their better experimental  
 819 recognition rate in the architecture. At each time step, the estimated audiovisual classes provided by  
 820 the overall HTM system is compared to the ground truth, and for each object. The resulting mean good  
 821 estimation rate  $\bar{\Gamma}_{MFI}[t]$ , computed over all objects, is plotted against time in Figure 11 (blue line). The  
 822 same is done for the naive robot, with a mean good estimation rate  $\bar{\Gamma}'_{R_n}[t]$  (red line in the same figure). As  
 823 expected, the proposed HTM system shows the best audiovisual classification rate. Indeed, one can see in  
 824 Figure 11 that the red line tends to the rate  $\bar{\Gamma}'_{R_n} = 37.9\%$  which is exactly the mean good classification  
 825 rate of the involved KS. In the same conditions, the MFI module converges to  $\bar{\Gamma}_{MFI} = 69.6\%$ . In the  
 826 very beginning of the experiment, both systems exhibit the same performances: the different smoothing  
 827 involved in the various computations (of the KS outputs, in the motor decisions, etc.) together with silences  
 828 in the sounds presented to the robot can explain this. But while the naive robot exhibits a constantly  
 829 decreasing good estimation rate of the audiovisual classes, the MFI module remains relatively robust to the  
 830 KS classification errors.

831 A direct consequence of these good performances of the HTM system can be observed in Figure 12 which  
 832 plots an histogram of all the audiovisual classes created by both systems (expressed in terms of number of  
 833 frames). The HTM system is able to considerably narrow the possible audiovisual classes existing in the  
 834 environment: from 22 by the naive robot, the HTM system narrows it down to only 5. However, one of the  
 835 class created is erroneous: PIANO, female has been mistaken with PIANO, male, but only for a short  
 836 period of time (two frames, i.e. 1s). This point has already been discussed in Section 3.3.4.



**Figure 12.** Number and labels of the audiovisual classes created by (blue) the MFI module, and (red) the naive robot. (Left panel:) Number of temporal frames (height of bars) during which the audiovisual classes have been categorized. (Light blue rectangles:) Audiovisual classes the two systems have in common. (Right panel:) Total number of different audiovisual classes created.

## 5 CONCLUSION

In this paper, a new system for the modulation of the exploratory behavior of a robot has been proposed. Based on the new notion of Congruence, it takes control of the head movements of a platform to put the robot attention towards audiovisual sources of interest. Additionally, it provides a robust description of the unknown environments explored all along the robot's life and following an unsupervised paradigm. This enriched representation consists, first, in the analysis of audiovisual objects through their relationship to the environments they are perceived in, and secondly, in how much the knowledge the system has about their actual audiovisual class is reliable and robust. Even in the case of classification errors by the audio or visual classifiers in the overall architecture, the system is then able to correctly infer the events' audiovisual classes by actively learning the interlink between the two modalities. All of this is achieved by the two constitutive modules of the HTM, namely the *Dynamic Weighting* module, and the *Multimodal Fusion & Inference* module. Each of them is able to trigger head movements that are used as an attentional reaction and as an active reaction to the need for additional data, respectively. Importantly, the extensive use of head movements is not limited to the sole benefit of the HTM system: audio localization algorithms such as (Nakashima and Mukai, 2005; Hornstein et al., 2006; Ma et al., 2017) relying also on head movements could be connected to the HTM as a top-down feedback unit, thus taking advantage from its motor commands to improve in parallel audio localization performances. The active self-supervised and online learning paradigm the MFI module relies upon, through the use of the *Multimodal Self- Organizing Map*, quickly provides the DW module with robust data while also offering inference abilities whenever a modality is missing (occlusion of the object, for instance). Whereas existing models provide audio-visual inference (Alameda-Pineda and Horaud, 2015) aiming at binding low-level cues of the audio and visual

857 data streams, the MFI module relies only on a higher level of representation of data, a representation  
858 that could be used as a top-down feedback to potentially enhance low-level audiovisual fusion algorithm.  
859 Additionally, the choice of learning the cross-modal relationship between auditory and visual data in  
860 an exclusively unsupervised way can be debated as not being powerful enough (Senocak et al., 2018).  
861 However, the results obtained here show significant improvements in the quality of the audiovisual data  
862 provided to the DW module without any inclusion of human knowledge. The system performances have  
863 been evaluated in realistic simulated conditions, but also on a real robot endowed with binaural audition  
864 and vision capabilities. Importantly, the overall architecture of the system, i.e. the TWO!EARS software,  
865 is made available online as an open source software<sup>2</sup>. The same applies for the proposed HTM system,  
866 entirely included inside this architecture<sup>3</sup>.

