

CONAIM: A Conscious Attention-Based Integrated Model for Human-Like Robots

Alexandre da Silva Simões, *Member, IEEE*, Esther Luna Colombini,
and Carlos Henrique Costa Ribeiro, *Member, IEEE*

Abstract—Understanding consciousness is one of the most fascinating challenges of our time. From ancient civilizations to modern philosophers, questions have been asked on how one is conscious of his/her own existence and about the world that surrounds him/her. Although there is no precise definition for consciousness, there is an agreement that it is strongly related to human cognitive processes such as attention, a process capable of promoting a selection of a few stimuli from a huge amount of information that reaches us constantly. In order to bring the consciousness discussion to a computational scenario, this paper presents conscious attention-based integrated model (CONAIM), a formal model for machine consciousness based on an attentional schema for human-like agent cognition that integrates: short- and long-term memories, reasoning, planning, emotion, decision-making, learning, motivation, and volition. Experimental results in a mobile robotics domain show that the agent can attentively use motivation, volition, and memories to set its goals and learn new concepts and procedures based on exogenous and endogenous stimuli. By performing computation over an attentional space, the model also allowed the agent to learn over a much reduced state space. Further implementation under this model could potentially allow the agent to express sentience, self-awareness, self-consciousness, autonoetic consciousness, mineness, and perspectivalness.

Index Terms—Attention, cognition, machine consciousness, memory, robotics.

I. INTRODUCTION

THROUGHOUT human history, it is possible to find reflections about the relation between body and mental phenomena. However, even after centuries of discussion, there is no common understanding about what consciousness is. For some authors, consciousness is everything we experience [1]. For others, it is the human ability to assign feelings to specialized mental capacities such as thinking, emotions, and speech [2]. Others define consciousness as the process that allows relevant information to remain online long enough so that it may be synchronously processed by multiple cortical networks [3]. Yet, some consider consciousness as the subjective experiences/awareness that one has when awake, following earlier suggestions that a more precise definition, given our currently inadequate scientific understanding of consciousness, is best

left to the future [4]. Another understanding is that consciousness results from the complex interaction of the multimodal association areas of, at least, the parietotemporal and limbic cortices in an animal's brain, and refers to the capacity to experience oneself as a being subject to the past, present, and future, including the reflection on oneself as a being that is aware of its surrounding environment [5]. In general lines, there is a common understanding that consciousness is related to subjective experiences of humans, and that it is intrinsically connected to other human cognitive subsystems such as thinking, emotions, intention, decision-making, and reasoning.

Machine consciousness is the field of artificial intelligence concerned with the production of conscious processes in engineering hardware or software devices [6]. According to some definitions, a machine is considered conscious if besides the required mechanisms for perception, action, learning, and associative memory, it has a central executive that controls all the processes—conscious or subconscious—of the machine, and this central executive itself is driven by the machine's motivation and goal selection, attention-switching, semantic and episodic memory, using cognitive perception and cognitive understanding of motivations, thoughts, or plans to control learning, attention, motivations, and monitor actions [7]. A useful picture that places artificial consciousness among other concepts of the artificial intelligence domain—such as intelligence, life, and awareness—is shown in Fig. 1 [7], [8].

Several computational models have been proposed for machine consciousness (for an overview see [4], [9]). In this scenario, two main aspects stand out. First, one of the key concepts in consciousness models is attention. In this context, the recent proposition of an attentional model [10]–[12] capable of dealing with several types of attention—sustained, selective, oriented, divided, overt, and covert—and that can deal with endogenous and exogenous stimuli suggests further investigation on the applicability of this model for machine consciousness purposes. As a second relevant aspect, in the vast majority of consciousness models, authors present only a brief concept of the modules that compose the system, without relevant details about the connections and operation of these modules. This lack of formalization makes their implementation and extension impossible. In general lines, these works are also not followed by computational experiments that can demonstrate their capacities, preventing comparisons among them.

In order to bring the consciousness discussion to a computational scenario, this paper presents CONAIM, a novel conscious attention-based integrated model that comprises: short- and

Manuscript received February 01, 2015; revised September 08, 2015; accepted October 20, 2015. Date of publication January 14, 2016; date of current version September 27, 2017.

A. S. Simões is with the Department of Control and Automation Engineering, Campus of Sorocaba, Univ Estadual Paulista (UNESP), Sorocaba, Brazil (e-mail: assimoes@sorocaba.unesp.br).

E. L. Colombini is with the University of Campinas (UNICAMP), Campinas, Brazil (e-mail: esther.colombini@gmail.com).

C. H. C. Ribeiro is with the Technological Institute of Aeronautics (ITA), São José dos Campos, Brazil (e-mail: carlos@ita.br).

Digital Object Identifier 10.1109/JSYST.2015.2498542

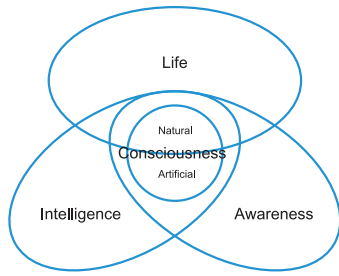


Fig. 1. Consciousness and its relation to intelligence, life, and awareness. Adapted from [7].

long-term memories, reasoning, planning, emotion, attention, decision-making, learning, motivation, and volition, and that is also able to accommodate other human-like cognitive processes. The formal model of the system is presented and, to the best of our knowledge, it represents the first known complete formalization of a model of cognitive processes that have been claimed to be linked to consciousness.

This paper is organized as follows. Section II introduces some of the basic concepts related to consciousness. Section III presents an overview of the computational models and theories for machine consciousness. Section IV discloses the proposed model for consciousness. Section V shows the experiments that were conducted. Section VI presents the experimental results and discussion. Finally, Section VII concludes this work.

II. FUNDAMENTALS OF CONSCIOUSNESS

Some of the terms adopted in the consciousness literature have close and, sometimes, conflicting meanings. In this way, this section presents usual definitions of related concepts. A first important concept is being awake, i.e., not asleep. Humans can be awake without being conscious [13]. Another concept is sentience, the consciousness in the generic sense of simply being a sentient creature, one capable of sensing and responding to its world [14], [15], being capable of suffering or feeling pleasure or happiness [5]. Awareness usually refers to the process that occurs as a result of the interaction of the nervous system of the animal (including the sensory apparatus) and its environment, resulting in the ability of the animal to react to received stimuli. Awareness is, hence, related to terms such as sentience, perception, and cognition and does not require consciousness. A worm can be aware without being conscious [5]. Self-awareness is related to the being aware of yourself, i.e., the ability to introspect and to recognize itself as an individual separated from the environment and from other individuals.

Self-consciousness is a term of more complex definition. Some authors understand it as a synonym for self-awareness or consciousness in its minimalist sense [7]. However, it is possible to identify at least four distinct meanings for self-consciousness as follows [5]:

- 1) consciousness of its own body sensations and about its own movements;
- 2) ability to recognize itself in a mirror (e.g., ability demonstrated by chimpanzees, dolphins, and human beings);
- 3) consciousness of something happening in your environment and how this can affect you under a proper context;

- 4) consciousness of some decision and its consequences.

The autonoetic consciousness can be defined as the ability of a human being to, mentally, think himself/herself in the past and future, using his/her capacity of reasoning, anticipating situations, and reasoning about his/her own thoughts. Another important concept is volition, the cognitive process whereby an individual decides to perform a certain action, i.e., the mental activity through what intention is implemented [16].

III. THEORIES AND MODELS FOR MACHINE CONSCIOUSNESS

One of the most accepted theories of consciousness is the integrated information theory (IIT) [1], [17], which claims that consciousness is integrated information, or, more precisely, that 1) the quantity of consciousness corresponds to the amount of integrated information (Φ) generated by a complex of elements and 2) the quality of experience is specified by the set of informational relationships generated within that complex. This theoretical framework can be applied to basic neurobiological observations as well as to artificial systems [18].

Computational models proposed for machine consciousness can be divided into four groups [4]: 1) models based on global workspace theory (GWT); 2) models based on internal self-models; 3) models based on higher level representations; and 4) models based on attention.

The GWT models are inspired by biological distributed and inter-communicating processes in regions of the cerebral cortex during mental efforts [19]–[21]. In these computational models, multiple specialized parallel processes cooperate and compete for the access to the global workspace. The coalition of winning processes can send information to all other processes through a broadcast system. The main hypothesis investigated by authors is that conscious experience can emerge through the collective interactions between the specialized processors via the global workspace. One of the first models for consciousness based on this paradigm was IDA [22], [23], a cognitive architecture that integrates several types of memory (perceptual, episodic, procedural, and working memory), consciousness, action selection, decision-making, and volition. One important contribution of IDA was the proposition of a cognitive cycle: perception, pre-conscious buffer, local association, competition for consciousness, broadcast, recruitment, setup of goal context hierarchy, action selection, and action execution. The adoption of a learning mechanism in IDA has led to the LIDA system [21]. Finally, the CERA-CRANIUM model [6], [24] presented a conscious system with attention and emotional learning based on the competition between specialized modules applied to a robotic agent.

