

2 **Learning emergent behaviours for a hierarchical Bayesian**
3 **framework for active robotic perception**

4 João Filipe Ferreira · Christiana Tsiourti ·
5 Jorge Dias

6
7 © Marta Olivetti Belardinelli and Springer-Verlag 2012

8 **Abstract** In this research work, we contribute with a
9 behaviour learning process for a hierarchical Bayesian
10 framework for multimodal active perception, devised to be
11 emergent, scalable and adaptive. This framework is com-
12 posed by models built upon a common spatial configura-
13 tion for encoding perception and action that is naturally
14 fitting for the integration of readings from multiple sensors,
15 using a Bayesian approach devised in previous work. The
16 proposed learning process is shown to reproduce goal-
17 dependent human-like active perception behaviours by
18 learning model parameters (referred to as “attentional
19 sets”) for different free-viewing and active search tasks.
20 Learning was performed by presenting several 3D audio-
21 visual virtual scenarios using a head-mounted display,
22 while logging the spatial distribution of fixations of the
23 subject (in 2D, on left and right images, and in 3D space),
24 data which are consequently used as the training set for the
25 framework. As a consequence, the hierarchical Bayesian
26 framework adequately implements high-level behaviour
27 resulting from low-level interaction of simpler building
28 blocks by using the attentional sets learned for each task,
29 and is able to change these attentional sets “on the fly,”
30 allowing the implementation of goal-dependent behaviours
31 (i.e., top-down influences).
32

Keywords Multisensory active perception · Hierarchical 33
Bayes models · Bioinspired robotics · Human–robot 34
interaction · Emergence · Scalability · Adaptive behaviour 35

Introduction 36

How should the uncertainty and incompleteness of the 37
environment be represented and modelled so as to increase 38
the autonomy of a robot? Can a robotic system perceive, 39
infer, decide and act more efficiently by using a probabilistic 40
framework? These are two of the challenging questions 41
robotics researchers are currently facing in the design of 42
more autonomous and intelligent artificial robotic systems. 43

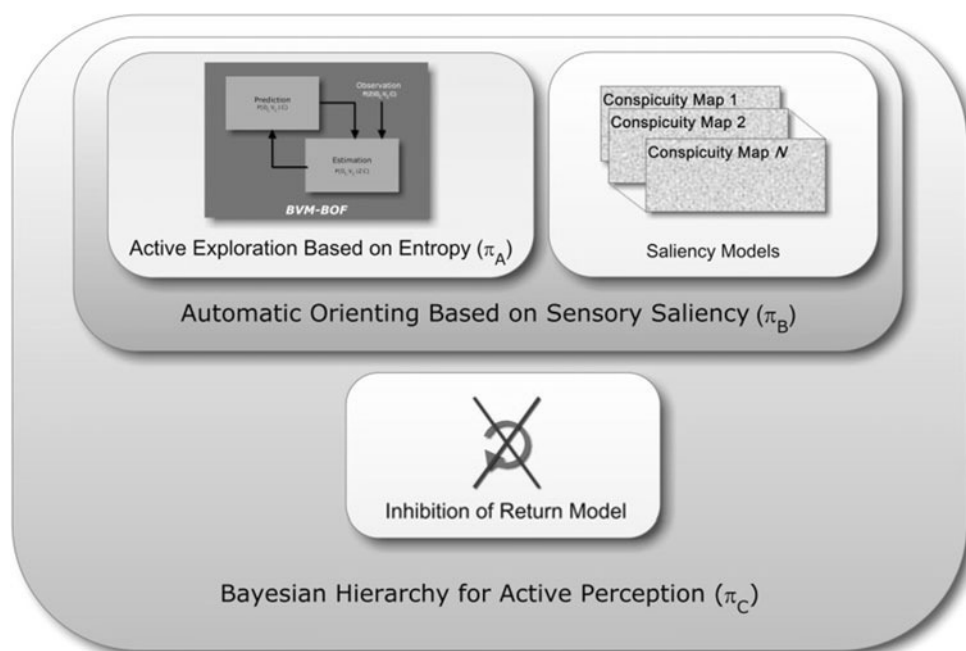
Previous work conducted in the Institute of Systems and 44
Robotics by Ferreira et al. (2012, 2011)—see also (Ferreira 45
2011)—has contributed to addressing these challenges by 46
providing the basis of a framework for artificial active 47
multimodal perception, comprised of a real-time GPU-based 48
implementation of a scalable, adaptive and emergent hier- 49
archical Bayesian active perception system that simulates 50
several bottom-up-driven human behaviours of attention 51
guidance—see Fig. 1. It was devised mainly to be used in 52
human–robot interaction (HRI) applications. 53