867 One of the main limitation of the current implementation is related to its high dependency to the  
868 localization experts. Indeed, the overall motor reactions are currently guided by each object azimuth  
869 localization, which have been shown precise enough to provide relevant results. Hopefully, binaural sound  
870 localization is a research topic by itself, and recent developments in the field show very robust algorithms,  
871 even in challenging acoustical conditions. Nevertheless, the robustness to localization errors could be  
872 enhanced by using tracking experts able to consolidate the sources position along time. For now, the HTM  
873 system is still being developed with the following improvements in mind. First, the definition of an object  
874 is currently limited to its audio and visual labels, while it could be enriched with additional information  
875 possibly coming from other modalities (emotions, audio pitch, forms and textures, etc.). Importantly,  
876 the proposed M-SOM has been designed to easily incorporate such additional parameters in the object  
877 definition: a subnetwork can be added for each of them together with their respective weights vectors.  
878 Concerning the Dynamic Weighting module, a significant improvement can be made by including the  
879 computation of a temporal habituation in order for the robot to not to be stuck in a deadlock kind of  
880 situation, as in Figure 8 where, if the scenario goes on forever, the robot would be keeping turning its head  
881 towards the CRYING, male. Finally, the coupling of the HTM system with other cognitive experts in the  
882 current framework is still under investigation. So far, the current version of the TWO!EARS software does  
883 not include others high- level cognitive experts. Nevertheless, the entire HTM system has been conceived  
884 with the idea that the motor exploration can also be guided by cognitive elements other than the ones  
885 implemented in the system. For instance, a model as the one recently proposed by Lanillos et al. (2015) on  
886 attention driven by social interaction, could easily be linked to the HTM, both benefiting from each other:  
887 one congruent source could still be focused because of its social interest, whereas a socially non-interesting  
888 object could still be focused for its high incongruence.

## ACKNOWLEDGMENTS

889 This work is supported by the European FP7 TWO!EARS project, ICT-618075, [www.twoears.eu](http://www.twoears.eu).

## REFERENCES

- 890 Alameda-Pineda, X. and Horaud, R. (2015). Vision-guided robot hearing. *The International Journal of*  
891 *Robotics Research* 34, 437–456
- 892 Baranes, A. and Oudeyer, P.-y. (2009). R-IAC: Robust Intrinsically Motivated Active Learning. *IEEE*  
893 *Transactions on Autonomous Mental Development* 1, 155–169