Models based on internal self-models are biologically inspired by neural activation patterns that can be found in the human brain, and typically have an internal representation of the spatial properties of their bodies (a body image). The main hypothesis that guide authors in these works are [4]: 1) agents who reason about themselves can reason about their own reasoning about themselves [25]; 2) phenomenal consciousness could emerge from introspective reasoning mechanisms about perception [26]; and 3) some key properties of a self-concept

(existence, continuity over time, supervenience on a physical substrate, etc.) can form the basis for self-modeling in intelligent agents, regardless of whether or not they are embodied [27]. Some examples of the implementation of these models in robots can be found in the literature [28]–[30].

Models based on higher level representations consider that conscious mental states are distinguished from unconscious mental states by their level of representation. CLARION [31], [32] is an architecture with declarative and procedural knowledge, control mechanisms, and goals with two levels: 1) a higher conscious level—with a symbolic knowledge representation—and 2) a lower unconscious level—where concepts are represented by neural networks activity patterns. Another architecture with five levels was proposed in [33]. The lowest level is a reactive level, whereas the highest one—where conscious emerges in tasks that require cognitive efforts—is a symbolic rules level concerned with the robot movements. The CERA-CRANIUM—discussed in GWT models—can also be classified as a higher level model since its architecture has distinct cognitive levels. Similar works are [34] and [35].

At each moment, one is conscious of only a fraction of the information that reaches his/her sensory system [4]. Attention is the mechanism responsible for actively selecting, which stimulus will receive the attentional focus. Consciousness and attention are, therefore, extremely related processes [36]–[40]. Although there is such proximity, they are not, at least apparently, the same process [40], [41], since they have distinct functions and neural mechanisms [42]–[46] that occur jointly. Recent studies suggest that stimuli not attended by the attentional mechanism are not able to reach consciousness [3]. In this way, consciousness could be positioned at a distinct hierarchical level, and would be related to attention, competition, and learning [47]. Some authors argue that attention and consciousness are distinct mechanisms that feed the same decision process that leads to behavior [46].

Given the undeniable proximity between consciousness and attention, several computational models have been proposed based on attentional mechanisms. Even in these models, however, consciousness is interpreted in different ways. A computational model concerned with the visual perception composed of layers of spiking neurons was presented in [48]. These neurons were capable of building supermaps based on changing in their activities due to changing in input image features. According to authors, patterns in these maps could be used as a representation of consciousness. A distinct approach was proposed in [49]–[51] where several distinct interconnected modules without a precise hierarchical structure are capable of recruiting other modules in a competitive process to attend to the same topic they are working on. The machine, hence, could reach consciousness of a topic when several modules cooperate focusing their attention on the perceptions related to this topic. The model proposed by [7] and [52] presents an architecture composed of three main blocks: 1) a sensory-motor block; 2) episodic memory and learning block; and 3) executive central block. The attention shift results from the competition between several types of signals (internal thinking, new perceptions, shifts in motivation, etc.). A model that integrates consciousness, attention, and creativity is presented in [41].

Creativity could be used by subjects to amplify their subjective experiences, and is reached by decreasing the intensity of the attentional process, allowing the generation of unconscious thoughts until the subject can reach a mental state suitable for solving the problem.

IV. PROPOSED MODEL

This section presents: 1) the architecture of the proposed model and its main aspects; 2) model formalization; and 3) model operation.

A. Model Architecture

The proposed model consists of an architecture and a formal description of its structures and corresponding operations. The architecture is presented in Fig. 2 and is composed of two main systems.

- 1) An *attentional system*, following the *selection for perception* components of the model proposed in [11], comprises sensory memory, feature maps, weights associated with feature maps, combined feature maps, saliency map, and attentional map. The course of the attentional dynamics also follows the one proposed in [11].
- 2) A *cognitive system*, which replaces the *selection for action* component proposed by [11], composed of:
 - a) *Decision-making* component that decides on the actions to be performed by the agent;
 - b) *Short-term memory* component (working memory) that stores a small amount of information—typically received from the saliency map and sensory memory—which can be used by various modules of the cognitive agent system, in particular for long-term memory storage and recall;
 - c) *Long-term memory* component stores, indefinitely, large amounts of information, consisting of:
 - i) Episodic memory component that stores data, in a declarative form, relative to specific events, specifically *what*, *when*, and *where*, using a suitable data structure.
 - ii) Semantic memory component that stores facts and knowledge about the world in a declarative way, using an appropriate data structure.
 - iii) Procedural memory component that stores knowledge regarding procedures that can be executed by the agent, typically in the form of state-action pairs.
 - d) *Evaluation* component: groups variables, lists, and functions related to the agent's performance, which composed of:
 - i) Goals that stores the agent's current goal(s);
 - ii) Evaluator: set of metrics designed to assess the agent's performance from different perspectives, typically acting over the internal state;
 - iii) Evaluation state: the current values of the agent's state for all the metrics considered;
 - iv) Task: current task under agent's execution;
 - v) Volition that represents the process of transforming the agent's intention into a goal.

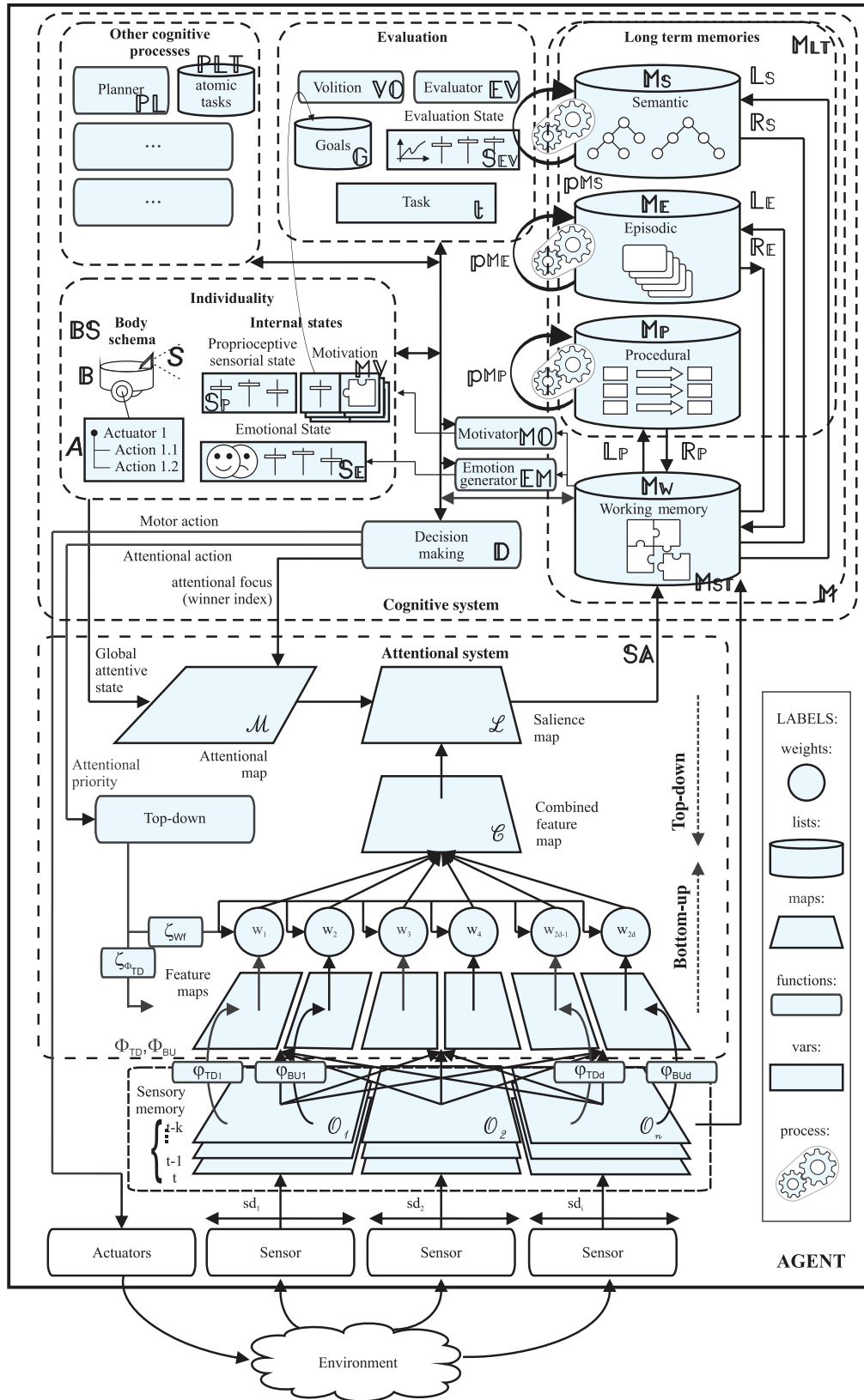


Fig. 2. CONAIM architecture, composed of an attentional system and a cognitive system: \mathcal{M} , memory; MLT , long-term memory; MS , semantic memory; ME , episodic memory; MP , procedural memory; Mw , working memory; D , decision-making; EM , emotion generator; MO , motivator; BS , body schema; B , body; A , actuators, S , sensors; MV , motivation; VO , volition; EV , evaluator; G , goals; t , task; PL , planner; PLT : planner list; p : processes; S_E : emotional state; S_P : proprioceptive state; S_{EV} : evaluation state; S_A : attentional state.