These emergent behaviours are implemented by combin- 54
ing simple behaviours using a set of weights, thus 55
implementing a process analogous to the attentional set as 56
defined by Corbetta and Shulman (2002). The bottom layer 57
of this framework consists of a log-spherical inference grid 58
updated using a Bayesian filter, the Bayesian volumetric 59
map or BVM. Ferreira et al. (2012, 2011) modelled visu- 60
oauditory perception using an approach that finds its 61
inspiration in the fast pathways believed to exist in the 62
human brain, which are closely linked to primal instincts of 63
survival and basic social interaction. 64

A1 J. F. Ferreira (✉) · C. Tsiourti · J. Dias
A2 ISR, University of Coimbra, Coimbra, Portugal
A3 e-mail: jfilipe@isr.uc.pt

A4 J. Dias
A5 Khalifa University of Science, Technology and Research
A6 (KUSTAR), Abu Dhabi, UAE

Fig. 1 Conceptual diagram for active perception model hierarchy (Ferreira et al. 2012)



65 In this text, we will give a brief overview of a behaviour
66 learning process for this framework, designed to estimate
67 its free parameters (identified as attentional sets) for dif-
68 ferent free-viewing and active search tasks.

69 Related work

70 During the last decades, researchers from different fields
71 (psychologists, neuroscientists and, more recently, com-
72 puter scientists) have investigated visual attention thor-
73 oughly and also, to a lesser extent, visuoauditory atten-
74 tion, both in terms of model analysis and in terms of synthesis
75 and implementation.

76 In the fields of neuroscience and psychology, for
77 example, research on these issues has ranged from the
78 early twentieth century, such as the work by Buswell
79 (1935), to recently Castelhanos et al. (2009), and sepa-
80 rately and subsequently also Mills et al. (2011), which
81 investigated the influence of task instruction on specific
82 parameters of eye movement control, such as the number
83 of fixations and gaze duration on specific objects. On the
84 other hand, research work in computational models of
85 artificial active perception has ranged from the seminal
86 work of Bajcsy (1985) and Aloimonos et al. (1987),
87 through Itti et al. (1998) in the bottom-up influence of
88 visual saliency and Breazeal et al. (2001) in active vision
89 for social robots, to important and recent work that has
90 attempted to implement learning by imitation for active

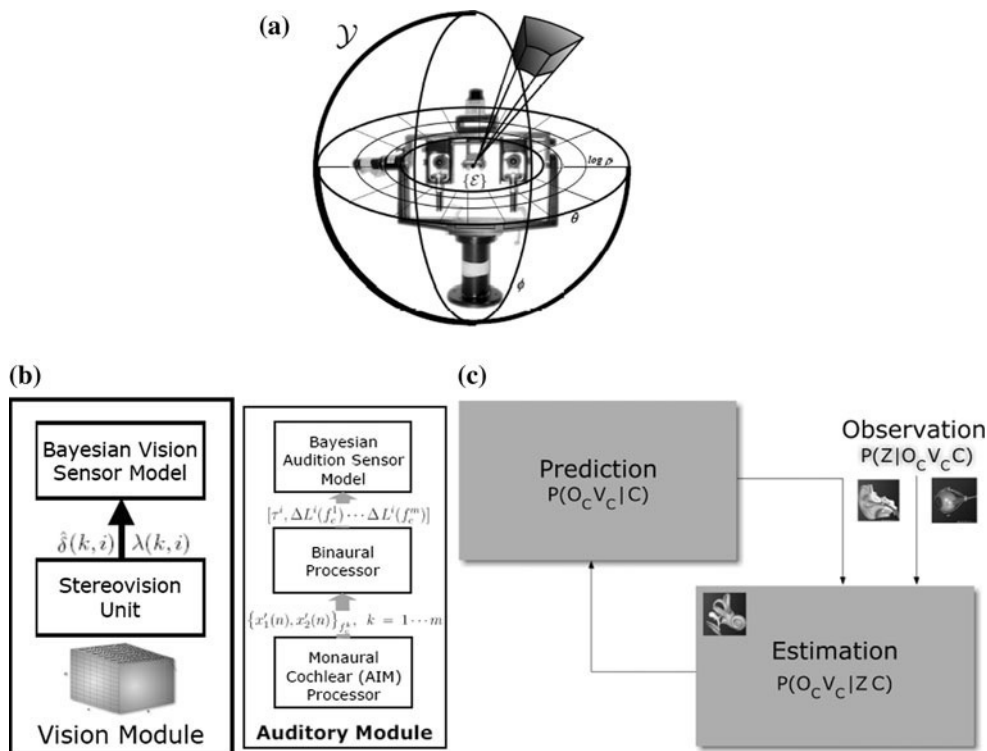
perception behaviours, such as the work by Belardinelli et al. (2007).