<sup>2</sup> <https://github.com/TWOEARS>

<sup>3</sup> <https://github.com/TWOEARS/audio-visual-integration>

- 894 Baranes, A. and Oudeyer, P.-Y. (2010). Intrinsically Motivated Goal Exploration for Active Motor Learning  
895 in Robots: A Case Study. In *Intelligent Robots and Systems (IROS), 2010 IEEE/RSJ International*  
896 *Conference on*. 1766–1773
- 897 Bauer, J. and Wermter, S. (2013). Self-organized neural learning of statistical inference from high-  
898 dimensional data. In *IJCAI*. 1226–1232
- 899 Berlyne, D. E. (1950). Novelty and Curiosity as Determinants of Exploratory Behavior. *British Journal of*  
900 *Psychology* 41, 68–80
- 901 Berlyne, D. E. (1965). *Structure and Direction in Thinking*
- 902 Braitenberg, V. (1986). *Vehicles: Experiments in Synthetic Psychology* (MIT Press)
- 903 Bustamante, G., Danès, P., Forgue, T., and Podlubne, A. (2016). Towards information-based feedback  
904 control for binaural active localization. In *2016 IEEE International Conference on Acoustics, Speech*  
905 *and Signal Processing (ICASSP)*. 6325–6329
- 906 Capdepuy, P., Polani, D., and Nehaniv, C. L. (2007). Maximization of Potential Information Flow as  
907 a Universal Utility for Collective Behaviour. In *IEEE Symposium on Artificial Life (Ieee)*, 207–213.  
908 doi:10.1109/ALIFE.2007.367798
- 909 Carrillo, H., Dames, P., Kumar, V., and Castellanos, J. A. (2015). Autonomous robotic exploration  
910 using occupancy grid maps and graph SLAM based on Shannon and Rényi entropy. In *Robotics and*  
911 *Automation (ICRA), 2015 IEEE International Conference on (IEEE)*, 487–494
- 912 Chapelle, O., Weston, J., and Schölkopf, B. (2002). Cluster Kernels for Semi-Supervised Learning.  
913 *Advances in Neural Information Processing Systems* 15 7, 1. doi:10.1016/S0090-3019(02)01037-6
- 914 Cohen-Lhyver, B. (2017). *Modulation de Mouvements de Tête pour l'Analyse Multimodale d'un*  
915 *Environnement Inconnu*. Ph.D. thesis, Université Pierre and Marie Curie
- 916 Cohen-Lhyver, B., Argentieri, S., and Gas, B. (2015). Modulating the Auditory Turn-to Reflex on the Basis  
917 of Multimodal Feedback Loops: the Dynamic Weighting Model. In *IEEE International Conference on*  
918 *Robotics and Biomimetics (ROBIO)*
- 919 Cohen-Lhyver, B., Argentieri, S., and Gas, B. (2016). Multimodal Fusion and Inference Using Binaural  
920 Audition and Vision. In *International Congress on Acoustics*
- 921 Corbetta, M., Patel, G., and Shulman, G. L. (2008). Review The Reorienting System of the Human Brain:  
922 From Environment to Theory of Mind. *Neuron* 58, 306–324. doi:10.1016/j.neuron.2008.04.017
- 923 Corbetta, M. and Shulman, G. L. (2002). Control of Goal-Directed and Stimulus-Driven Attention in the  
924 Brain. *Nature Reviews Neuroscience* 3, 201–215
- 925 Corneil, B. D., Munoz, D. P., Chapman, B. B., Admans, T., and Cushing, S. L. (2008). Neuromuscular  
926 consequences of reflexive covert orienting. *Nature neuroscience* 11, 13
- 927 Cuperlier, N., Quoy, M., and Gaussier, P. (2007). Neurobiologically inspired mobile robot navigation and  
928 planning. *Frontiers in Neurorobotics* 1. doi:10.3389/neuro.12
- 929 Downar, J., Crawley, A. P., Mikulis, D. J., and Davis, K. D. (2000). A Multimodal Cortical Network for  
930 the Detection of Changes in the Sensory Environment. *Nature Neuroscience* 3, 277–283
- 931 Driver, J. and Spence, C. (1998). Attention and crossmodal construction of space. *Trends in Cognitive*  
932 *Sciences* 2, 254–262
- 933 Droniou, A., Ivaldi, S., and Sigaud, O. (2015). Deep unsupervised network for multimodal perception,  
934 representation and classification. *Robotics and Autonomous Systems* 71, 83–98
- 935 Duangudom, V. and Anderson, D. V. (2007). Using auditory saliency to understand complex auditory  
936 scenes. In *European Signal Processing Conference, 2017 (EUSIPCO)*. 15th
- 937 Ferreira, J. F. and Dias, J. (2014). Attentional mechanisms for socially interactive robots—a survey. *IEEE*  
938 *Transactions on Autonomous Mental Development* 6, 110–125. doi:10.1109/TAMD.2014.2303072