e) *Cognitive processes* that combine any number of the agent's cognitive functions, such as planning and state estimation;

f) *Individuality* component that composed of a set of inner states, such as proprioceptive, motivation, emotional state, and a body schema of the agent.

3) A set of *processes* with diverse purposes, running in background, which can act over the cognitive system components, not necessarily synchronized to its remaining components.

The system operates considering that any function or cognitive system process can request information from any other module and/or change the agent's internal states at any time. Therefore, the flow of information through the decision-making module is not mandatory.

B. Formal Model

This section presents a formal description of the proposed model. It uses the variables and functions defined in the *attentional system* formalized in [11]. The *cognitive system* can be described as follows.

- 1) A set of actions \mathbb{A} with an arbitrary number of elements a_i termed *actions*, with $\{a_1, a_2, \dots, a_n\} \in \mathbb{A}$, and $\mathbb{A} = \mathbb{A}_A \cup \mathbb{A}_M$.
- 2) A set of attentional actions \mathbb{A}_A with an arbitrary number of elements a_{Ai} termed *attentional actions*, with $\mathbb{A}_A = \{a_{A1}, a_{A2}, \dots, a_{An}\}$.
- 3) A set of motor actions \mathbb{A}_M with an arbitrary number of elements a_{Mi} termed *motor actions*, with $\mathbb{A}_M = \{a_{M1}, a_{M2}, \dots, a_{Mn}\}$.
- 4) A list of *goals* \mathbb{G} with an arbitrary number of elements g_i , with $\mathbb{G} = \{g_1, g_2, \dots, g_n\}$.
- 5) A set of tasks \mathbb{T} with an arbitrary number of elements t_i termed *tasks*, with $\mathbb{T} = \{t_1, t_2, \dots, t_n\}$, and the specific task t named agent's *current task*.
- 6) A set of states \mathbb{S} with an arbitrary number of elements s_i termed *states*, where $\mathbb{S} = \{s_1, s_2, \dots, s_n\}$, with $\mathbb{S} = \mathbb{S}_A \cup \mathbb{S}_E \cup \mathbb{S}_{EV} \cup \mathbb{S}_P$.
- 7) A set of attentional states \mathbb{S}_A with an arbitrary number of elements such that $\mathbb{S}_A = \{s_{A1}, s_{A2}, \dots, s_{An}\}$, where each element s_{Ai} , denoted *attentional state*, consists of a *saliency map* of m dimension, termed $\mathcal{L} \in \mathbb{R}^m$, with elements represented by $\{l_1, l_2, \dots, l_m\}$, according to the model proposed by [11].
- 8) A set of emotional states \mathbb{S}_E with an arbitrary number of elements, where each term s_{Ei} , denoted *emotional state*, is composed of an arbitrary number *emotional variables* v_E , i.e., $\mathbb{S}_E = \{v_{E1}, v_{E2}, \dots, v_{En}\}$, where $v_{Ei} \in \mathbb{R}$.
- 9) A set of evaluation states \mathbb{S}_{EV} with an arbitrary number of elements, where each element s_{EVi} , denoted *evaluation state*, is constituted of an arbitrary number of *evaluation variables* v_{EV} , in other words, $\mathbb{S}_{EV} = \{v_{EV1}, v_{EV2}, \dots, v_{EVn}\}$, where $v_{EVi} \in \mathbb{R}$.
- 10) A set of proprioceptive states \mathbb{S}_P composed of an arbitrary number of elements, with each element s_{Pi} ,

termed *proprioceptive state*, constituted of an arbitrary number of *proprioceptive variables* v_P , i.e., $\mathbb{S}_P = \{v_{P1}, v_{P2}, \dots, v_{Pn}\}$, with $v_{Pi} \in \mathbb{R}$.

- 11) A set of motivations \mathbb{M} of the agent with an arbitrary number of elements m_i denoted *motivation*, where each element $m_i \in \mathbb{M}$ is composed by a pair that contains $v_{M} \in \mathbb{R}$ denoted *motivation value*, and a memory element $m \in \mathbb{M}$, in other words, $m_i = \{v_{M}, m_i\}$.
- 12) A superset named *memory* \mathbb{M} , with $\mathbb{M} = \mathbb{M}_{LT} \cup \mathbb{M}_{ST}$.
- 13) A set named *short-term memory* \mathbb{M}_{ST} , with $\mathbb{M}_{ST} = \mathbb{M}_W$.
- 14) A set named *long-term memory* \mathbb{M}_{LT} , with $\mathbb{M}_{LT} = \mathbb{M}_E \cup \mathbb{M}_P \cup \mathbb{M}_S$.
- 15) A doubly linked list of elements named *working memory* \mathbb{M}_W that can able to admit an arbitrary number of elements m_W , i.e., $\mathbb{M}_W = \{m_{W1}, m_{W2}, \dots, m_{Wn}\}$. The dynamics of insertions and removals in the list follows the first-in-first-out (FIFO) principle. It is recommended modeling elements m_W as proposed by [27].
- 16) A doubly linked list of elements named *semantic memory* \mathbb{M}_S that can comprise an arbitrary number of elements m_S , with $\mathbb{M}_S = \{m_{S1}, m_{S2}, \dots, m_{Sn}\}$, where each element of the semantic memory m_{Si} contains an arbitrary number of n -tuples d_{MS} , such that $m_{Si} = \{d_{MS1}, d_{MS2}, \dots, d_{MSn}\}$. The ordering of the elements m_S in the list \mathbb{M}_S is given by the value of d_{MSi} denoted *semantic memory index element*.
- 17) A doubly linked list of elements named *episodic memory* \mathbb{M}_E that can able to admit an arbitrary number of elements m_E , with $\mathbb{M}_E = \{m_{E1}, m_{E2}, \dots, m_{En}\}$, where each element of the episodic memory m_{Ei} contains an arbitrary number of n -tuples d_{ME} , such that $\{d_{ME1}, d_{ME2}, \dots, d_{MEN}\} \in m_{Ei}$. The ordering of the elements m_E in the list \mathbb{M}_E is given by the increasing value of d_{MEi} denoted *episodic memory index element*. It is recommended modeling the elements m_E as proposed by [53].
- 18) A doubly linked list of elements termed *procedural memory* \mathbb{M}_P that can comprise an arbitrary number of elements m_P , with $\mathbb{M}_P = \{m_{P1}, m_{P2}, \dots, m_{Pn}\}$, where each element of the procedural memory m_{Pi} is a triple composed of a *task* $t \in \mathbb{T}$, a set of qualifiers d_{MP} and a *state-action set* \mathbb{C}_{SA} , with an arbitrary number of elements $c_{SAi} \in \mathbb{C}_{SA}$, with each element c_{SA} formed by state $s \in \mathbb{S}$ and action $a \in \mathbb{A}$ pairs, with $c_{SAi} = \{s_i, a_i\}$, such that $m_{Pi} = \{t_i, d_{MP}, \mathbb{C}_{SA}\}$.
- 19) A *planning function* \mathbb{PL} that maps a goal $g \in \mathbb{G}$ to a set with an arbitrary number of tasks $t \in \mathbb{PLT}$ with \mathbb{PLT} represented by a list of tasks internal to the planner. In other words,

$$\mathbb{PL}$$

$$g \longrightarrow \{t_1, t_2, \dots, t_n\} \in \mathbb{PLT}.$$

- 20) A *volition function* \mathbb{VO} that maps motivations $m \in \mathbb{M}$ into goals $g \in \mathbb{G}$ is given by

$$\mathbb{VO}$$

$$\{\mathbb{M}\} \longrightarrow \{\mathbb{G}\}.$$

- 21) An *evaluation function* \mathbb{E}_V that maps the agent's current state \mathcal{S} its goals \mathcal{G} and task $\mathcal{t} \in \text{PLT}$ to the agent's evaluation state \mathcal{S}_{EV} . In other words,

$$\{\mathcal{S}_A, \mathcal{S}_E, \mathcal{S}_{\text{EV}}, \mathcal{S}_P, \mathcal{G}, \mathcal{t}, \text{PLT}\} \xrightarrow{\mathbb{E}_V} \{\mathcal{S}_{\text{EV}}\}.$$