We improve the current state of the art in all of these fields by contributing in two specific fronts: (1) a learning process integrated within psychophysical experiments, representing an important tool for both model analysis and synthesis in human studies and robot development, respectively; and (2) a fully integrated probabilistic framework that closely follows human behaviour, formally and explicitly dealing simultaneously with perceptual uncertainty, multisensory fusion and the perception–action loop.

Overview of the hierarchical Bayesian framework for active robotic perception

In the BVM framework, cells of a partitioning grid on the BVM log-spherical space Y associated with the egocentric coordinate system $\{E\}$ are indexed through $C \in Y$, representing the subset of positions in Y corresponding to the “far corners” ($\log_b \rho_{\max}, \theta_{\max}, \varphi_{\max}$) of each cell C ; O_C is a binary variable representing the state of occupancy of cell C (as in the commonly used occupancy grids—see Elfes (1989)), and V_C is a finite vector of random variables that represent the state of all local motion possibilities used by the prediction step of the Bayesian filter associated with the BVM for cell C , assuming a constant velocity hypothesis, as depicted on Fig. 2. Sensor measurements (i.e., the result of visual and auditory processing) are denoted by Z —

Fig. 2 Multisensory perception framework details (Ferreira et al. 2012). **a** The Bayesian volumetric map (BVM) referred to the egocentric coordinate frame of the robotic active perception system; **b** BVM sensor models; **c** BVM Bayesian occupancy filter



117 observations $P(Z | O_C V_C C)$ are given by the Bayesian
 118 sensor models of Fig. 2, which yield results already inte-
 119 grated within the log-spherical configuration.

120 The BVM framework is extensible in such a way that
 121 other properties characterised by additional random vari-
 122 ables and corresponding probabilities might be represented,
 123 other than the already implemented occupancy and local
 124 motion properties of the BVM, by augmenting the hier-
 125 archy of operators through Bayesian subprogramming (Bess-
 126 ière et al. 2008). This ensures that the framework is scalable.
 127 On the other hand, the combination of these strategies to
 128 produce a coherent behaviour ensures that the framework is
 129 emergent.

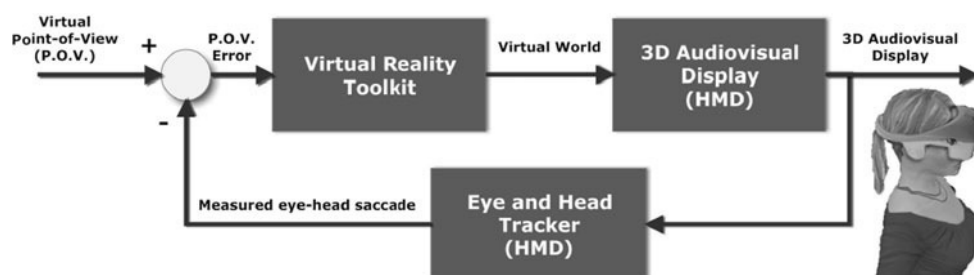
130 Three decision models were proposed by Ferreira et al.
 131 (2012): π_A , which implements entropy-based active
 132 exploration based on the BVM; π_B , which uses entropy and
 133 saliency together for active perception; and finally, π_C ,
 134 which adds a simple inhibition of return (IoR) mechanism
 135 based on the fixation point of the previous time step. In
 136 other words, each model incorporates its predecessor
 137 through Bayesian fusion, therefore constituting a model
 138 hierarchy—see Fig. 1. Each decision model will infer a
 139 probability distribution on the next point of fixation for the
 140 next desired gaze shift represented by a random variable G_t
 141 $\in Y$ at each time $t \in [1, t_{max}]$. For more details, refer to
 142 (Ferreira et al. 2012).