- 939 Girard, B., Cuzin, V., Guillot, A., Gurney, K. N., and Prescott, T. J. (2002). Comparing a Brain-  
940 Inspired Robot Action Selection Mechanism With 'Winner-Takes-All'. In *From Animals to Animats*  
941 7: Proceedings of the seventh international conference on simulation of adaptive behavior (MIT Press),  
942 vol. 7, 75
- 943 Gurney, K., Prescott, T. J., and Redgrave, P. (2001a). A computational model of action selection in the  
944 basal ganglia. I. A new functional anatomy. *Biological cybernetics* 84, 401–10
- 945 Gurney, K., Prescott, T. J., and Redgrave, P. (2001b). A Computational Model of Action Selection in the  
946 Basal Ganglia. II. Analysis and Simulation of Behaviour. *Biological cybernetics* 84, 411–23
- 947 Henneberger, G., Brunsbach, B. J., and Klepsch, T. (1992). Field Oriented Control of Synchronous and  
948 Asynchronous Drives Without Mechanical Sensors Using a Kalman-Filter. In *European Conference on*  
949 *Power Electronics and Applications, 1992 (ECPEA)*. vol. 3, 664
- 950 Hinterstoesser, S., Cagniart, C., Ilic, S., Sturm, P., Navab, N., Fua, P., et al. (2012). Gradient response  
951 maps for real-time detection of textureless objects. *IEEE Transactions on Pattern Analysis and Machine*  
952 *Intelligence* 34, 876–888
- 953 Hopfinger, J. B., Buonocore, M. H., and Mangun, G. R. (2000). The Neural Mechanisms of Top-Down  
954 Attentional Control. *Nature neuroscience* 3, 284–291
- 955 Hornstein, J., Lopes, M., Santos-Victor, J., and Lacerda, F. (2006). Sound localization for humanoid  
956 robots-building audio-motor maps based on the hrtf. In *Intelligent Robots and Systems, 2006 IEEE/RSJ*  
957 *International Conference on (IEEE)*, 1170–1176
- 958 Huang, G.-b., Member, S., Zhu, Q.-y., and Siew, C.-k. (2006). Real-Time Learning Capability of Neural  
959 Networks. *Transactions on Neural Networks, 2006 IEEE* 17, 863–878
- 960 Huang, X. and Weng, J. (2002). Novelty and Reinforcement Learning in the Value System of Developmental  
961 Robots. In *Lund University Cognitive Studies*. 47–55
- 962 Huang, X. and Weng, J. (2004). Motivational system for human-robot interaction. In *International*  
963 *Workshop on Computer Vision in Human-Computer Interaction* (Springer), 17–27
- 964 Hunt, J. (1965). Intrinsic motivation and its role in psychological development. In *Nebraska symposium*  
965 *on motivation* (University of Nebraska Press), vol. 13, 189–282
- 966 Indovina, I. and Macaluso, E. (2007). Dissociation of Stimulus Relevance and Saliency Factors during  
967 Shifts of Visuospatial Attention. *Cerebral Cortex* 17, 1701–1711. doi:10.1093/cercor/bhl081
- 968 Itti, L., Koch, C., and Niebur, E. (1998). A Model of Saliency-Based Visual Attention for Rapid  
969 Scene Analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 20, 1254–1259.  
970 doi:10.1109/34.730558
- 971 James, W. (1890). *The Principles of Psychology* (Read Books Ltd)
- 972 Kayser, C., Petkov, C. I., Lippert, M., and Logothetis, N. K. (2005). Mechanisms for Allocating Auditory  
973 Attention: An Auditory Saliency Map. *Current Biology* 15, 1943–1947. doi:10.1016/j.cub.2005.09.040
- 974 Kohonen, T. (1982). Self-Organized Formation of Topologically Correct Feature Maps. *Biological*  
975 *Cybernetics* 43, 59–69. doi:10.1007/BF00337288
- 976 Kohonen, T. (1990). The Self-Organizing Map. *Proceedings of the IEEE* 78, 1464–1480
- 977 Kohonen, T. (2013). Essentials of the Self-Organizing Map. *Neural Networks* 37, 52–65. doi:10.1016/j.
- 978 neunet.2012.09.018
- 979 Lanillos, P., Ferreira, J. F., and Dias, J. (2015). Designing an artificial attention system for social robots  
980 (Institute of Electrical and Electronics Engineers)
- 981 Le Meur, O., Le Callet, P., Barba, D., and Thoreau, D. (2006). A Coherent Computational Approach  
982 to Model Bottom-Up Visual Attention. *Transactions on Pattern Analysis and Machine Intelligence* 28,  
983 802–817