- 22) A *motivation function* \mathbb{M}_O that maps the content of the working memory \mathcal{M}_W to the agent's motivation \mathcal{M}_V , i.e., $\mathbb{M}_O: \mathcal{M}_W \rightarrow \mathcal{M}_V$.
- 23) An *emotion generator* \mathbb{E}_M that maps the content of the working memory \mathcal{M}_W into an emotional state \mathcal{S}_E of the agent: $\mathbb{E}_M: \mathcal{M}_W \rightarrow \mathcal{S}_E$.
- 24) A *decision-making function* \mathbb{D} that typically coordinates all cognitive system lists and elements choosing an action $a \in A$, which can be a motor action $\mathcal{O}_M \in \mathcal{A}_M$ or an attentional action $\mathcal{O}_A \in \mathcal{A}_A$, and a number i denoted *attentional process winner index* or *attentional focus*. In other terms,

$$\{\mathcal{S}_A, \mathcal{S}_E, \mathcal{S}_{\text{EV}}, \mathcal{S}_P, \mathcal{G}, \mathcal{T}, \mathcal{M}_W, \mathcal{M}_V, \text{PLT}, \mathcal{BS}\} \xrightarrow{\mathbb{D}} \{\mathcal{O}_M, \mathcal{O}_A, i, \mathcal{G}, \mathcal{T}, \mathcal{M}_W, \mathcal{M}_V, \mathcal{S}_E\}.$$

- 25) A *body schema* \mathcal{BS} of the agent that is a n -tuple that comprises the agent's set of sensors \mathcal{S} actuators \mathcal{A} , and a n -tuple \mathcal{B} that contains a set of descriptors of the physical characteristics of the agent structure (weight, volume, center of gravity, external form, etc.), such that $\mathcal{BS} = \{\mathcal{S}, \mathcal{A}, \mathcal{B}\}$.
- 26) A set of unconscious asynchronous *processes* \mathcal{P} , where each $\{\mathcal{P}_1, \mathcal{P}_2, \dots, \mathcal{P}_n\} \in \mathcal{P}$ is represented by a process responsible for performing a mapping from a specific system component to updates in the agent's state and/or lists

$$\{\mathcal{S}_A, \mathcal{S}_E, \mathcal{S}_{\text{EV}}, \mathcal{S}_P, \mathcal{M}_V, \mathcal{M}, \mathcal{M}_W, \mathcal{G}, \mathcal{T}\} \xrightarrow{\mathcal{P}} \{\mathcal{S}_A, \mathcal{S}_E, \mathcal{S}_{\text{EV}}, \mathcal{S}_P, \mathcal{M}_V, \mathcal{M}, \mathcal{M}_W, \mathcal{G}, \mathcal{T}\}.$$

- 27) A *semantic learning* \mathbb{L}_S that can perform a mapping usually is given by

$$\{\mathcal{m}_W, \mathcal{M}_W, \mathcal{M}_S\} \xrightarrow{\mathbb{L}_S} \{\mathcal{m}_S\}.$$

- 28) An *episodic learning* \mathbb{L}_E that can perform a mapping usually is given by

$$\{\mathcal{m}_W, \mathcal{M}_W, \mathcal{M}_E\} \xrightarrow{\mathbb{L}_E} \{\mathcal{m}_E\}.$$

- 29) A *procedural learning* \mathbb{L}_P that can perform a mapping usually is given by

$$\{\mathcal{m}_W, \mathcal{M}_W, \mathcal{M}_P\} \xrightarrow{\mathbb{L}_P} \{\mathcal{m}_P\}.$$

- 30) A *semantic recall* \mathbb{R}_S that can perform a mapping usually is given by

$$\{\mathcal{m}_W, \mathcal{M}_W, \mathcal{M}_S\} \xrightarrow{\mathbb{R}_S} \{\mathcal{m}_W\}.$$

- 31) An *episode recall* \mathbb{R}_E that can perform a mapping usually is given by

$$\{\mathcal{m}_W, \mathcal{M}_W, \mathcal{M}_E\} \xrightarrow{\mathbb{R}_E} \{\mathcal{m}_W\}.$$

- 32) A *procedure recall* \mathbb{R}_P that can perform a mapping usually is given by

$$\{\mathcal{m}_W, \mathcal{M}_W, \mathcal{M}_P\} \xrightarrow{\mathbb{R}_P} \{\mathcal{m}_W\}.$$

C. Model Operation

The operation of the model proposed in Section IV-B is given by:

- 1) *Start a new attentional cycle:*

- Compute the selection for perception component of the model proposed by [11] until a *saliency map* \mathcal{L} is available;
- Compute the *winner* of the competitive process, in other words, the index of the element with maximum saliency, according to

$$\text{winner} = \arg(\max(\mathcal{L})).$$

- 2) *Start a new cognitive cycle:*

- Create a new element in the working memory $\mathcal{m}_W \in \mathcal{M}_W$, inserting in \mathcal{m}_W the index of winner, the observations $\mathcal{O}_{\text{winner}}$ and the saliency map \mathcal{L} , such that

$$\{\text{winner}, \mathcal{O}_{\text{winner}}, \mathcal{L}\} \longrightarrow \mathcal{m}_W \in \mathcal{M}_W.$$

- Compute the *semantic learning* \mathbb{L}_S over \mathcal{m}_W aiming at the insertion or extension of the concepts in \mathcal{M}_S , i.e., $\mathbb{L}_S: \{\mathcal{m}_W, \mathcal{M}_W, \mathcal{M}_S\} \rightarrow \mathcal{m}_S$.
- Compute the *episodic learning* \mathbb{L}_E over \mathcal{m}_W aiming at the insertion or extension of the episodes in \mathcal{M}_E , i.e., $\mathbb{L}_E: \{\mathcal{m}_W, \mathcal{M}_W, \mathcal{M}_E\} \rightarrow \mathcal{m}_E$.
- Compute the *procedural learning* \mathbb{L}_P over \mathcal{m}_W aiming at the insertion or extension of the procedures in \mathcal{M}_P , i.e., $\mathbb{L}_P: \{\mathcal{m}_W, \mathcal{M}_W, \mathcal{M}_P\} \rightarrow \mathcal{m}_P$.
- Compute a *semantic recall* \mathbb{R}_S over \mathcal{m}_W , inserting in \mathcal{m}_W the content recalled \mathcal{m}_S or *null*, i.e., $\mathbb{R}_S: \{\mathcal{M}_W, \mathcal{M}_S, \emptyset\} \rightarrow \mathcal{m}_W$. \mathbb{R}_S can be performed based on many distinct aspects $\mathcal{d}_{\mathcal{M}_S}$ of elements \mathcal{m}_S (shape, color, smell, purpose, number of recalls, etc.). For each $\mathcal{d}_{\mathcal{M}_S}$, a distinct matching algorithm must perform a full search in the memory elements.
- Compute an *episode recall* \mathbb{R}_E over \mathcal{m}_W , inserting in \mathcal{m}_W the content recalled \mathcal{m}_E or *null*, i.e., $\mathbb{R}_E: \{\mathcal{M}_W, \mathcal{M}_E, \emptyset\} \rightarrow \mathcal{m}_W$. \mathbb{R}_E can be performed based on many distinct aspects $\mathcal{d}_{\mathcal{M}_E}$ of elements \mathcal{m}_E (time of event, environment conditions, objects, actions performed, etc.). For each $\mathcal{d}_{\mathcal{M}_E}$, a distinct matching algorithm must perform a full search in the memory elements.
- Compute a *procedure recall* \mathbb{R}_P over \mathcal{m}_W , inserting in \mathcal{m}_W the content recalled \mathcal{m}_P or *null*, i.e., $\mathbb{R}_P: \{\mathcal{M}_W, \mathcal{M}_P, \emptyset\} \rightarrow \mathcal{m}_W$. \mathbb{R}_P can be performed based on

distinct qualifiers d_{MP} of a task (purposes, restrictions, conditions, etc.). Some reasoning can be necessary to check the applicability of a procedure in a given situation to achieve an expected result.