The complete set of variables that set up the framework
 and its extensions is summarised in the following list (ref-
 erences to temporal properties removed for easier reading):

- C : cell index on the BVM occupancy grid given by the 3D coordinates of its “far corner;”
- Z : generic designation for either visual or auditory sensor measurements;
- O_C : binary value signalling the fact that a cell C is either empty or occupied by an object;
- V_C : discrete variable indicating instantaneous local motion vector for objects occupying cell C ;
- G : fixation point for next gaze shift in log-spherical coordinates;
- U_C : entropy gradient-based variable ranging from 0 to 1, signalling the potential interest (i.e., 0 and 1, meaning minimally and maximally interesting, respectively) of cell C as future focus of attention given the uncertainty on its current state given by (O_C, V_C) , thus promoting an active exploration behaviour;
- S_C^i : binary value describing the i th of N sensory saliency of cell C ;
- $Q_C^i = P([S_C^i = 1] | Z^i C)$: probability of a perceptually salient object occupying cell C ;
- R_C : inhibition level for cell C as a possible future focus of attention modelling the inhibition of return behaviour, ranging from no inhibition (0) to full inhibition (1).

Author Proof

Fig. 3 Virtual point-of-view generator set-up that allows the updating of audiovisual stimuli presentation according to the monitored subjects' gaze direction



169 **Learning automatic emergent behaviours for active**
170 **multisensory perception**

171 Any of the three decision models, π_A , π_B and π_C , results in
172 an inference result similar to the following equation
173 (which, in fact, corresponds to model π_B),
174

175 The automatic multisensory active perception behav-
176 iours emerge from the distributions $P(Q_C^{i,t} | G^t \pi_B)$, which
177 are either beta distributions $B(\alpha_Q, \beta_Q)$ for

perception, namely active exploration and automatic ori- 208
enting using sensory saliency, as valid strategies in human 209
behaviour regarding saccade generation. 210

Discussion 211

At the time of writing this text, pilot experiments have 212
already been conducted, validating the learning procedure 213
and already displaying promising results. 214

$$P(G^t | V^{1 \rightarrow t} O^{1 \rightarrow t} S^t \pi_B) \propto P(G^t | V^{1 \rightarrow t} O^{1 \rightarrow t} \pi_A) \prod_C \left[\prod_{i=1}^N P(Q_C^{i,t} | G^t \pi_B) \right]. \quad (1)$$

183 [$G^t = C$] expressing that, for a given point of fixation
184 proposal for the next gaze shift, $Q_C^{i,t}$ is more likely near 1,
185 or a uniform distribution on $Q_C^{i,t}$ for [$G^t \neq C$]. In this
186 equation, the probability of a perceptually salient object
187 occupying cell C , given by $Q_C^{i,t}$, is to be replaced by U_C^t or
188 R_C^t , depending on which model besides π_B one is referring
189 to, either π_A or π_C , respectively.

190 Therefore, the learning process in this context is defined
191 as supervised learning through the maximum likelihood
192 estimation (MLE) of the free parameters of the respective
193 beta distributions. The training data to perform this learn-
194 ing is gathered from psychophysical experiments, in which
195 human subjects using a head-mounted device are presented
196 with realistic 3D, audiovisual, virtual-reality scenarios. The
197 subjects' tracked head-eye gaze shifts control the virtual
198 stereoscopic-binaural point of view and hence the pro-
199 gression of each stimulus movie—see Fig. 3—while
200 audiovisual stimuli and corresponding fixation points are
201 logged. This way, controlled conditions will be enforced by
202 proposing both free-viewing and active search tasks to the
203 subjects, thus enabling as systematic estimation of distri-
204 bution parameters to promote the appropriate human-like
205 emergent behaviour depending on the robot's goal. On the
206 other hand, this learning process will allow testing both of
207 our primary hypotheses for active visuoauditory