- 984 Le Meur, O. and Liu, Z. (2015). Saccadic model of eye movements for free-viewing condition. *Vision*  
985 *research* 116, 152–164
- 986 Ma, N., Brown, G. J., and May, T. (2015). Exploiting deep neural networks and head movements for  
987 binaural localisation of multiple speakers in reverberant conditions. In *Interspeech*
- 988 Ma, N., May, T., and Brown, G. J. (2017). Exploiting deep neural networks and head movements for  
989 robust binaural localization of multiple sources in reverberant environments. *IEEE/ACM Transactions*  
990 *on Audio, Speech and Language Processing (TASLP)* 25, 2444–2453
- 991 Makarenko, A. A., Williams, S. B., Bourgault, F., and Durrant-Whyte, H. F. (2002). An Experiment in  
992 Integrated Exploration. In *Intelligent Robots and Systems, 2002. IEEE/RSJ International Conference*  
993 *on*, 534–539
- 994 May, T., van de Par, S., and Kohlrausch, A. (2011). A Probabilistic Model for Robust Localization Based  
995 on a Binaural Auditory Front-End. *Audio, Speech, and Language Processing, IEEE Transactions on* 19,  
996 1–13. doi:10.1109/TASL.2010.2042128
- 997 Meyer, J.-A., Guillot, A., Girard, B., Khamassi, M., Pirim, P., and Berthoz, A. (2005). The Psikharpx  
998 Project: Towards Building an Artificial Rat. *Robotics and Autonomous Systems* 50, 211–223. doi:10.  
999 1016/j.robot.2004.09.018
- 1000 Montemerlo, M., Thrun, S., Koller, D., and Wegbreit, B. (2002). FastSLAM: A Factored Solution to  
1001 the Simultaneous Localization and Mapping Problem. *Proc. of 8th National Conference on Artificial*  
1002 *Intelligence/14th Conference on Innovative Applications of Artificial Intelligence* 68, 593–598. doi:10.  
1003 1.1.16.2153
- 1004 Näätänen, R., Gaillard, A., and Mäntysalo, S. (1978). Early Selective-Attention Effect on Evoked Potential  
1005 Reinterpreted. *Acta Psychologica* 42, 313–329
- 1006 Nakashima, H. and Mukai, T. (2005). 3d sound source localization system based on learning of binaural  
1007 hearing. In *Systems, Man and Cybernetics, 2005 IEEE International Conference on (IEEE)*, vol. 4,  
1008 3534–3539
- 1009 Nothdurft, H.-C. (2006). Salience and Target Selection in Visual Search. *Visual Cognition* 14, 514–542
- 1010 O’Keefe, J. and Nadel, L. (1978). *The Hippocampus as a Cognitive Map*. 9 (Oxford: Clarendon Press)
- 1011 Oliva, A., Torralba, A., Castelhano, M. S., and Henderson, J. M. (2003). Top-Down Control of Visual  
1012 Attention in Object Detection. In *Image processing, 2003. icip 2003. proceedings. 2003 international*  
1013 *conference on*, vol. 1, I–253. doi:10.1109/ICIP.2003.1246946
- 1014 Oudeyer, P.-Y. and Kaplan, F. (2008). How can we define intrinsic motivation? In *Proceedings of the*  
1015 *8th International Conference on Epigenetic Robotics: Modeling Cognitive Development in Robotic*  
1016 *Systems, Lund University Cognitive Studies, Lund: LUCS, Brighton (Lund University Cognitive*  
1017 *Studies, Lund: LUCS, Brighton)*
- 1018 Oudeyer, P. Y., Kaplan, F., and Hafner, V. V. (2007). Intrinsic Motivation Systems for Autonomous Mental  
1019 Development. *IEEE Transactions on Evolutionary Computation* 11, 265–286. doi:10.1109/TEVC.2006.  
1020 890271
- 1021 Papliński, A. P. and Gustafsson, L. (2005). Multimodal feedforward self-organizing maps. In *International*  
1022 *Conference on Computational and Information Science* (Springer), 81–88
- 1023 Petersen, S. and Posner, M. (2012). The Attention System of the Human Brain: 20 Years After. *Annual*  
1024 *Review of Neuroscience* 21, 73–89. doi:10.1146/annurev-neuro-062111-150525.The
- 1025 Redgrave, P., Prescott, T. J., and Gurney, K. (1999). The Basal Ganglia: A Vertebrate Solution to the  
1026 Selection Problem? *Neuroscience* 89, 1009–1023
- 1027 Roy, N., McCallum, A., and Com, M. W. (2001). Toward Optimal Active Learning through Monte Carlo  
1028 Estimation of Error Reduction. In *International Conference on Machine Learning*