- h) Compute the *motivation function* MO . Motivation functions can be proposed based on a large number of psychological theories [16]. Typically, the motivation is reduced when certain activity is carried out over a long time, or if it is perceived as boring, difficult, or if the agent has more enjoyable activities. It also depends on the agent as an individual, on its emotional state and memories or on its curiosity regarding new stimuli. If the item in m_W indicates that:

- i) M_S , M_E , or M_P have a concept, episode, or procedure under significant evolution that should be investigated by the agent, or
- ii) Elements recalled have relevant emotions associated, or
- iii) Elements are frequently recalled from memories

then the agent is *motivated* to relate to the concept, episode or procedure by doing:

- if $m_W \subset MV$, then update the value of the associated motivation v_{MV} .
- else, create a new item $m_{v_i} \in MV$, insert a high value of motivation in $v_{MV_i} \in m_{v_i}$ and insert m_W in $m_i \in m_{v_i}$.

and also to *update* the values associated to motivations $m_W \subset MV$ according to a time decay policy, eventually removing items from the list.

- i) Compute the *Emotion Generator* EM based on m_W (data observed and recalled memories) to update the agent's *emotional state* S_E composed by a number of emotions v_E . Although a large range of basic emotions can be investigated, some candidates are anger, disgust, fear, happiness, and sadness, surprise [54], [55]. Emotion functions must affect the intensity of the emotions trying to mimic the complex biological regulation that takes place in the cortex, hippocampus, thalamus, and related regions of the human brain. For an overview, see [7], [56], [57].
- j) Compute the agent's *volition function* VO , evaluating MV and G , to decide whether a new goal $g \in G$ should be defined for the agent. VO is assumed to be a process distinct from MO , since VO is more related to a deliberative process. It can be divided into preactional and actional phases. The preactional phase is concerned with forming interpretations and planning when, where, and how to initiate the action. VO must choose considering the real possibility to accomplish the goal, the real benefits on achieving it, environmental conditions, strength of motivation, etc. The actional phase is related to the action's execution evaluation and the reasons to sustain it. Some factors must be considered: ability to resist distractions, environment capacity to capture

agent's attention, unpredicted obstacles, discouragement from performing action, emotional state [16].

- k) If needed, execute the *planning function* PL to determine the atomic tasks $\{\ell_1, \ell_2, \dots, \ell_n\} \in PLT$ that are required to fulfil the agent's goal G

$$G \xrightarrow{PL} \{\ell_1, \ell_2, \dots, \ell_n\} \in PLT.$$

- l) Evaluate the goals G , the agent's current task ℓ , its state S and the tasks that should be executed PLT , and compute the agent's *evaluation state* S_D through

$$\{S_A, S_{EV}, S_E, S_P, G, \ell, PLT\} \xrightarrow{EV} S_{EV}.$$

- m) Execute a *decision-making* D function by doing:

- i) Considering the list of goals G , the current agent's task ℓ , the agent's general state S , and the tasks that should be executed PLT , decide on the maintenance or *modification of the current task* ℓ . In other words, compute

$$\{S_A, S_E, S_{EV}, S_P, G, \ell, PLT\} \xrightarrow{D} \ell.$$

- ii) Insert ℓ as a working memory item m_W and wait for a recall of the R_P task from M_P , by doing:

- If the procedure to solve the task is known, i.e., if $R_P: \{m_P, MP\} \rightarrow m_W$ results in a $m_W \neq \emptyset$ then execute the task ℓ using the procedure described by the state-action pairs $C_{SA} \in m_W$.
- Otherwise, if the recall returns a *fail*, i.e., $m_W = \emptyset$, then evaluate if the agent should start a learning procedure for the task ℓ . If yes, a new element $m_{p_i} \in M_P$ with $m_{p_i} = \{\ell_i, C_{SA_i}\}$ should be created and the learning should be started.

- 3) Start a new attentional cycle (step 1) for time = time + 1.

Besides the model operation described, at least three unconscious asynchronous processes can be running in background.

- a) A process p_{M_S} that permanently scans the semantic memory M_S seeking to derive or to unify concepts, with a mapping typically represented by: $p_{M_S}: M_S \rightarrow M_S$.
- b) A process p_{M_E} that permanently scans the episodic memory M_E seeking to derive or to unify episodes, with a mapping typically represented by: $p_{M_E}: M_E \rightarrow M_E$.
- c) A process p_{M_P} that permanently scans the procedural memory M_P seeking to derive or to unify processes, with a mapping typically represented by: $p_{M_P}: M_P \rightarrow M_P$.

V. MATERIALS AND METHODS

A. Simulation Environment

The system modules were implemented in C++ with wxWidgets and the experiments were conducted on the USARSim

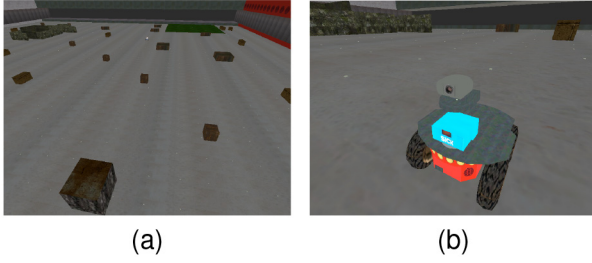


Fig. 3. Screenshots of the USARSim platform. (a) Overview of the DM-TestRoom250 scenario, with boxes of distinct sizes, walls, and victims, with wide open space. (b) Pioneer P2-AT virtual robot.

platform [58]. A P2AT simulated mobile robot with eight sonars and range scanner, both covering the frontal 180° of the robot was deployed in the environment. Sonars and laser have, respectively, 15 cm–5 m and 20 m range. All sensors include a random noise generator, configured to add up to 5% noise with a standard deviation $\sigma = 0.1$. The simulation scenario, shown in Fig. 3, was DM-Testroom-250, with large clear spaces, walls, irregular ground, and large fixed boxes. The robot was initially placed close to a box.

B. Experiment Description

Two distinct experiments were performed. EXP01 is an experiment where a robot with the attentional system was placed in the environment in front of a box. No cognitive system was adopted, and a simple controller drives the robot to the direction of the most salient stimulus. In EXP02, the robot was placed in the environment and controlled by the proposed conscious model, the robot was free to decide its own trajectory and to learn. The aim of this experiment is to show that consciousness-related aspects can emerge in this architecture due to the fact that a central executive \mathbb{D} can coordinate learning, attention, motivation, and actions driven by its own motivations, goals, perception, attention-shifting, and memory-stored data. The attentional and cognitive systems setup as well as the system operation are detailed in the next sections.

C. Attentional System Setup

The observation spaces adopted were defined as follows. O_1 represents the observation space for sonar readings, defined by: $o_{1n_t} = \text{sonar}_{n_t}$ with $n \in [1, 8]$; O_2 represents the observation space for range scanner readings, defined by: $o_{2n_t} = \text{range}_{n_t}$ with $n \in [1, 180]$.

Three bottom-up feature maps were adopted.

- 1) \mathcal{F}_1 represents the level of discrepancy of the intensity of motion detected in the environment;
- 2) \mathcal{F}_2 represents the level of discrepancy of the direction of moving objects in the environment relative to the agent;
- 3) \mathcal{F}_3 represents the level of discrepancy of the distance between obstacles (static or moving) and the agent.

Details on the functions that perform the mapping from the observation spaces \mathcal{O} to the feature maps spaces \mathcal{F} can be found in [11]. Two top-down feature maps were adopted.

- 1) \mathcal{F}_4 represents the interest of the agent for elements disposed at a specific distance from it.
- 2) \mathcal{F}_5 represents the interest of the agent for specific regions of the space.

The saliency map \mathcal{L} is an 8-dimensional state space and is used as the basis for the learning process instead of the observation space, composed by an 8-dimensional \mathcal{O}_1 and a 180-dimensional \mathcal{O}_2 . This represents a 95% reduction in the usual state space.

D. Cognitive System Setup

In this experiment, a conscious robot was placed in the environment with the following configuration.