In the following months, the final experiments will be 215
conducted and a robotic demonstrator will be set up, using 216
the learned attentional sets in implementing different tasks. 217
This will further prove that the Bayesian hierarchical 218
framework adequately follows human-like active percep- 219
tion behaviours, namely by exhibiting the following 220
desirable properties: 221

Emergence—High-level behaviour results from low- 222
level interaction of simpler building blocks. 223

Scalability—Seamless integration of additional inputs is 224
allowed by the Bayesian programming formalism used 225
to state the models of the framework. 226

Adaptivity—Initial “genetic imprint” of distribution 227
parameters may be changed “on the fly” through 228
parameter manipulation, thus allowing for the imple- 229
mentation of goal-dependent behaviours (i.e., top-down 230
influences). 231

Acknowledgments The authors would particularly like to thank, at 232
the Institute of Biomedical Research in Light and Image of the 233
University of Coimbra (IBILI/UC), Prof. Miguel Castelo-Branco, 234
João Castelhamo, Carlos Amaral and Marco Simões for their help with 235
the psychophysical experiments. 236

Conflict of interest This supplement was not sponsored by outside 237
commercial interests. It was funded entirely by ECONA, Via dei 238
Marsi, 78, 00185 Roma, Italy. 240

241 **References**

- 242 Aloimonos J, Weiss I, Bandyopadhyay A (1987) Active vision. *Int J*
 243 *Comput Vis* 1:333–356 264
- 244 Bajcsy R (1985) Active perception vs passive perception. In: Third
 245 IEEE workshop on computer vision, Bellair, Michiganm,
 246 pp 55–59 265
- 247 Belardinelli A, Pirri F, Carbone A (2007) Bottom-up gaze shifts and
 248 fixations learning by imitation. *IEEE Trans Syst Man Cybern B*
 249 37(2):256–271 266
- 250 Bessière P, Laugier C, Siegwart R (eds) (2008) Probabilistic
 251 reasoning and decision making in sensory-motor systems,
 252 volume 46 of Springer tracts in advanced robotics, Springer.
 253 ISBN: 978-3-540-79006-8 267
- 254 Breazeal C, Edsinger A, Fitzpatrick P, Scassellati B (2001) Active
 255 vision for sociable robots. *IEEE Trans Syst Man Cybern A Syst*
 256 *Hum* 31(5):443–453 268
- 257 Buswell GT (1935) How people look at pictures: a study of the
 258 psychology and perception in art. University Chicago Press,
 259 Chicago 269
- 260 Castelhana MS, Mack ML, Henderson JM (2009) Viewing task
 261 influences eye movement control during active scene perception.
 262 *J Vis* 9:1–15 270
- Corbetta M, Shulman GL (2002) Control of goal-directed and
 stimulus-driven attention in the brain. *Nat Rev Neurosci*
 3(3):201–215 271
- Elfes A (1989) Using occupancy grids for mobile robot perception
 and navigation. *IEEE Comput* 22(6):46–57 272
- Ferreira JF (2011) Bayesian cognitive models for 3D structure and
 motion multimodal perception. PhD thesis, Faculty of Sciences
 and Technology of the University of Coimbra (FCTUC) 273
- Ferreira JF, Lobo J, Dias J (2011) Bayesian real-time perception
 algorithms on GPU—Real-time implementation of Bayesian
 models for multimodal perception using CUDA. *J Real-Time*
Image Proc 6(3):171–186 274
- Ferreira JF, Castelo-Branco M, Dias J (2012) A hierarchical Bayesian
 framework for multimodal active perception. *Adapt Behav*
 20(3):172–190 Published online ahead of print, March 1 275
- Itti L, Koch C, Niebur E (1998) A model of saliency-based visual
 attention for rapid scene analysis. *IEEE Trans Pattern Anal Mach*
Intell 20(11):1254–1259 276
- Mills M, Hollingworth A, Dodd MD (2011) Examining the influence
 of task set on eye movements and fixations. *J Vis* 11:1–15 277
- 280 281 282 283

UNCORRECTED PROOF