- 1029 Ruesch, J., Lopes, M., Bernardino, A., Hörnstein, J., Santos-Victor, J., and Pfeifer, R. (2008). Multimodal  
1030 Saliency-Based Bottom-Up Attention a Framework for the Humanoid Robot iCub. *Robotics and*  
1031 *Automation, 2008. ICRA 2008. IEEE International Conference on*, 962–967doi:10.1109/ROBOT.2008.  
1032 4543329
- 1033 Rushworth, M. F., Noonan, M. A. P., Boorman, E. D., Walton, M. E., and Behrens, T. E. (2011). Frontal  
1034 Cortex and Reward-Guided Learning and Decision-Making. *Neuron* 70, 1054–1069. doi:10.1016/j.  
1035 neuron.2011.05.014
- 1036 Ryan, R. M. and Deci, E. L. (2000). Intrinsic and Extrinsic Motivations: Classic Definitions and New  
1037 Directions. *Contemporary Educational Psychology* 25, 54–67. doi:10.1006/ceps.1999.1020
- 1038 Schillaci, G., Hafner, V. V., and Lara, B. (2014). Online learning of visuo-motor coordination in a  
1039 humanoid robot. a biologically inspired model. In *Development and Learning and Epigenetic Robotics*  
1040 (*ICDL-Epirob*), 2014 Joint IEEE International Conferences on (IEEE), 130–136
- 1041 Schmidhuber, J. (1991). Curious Model-Building Control Systems. In *Neural Networks, 1991. 1991 IEEE*  
1042 *International Joint Conference on* (Singapore), 1458–1463
- 1043 Schymura, C., Walther, T., Kolossa, D., Brown, G. J., Ma, N., and Brown, G. J. (2014). Binaural  
1044 Sound Source Localisation using a Bayesian-network-based Blackboard System and Hypothesis-driven  
1045 Feedback. *Forum Acusticum* doi:10.13140/2.1.4026.4966
- 1046 Senocak, A., Oh, T.-H., Kim, J., Yang, M.-H., and Kweon, I. S. (2018). Learning to localize sound source  
1047 in visual scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.  
1048 4358–4366
- 1049 Smith, R., Self, M., and Cheeseman, P. (1988). A Stochastic Map for Uncertain Spatial Relationships. *4th*  
1050 *International Symposium on Robotic Research*, 467–474
- 1051 Thompson, G. C. and Masterton, R. B. (1978). Brain stem auditory pathways involved in reflexive head  
1052 orientation to sound. *Journal of neurophysiology* 41, 1183–1202
- 1053 Treisman, A. M. and Gelade, G. (1980). A Feature-Integration Theory of Attention. *Cognitive psychology*  
1054 12, 97–136. doi:10.1016/0010-0285(80)90005-5
- 1055 Tsiami, A., Katsamanis, A., Maragos, P., and Vatakis, A. (2016). Towards a behaviorally-validated  
1056 computational audiovisual saliency model. *Acoustics, Speech and Signal Processing (ICASSP), 2016*  
1057 *IEEE International Conference on*, 2847–2851doi:10.1109/ICASSP.2016.7472197
- 1058 Two!Ears (2016a). *Report on Evaluation of the Two!Ears Expert System*. Tech. Rep. June  
1059 [Dataset] Two!Ears (2016b). The Two!Ears project
- 1060 Walter, W. G. (1951). A Machine that Learns. *Scientific American* 185, 60–63
- 1061 Walther, D., Rutishauser, U., Koch, C., and Perona, P. (2005). Selective visual attention enables learning  
1062 and recognition of multiple objects in cluttered scenes. *Computer Vision and Image Understanding* 100,  
1063 41–63. doi:10.1016/j.cviu.2004.09.004
- 1064 Walther, T. and Cohen-Lhyver, B. (2014). Multimodal Feedback in Auditory-Based Active Scene  
1065 Exploration. In *Forum Acusticum*
- 1066 Zhou, D., Bousquet, O., Lal, T. N., Weston, J., and Schölkopf, B. (2004). Learning with Local and Global  
1067 Consistency. In *Advances in Neural Information Processing Systems* 17, 2004. *Proceedings of Neural*  
1068 *Information Processing Systems*. vol. 1. doi:10.1.1.115.3219