- 1) Actions set $\mathbb{A} = \mathbb{A}_A \cup \mathbb{A}_M$ with the following actions.
 - a) \mathcal{O}_{M1} rotates toward the attentional focus.
 - b) \mathcal{O}_{M2} moves forward.
 - c) \mathcal{O}_{A1} drives the attentional focus through the top-down feature \mathcal{F}_4 to elements very close to the agent.
 - d) \mathcal{O}_{A2} drives the attentional focus through the top-down feature \mathcal{F}_4 to elements close to the agent.
 - e) \mathcal{O}_{A3} drives the attentional focus through the top-down feature \mathcal{F}_4 to elements far from the agent.
 - f) \mathcal{O}_{A4} drives the attentional focus through the top-down feature \mathcal{F}_5 to the region on the agent extreme right.
 - g) \mathcal{O}_{A5} drives the attentional focus through the top-down feature \mathcal{F}_5 to the region on the agent right side.
 - h) \mathcal{O}_{A6} drives the attentional focus through the top-down feature \mathcal{F}_5 to the region on the agent front.
 - i) \mathcal{O}_{A7} drives the attentional focus through the top-down feature \mathcal{F}_5 to the region on the agent left side.
 - j) \mathcal{O}_{A8} drives the attentional focus through the top-down feature \mathcal{F}_5 to the region on the agent extreme left side.
- 2) Evaluation: List of goals \mathbb{G} and Task \mathbb{T} , all initially empty.
- 3) Motivation \mathbb{M}_V , initially empty.
- 4) Motivation function \mathbb{M}_O : every new concept in \mathbb{M}_W that could not be located in the agent's semantic memory adds a motivation element in \mathbb{M}_V with a certain level of intensity that changes across time.
- 5) Semantic memory \mathbb{M}_S : initially empty. Shape was adopted as the key element $\mathcal{O}_{MS} \in \mathbb{M}_S$. The object shape was initially represented using a point cloud. These points can typically evolve to a polygon using a *convex hull* algorithm performed by a semantic memory background process \mathbb{P}_{MS} , as soon as points close enough are available.
- 6) Semantic recall \mathbb{R}_S : A full search in the semantic memory is performed looking for elements that have points close to those currently observed by the agent according to an Euclidian distance criteria.
- 7) Semantic learning \mathbb{L}_S : the following rules were applied:
 - i) if recall $\mathbb{R}_S = \emptyset$ for \mathbb{M}_W , then create a new entry \mathbb{M}_S in \mathbb{M}_S and store data \mathbb{M}_W in $\mathcal{O}_{MS} \in \mathbb{M}_S$; ii) if object already exists in \mathbb{M}_S , then append \mathbb{M}_W data in \mathbb{M}_S .
- 8) Procedural memory \mathbb{M}_P : initially empty.
- 9) Procedural recall \mathbb{R}_P : a full search in the procedural memory is performed looking for a task \mathbb{T} that can be applied over object \mathbb{M}_S .

- 10) Procedural learning \mathbb{L}_P : implemented through reinforcement learning with the following structure:
 - a) algorithm: Q-learning;
 - b) agent's state \mathcal{L} and actions \mathbb{A} , as defined in Section V-D;
 - c) state space discretization: five levels for each state (5^8 states);
 - d) exploration rate: 0.95, decaying linearly with time to 0;
 - e) learning rate: 0.9;
 - f) temporal discount factor: 0.99;
 - g) maximum number of trials: 100. Max 500 steps per trial;
 - h) reward structure: $r = r + 10$ for each new data inserted into the concept under investigation in the semantic memory; $r = r + 3/d$, where d is the distance from the agent to the first element that generated the concept; $r = r - 2$ if the robot is close to collision; $r = r - 10$ if the robot already collided; $r = r - 50$ if the robot is permanently stuck.
- 11) Decision-making \mathbb{D} : The following rules were adopted:
 - 1) if task \mathbb{T} and goals list \mathbb{G} are both empty, then drive the robot to the most salient stimulus;
 - 2) if task \mathbb{T} is empty and goals list \mathbb{G} is not empty, then remove the highest priority goal g from \mathbb{G} , and make it the new task \mathbb{T} ;
 - 3) if task \mathbb{T} requires a known procedure, then verify state s and execute corresponding action a ;
 - 4) if task \mathbb{T} requires an unknown procedure, then apply RL learning algorithm and execute referred action a .

E. System Operation

The attentional and conscious courses are executed as depicted in Fig. 4. The resulting saliency map \mathcal{L} , the winner of the competitive process and its observations $\mathcal{O}_{\text{winner}}$ are stored in the working memory \mathbb{M}_W . The arrival of new data in \mathbb{M}_W stimulates recalls and learning at procedural and semantic memories. Since memories were initially empty, no recalls are expected. In this case, semantic memory is expected to create new concepts based on the data received from the environment through the learning function \mathbb{L}_S . One of the available mechanisms that can stimulate the agent to complete a concept is motivation. In this way, as soon as the first new data reaches \mathbb{M}_S , a motivation is expected to be generated in \mathbb{M}_V to foster the investigation of this concept. The volition function $\mathbb{V}\mathbb{O}$ typically transforms this motivation into a goal, adding it to the agent's goals list \mathbb{G} , which was initially empty. The decision-making function \mathbb{D} is responsible for turning this goal into a task, by inserting it at the working memory \mathbb{M}_W to recall a procedure for this task. Since the procedural memory \mathbb{M}_P has no stored procedures, \mathbb{D} can decide to start a new procedural learning \mathbb{L}_P for the task. This process is accomplished by choosing attentional actions $a_A \in \mathbb{A}_A$ that can emphasize attention in the area close to the recent points through top-down actions. Learning can continue typically until the point cloud allow the formulation of the geometrical concept of the object. At each learning step, the state-action pairs are stored in the procedural memory as a recipe for the task of investigating a new object.

VI. RESULTS AND DISCUSSION

The dynamics of the attentional system guided by CONAIM after the procedural learning is shown in Fig. 5. Fig. 5(m) shows the successive application of the top-down actions a_{05} and a_{07} (emphasize stimuli from the left and right sides of the robot) as decided by the cognitive system to guide the agent to investigate the box. The robot trajectory is shown in Fig. 5(a). Fig. 6 compares the stimulus-guided and the conscious-guided agent. Fig. 6(a) presents the trajectory of 100 trials from 8 distinct starting poses considering that the agent was always driven to the most salient stimulus. Fig. 6(b) presents the consciousness-guided agent trajectories during 100 learning trials from 8 distinct starting poses. The agent learned to walk around the object under interest collecting new information about it, until the geometrical concept of the box was formulated in the semantic memory. Fig. 6(c) and (d) presents the number of points collected from the box in both approaches.

A. Model Summary

As shown in the conducted experiment, the proposed model considers:

- 1) Working with *multiple sensors* \mathcal{S} of different categories and dimensions, which are automatically merged by the sensory-attentional system [Fig. 5(n) and (o)].
 - 2) *Attentiveness*, since the agent is able to react to environmental stimuli received through its sensors and modulated according to attentional and cognitive systems (Fig. 5).
 - 3) *Volition* $\mathbb{V}\mathbb{O}$ here modeled as a function able to turn motivations $m_v \in \mathbb{M}_V$ into goals $g \in \mathbb{G}$.
 - 4) Bottom-up *exogenous* and top-down *endogenous* stimuli and the attentional course of selection;
 - 5) The *sustained component of attention*, i.e., to maintain the focus of attention based on the agent's complete conscious cognitive scheme, now guided by the decision-making block (\mathbb{D}).
 - 6) The *selective component of attention*, which raises to memory only a small subset of sensory information located in the vicinity of the attentional focus, dramatically reducing the sensory input space.
 - 7) The *orienting component of attention*, which results on changing the focus of attention on the current cycle.
 - 8) The *overt attention*, as the decision-making module \mathbb{D} can directly adjust the sensory structures through the execution of a motor action $a_M \in \mathbb{A}_M$.
 - 9) The *covert attention*, as the decision-making module \mathbb{D} can choose, if necessary, just for an attentional action $a_A \in \mathbb{A}_A$, without a motor action $a_M \in \mathbb{A}_M$ associated to it.
- The proposed formal model points out the potential for further developments that could support:
- 1) *Sentience*, since the agent can have its emotional state variables $v_E \in \mathbb{S}_E$ altered by the conscious experience;
 - 2) *Self-awareness*, since the model can accommodate an individuality schema that contains internal states and body image, as well as planning, decision-making, and learning mechanisms that support insight and possible comparisons with other individuals;

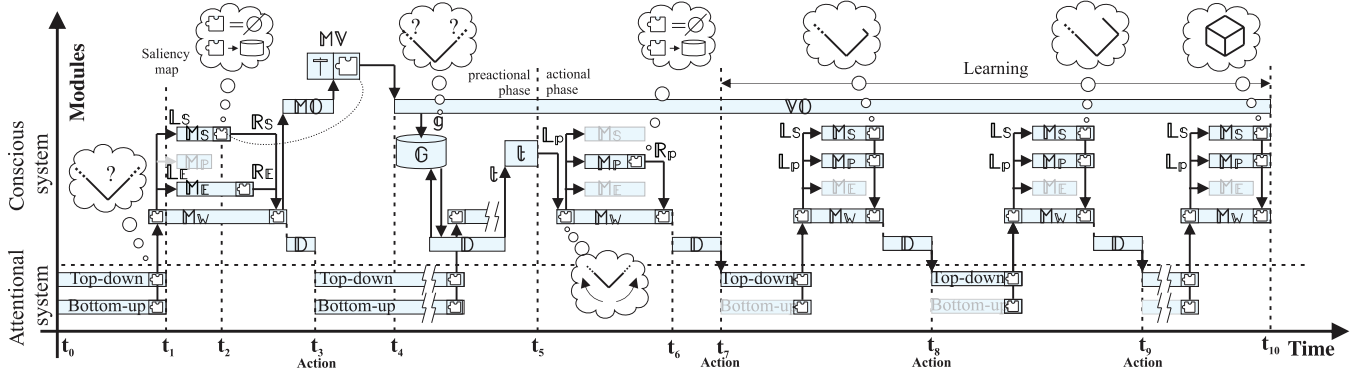


Fig. 4. CONAIM system dynamics during the experiment. t_0 : new stimuli arrive from the environment; t_1 saliency map \mathcal{L} is computed and working memory M_W receives data regarding the new object; t_2 : recalls in memories are stimulated by m_W . Semantic memory recall (R_S) fails indicating that this is a concept new to the agent. A new entry m_S is created in M_S with the collected data. t_3 : after memory recall (R_S and R_E) fails, decision-making D drives the agent to the most salient stimulus. A new attentional cycle is started. Motivation function M_O makes the agent interested on the new concept. t_4 : volition V_O in preactional phase inserts a new goal g at goals list G , and decision-making D turns it into the new agent task l ; t_5 : in order to investigate the new concept, the agent looks for an investigation procedure in procedural memory M_P , and the recall R_P fails. A new entry m_P is created in M_P ; t_6 : decision-making D starts a new learning for m_P by choosing attentional actions a_A that can emphasize stimuli in areas close to recently collected points (top-down influence); t_7-10 : for several interactions, this investigation procedure will be learned by the agent. New collected data about the concept will be stored in m_S . The learning procedure of $m_P \in M_P$ can continue until a complete concept $m_S \in M_S$ can be considered fully formulated.

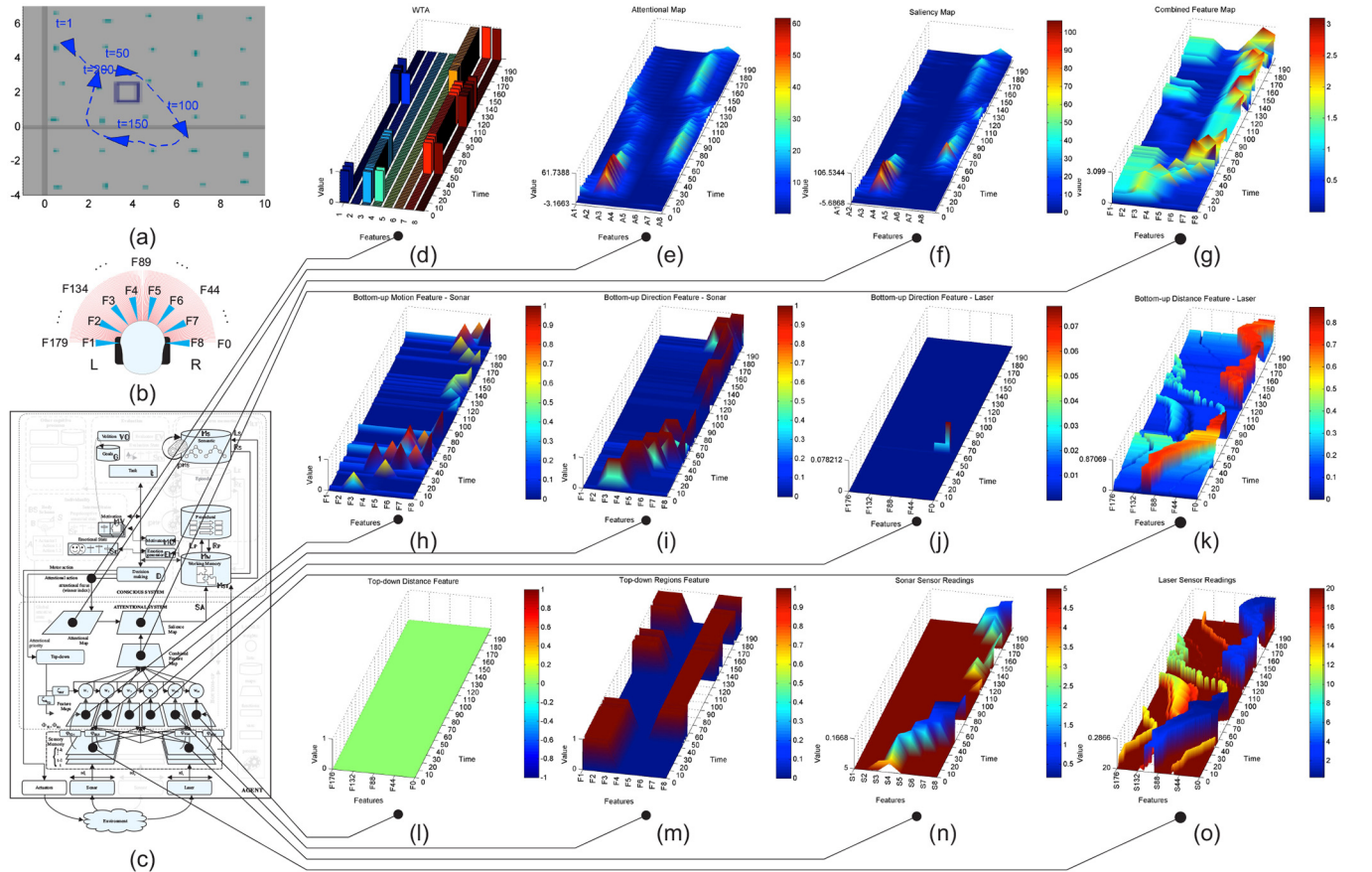


Fig. 5. Dynamics of the attentional system guided by the cognitive system in one of the experiments. (a) Trajectory of the agent around a box from $t = 0$ to $t = 200$ steps (dimensions in meters). (b) Disposition and labels for the eight sonar measurements [F1-F8] (blue) and 180 range scanner measurements [F0-F179] (red). (c) CONAIM architecture adopted in the the experiment (for details see Fig. 2). (d) Index of the *winner* of the attentional focus. (e) Attentional map \mathcal{M} . (f) Saliency Map \mathcal{L} . (g) Combined feature map \mathcal{C} . (h) *Bottom-up* feature map \mathcal{F}_1 (motion discrepancy). (i) *Bottom-up* feature map \mathcal{F}_2 sonar component (moving directions discrepancy). (j) *Bottom-up* feature map \mathcal{F}_2 laser component (moving directions discrepancy). (k) *Bottom-up* feature map \mathcal{F}_3 (distance discrepancy). (l) *Top-down* feature map \mathcal{F}_4 (interest on measurements with a specific distance). (m) *Top-down* feature map \mathcal{F}_5 (interest on specific regions). (n) Sonar readings \mathcal{O}_1 (meters). (o) Range scanner readings \mathcal{O}_2 (m). As shown in (m), after learning, the attentional actions a_{05} (drive attentional focus to right side) and a_{07} (drive the attentional focus to left side) are successively recovered from procedural memory M_P , guiding the robot to investigate the box.

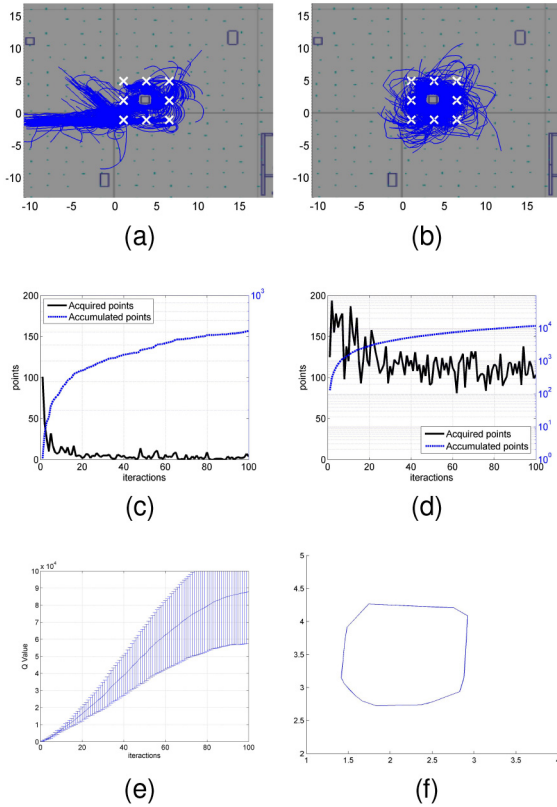


Fig. 6. Agent trajectories when exploring an environment with several boxes starting from eight distinct poses around a central box. (a) Trajectories of the robot guided by the most salient stimulus. (b) Trajectories of the robot guided by the CONAIM system performing the semantic and procedural learning. (c) and (d) Number of points of the central box acquired by the agents in “a” and “b,” respectively, during computational steps and total amount of accumulated points. (e) Average Q values and standard deviation of the agent in “b” along learning trials for all experiments. (f) Example of resulting convex hull applied in the semantic memory of the agent in “b” over the point cloud.

- 3) *Autonoetic consciousness*, since the model can accommodate elements that enable reasoning about the past, present and future, especially when considering the modeling of the working memory proposed by [27];
- 4) *Mineness*, since objects frequently used by the agent can make use of one of its indicators $dms \in Ms$, in order to represent an ownership relationship;
- 5) *Perspectivalness*, since the mental mechanism can be expanded to reason about the perspective of the *self* compared to the other. The adoption of a point of view will depend on the capacity of building concepts in the semantic memory, and the agent’s reasoning over them;
- 6) *Self-consciousness*, since the model can be expanded to allow the formulation of concepts about the surroundings of the agent, the assimilation of experiences and the analysis of the consequences of the agent decisions, usually based on learning and experience. The model could also allow the agent to be aware of its body, to recognize a similar and/or itself and to reason about the consequences of its decisions;
- 7) Other specialized modules such as language processing can be inserted in the other cognitive processes block.

VII. CONCLUSION

This paper proposed a novel model based on an attentional paradigm for machine consciousness consistent with various definitions of the concept available in the literature. The attentional system allows the agent to learn over the saliency map space instead of the traditional sensorial space, promoting a reduction of 95% on the original state-space. The conscious architecture includes cognitive elements such as memories, reasoning, planning, emotion, decision-making, learning, motivation, and volition. The model also provides a methodology for performing computation over these elements. Unlike other studies in the literature, where in most cases the authors briefly present the main concepts behind the modules that compose the system, this work details such modules, their connections, and their ways of operation. The formal model of the system is presented and, to our knowledge, it represents the first known complete formalization of a consciousness integrated model. The organization of the elements and processes involved in consciousness in an architecture also allows further comparison among distinct algorithms, as well as offers a base for testing consciousness, psychological, neuroscience, and machine learning theories, helping us to better understand human and machine consciousness. Experimental results suggest that this model is capable of supporting attentiveness, sentience, and self-awareness based on exogenous and endogenous stimuli. Although in these initial experiments only some basic components were engaged, there is room for the proposition of new functions (motivation, emotion, etc.) and modules that can far improve the agent behavior. Finally, we hypothesise that the model could be extended to demonstrate autonoetic consciousness, mineness, perspectivalness, and self-consciousness.

REFERENCES

- [1] G. Tononi, “An information integration theory of consciousness,” *BMC Neurosci.*, vol. 5, no. 42, pp. 1–22, 2004.
- [2] M. S. Gazzaniga, R. B. Ivry, and G. R. Mangun, Eds., *Cognitive Neuroscience: The Biology of Mind*. New York, NY, USA: Norton, 2002.
- [3] M. A. Cohen, P. Cavanagh, M. M. Chun, and K. Nakayama, “The attentional requirements of consciousness,” *Trends Cognit. Sci.*, vol. 16, pp. 411–417, Aug. 2012.
- [4] J. A. Reggia, “The rise of machine consciousness: Studying consciousness with computational models,” *Neural Netw.*, vol. 44, pp. 112–131, 2013.
- [5] R. Arp, “Consciousness and awareness: Switched-on Rheostats: Response to de Quincey,” *J. Conscious. Stud.*, vol. 14, no. 3, pp. 101–106, 2007.
- [6] R. A. Moreno, A. L. Espino, and A. S. de Miguel, “Modeling consciousness for autonomous robot exploration,” in *Proc. Int. Work-Conf. Interplay Nat. Artif. Comput. (IWINAC)*, 2007, vol. 4527, pp. 51–60.
- [7] J. A. Starzyk and D. K. Prasad, “A computational model of machine consciousness,” *Int. J. Mach. Conscious.*, vol. 3, pp. 255–281, 2011.
- [8] M. Nisargadatta, *I Am That*. Bombay, India: Chetana, 1973.
- [9] D. Gamez, “Progress in machine consciousness,” *Conscious. Cognit.*, vol. 17, pp. 887–910, Sep. 2008.
- [10] E. L. Colombini and C. C. Ribeiro, “An attentive multi-sensor based system for mobile robotics,” in *Proc. IEEE Int. Conf. Syst. Man Cybern. (SMC)*, 2012, pp. 1509–1514.
- [11] E. L. Colombini, *An attentional model for intelligent robotics agent*. Ph.D. dissertation, Instituto Tecnológico de Aeronáutica (ITA), São José dos Campos, SP, Brazil, 2014.
- [12] E. L. Colombini, A. S. Simões, and C. C. Ribeiro, “Top-down and bottom-up features combination for multi-sensor attentive robots,” in *Proc. Int. Symp. Atten. Cognit. Syst. (ISACS); 23rd Int. Joint Conf. Artif. Intell. (IJCAI)*, Beijing, China, pp. 1–14, 2013.

- [13] S. Laureys *et al.*, "Brain function in the vegetative state," *Acta Neurol. Belgica*, vol. 102, pp. 177–185, Dec. 2002.
- [14] R. Van Gulick, "Consciousness," in *The Stanford Encyclopedia of Philosophy*, Zalta E. N., Ed. Stanford, CA, USA: Stanford Univ., Spring 2014.
- [15] D. Armstrong, *What is Consciousness? In the Nature of Mind*. Ithaca, NY, USA: Cornell Univ. Press, 1981.
- [16] J.-P. Broonen *et al.*, "Is volition the missing link in the management of low back pain?," *Joint Bone Spine*, vol. 78, pp. 364–367, Jul. 2011.
- [17] G. Tononi, "Consciousness as integrated information: A provisional manifesto," *Biol. Bull.*, vol. 215, pp. 216–242, 2008.
- [18] D. Gamez, "Information integration based predictions about the conscious states of a spiking neural network," *Conscious. Cognit.*, vol. 19, pp. 294–310, Mar. 2010.
- [19] B. J. Baars, *A Cognitive Theory of Consciousness*. Cambridge, U.K.: Cambridge Univ. Press, 1988.
- [20] B. Baars, "The conscious aspects hypothesis," *Trends Cognit. Sci.*, vol. 6, pp. 47–52, 2002.
- [21] B. Baars and S. Franklin, "Conscious is computational: The LIDA model," *Int. J. Mach. Conscious.*, vol. 1, pp. 23–32, 2009.
- [22] S. Franklin, A. Kelemen, and L. McCauley, "IDA: A cognitive agent architecture," in *Proc. IEEE Int. Conf. Syst. Man Cybern. (SMC)*, 1998, pp. 2646–2651.
- [23] S. Franklin *et al.*, "Lida: A computational model of global workspace theory and developmental learning," in *Proc. AAAI Fall Symp. AI Conscious. Theor. Found. Curr. Approaches*, 2007, pp. 1–6.
- [24] R. Arrabales, A. Ledezma, and A. Sanchis, "Towards conscious-like behaviour in computer game characters," in *Proc. IEEE Symp. Comput. Intell. Games*, 2009, pp. 217–224.
- [25] D. Perlis, "Consciousness as self-function," *J. Conscious. Stud.*, vol. 4, pp. 509–525, 1997.
- [26] D. McDermott, "Artificial intelligence and consciousness," in *Cambridge Handbook of Consciousness*, Cambridge, U.K.: Cambridge Univ. Press, 2007, pp. 117–150.
- [27] A. Samsonovich and L. Nadel, "Fundamental principles of mechanisms of the conscious self," *Cortex*, vol. 41, pp. 669–689, 2005.
- [28] O. Holland, "A strong embodied approach to machine consciousness," *J. Conscious. Stud.*, vol. 14, pp. 97–110, 2007.
- [29] J. Takeno, "A robot succeeds in 100% mirror image cognition," *Int. J. Smart Sens. Intell. Syst.*, vol. 1, pp. 891–911, 2008.
- [30] J. Takeno, *Creation of a Conscious Robot*. Singapore: Pan Stanford, 2013.
- [31] R. Sun, "Accounting for the computational basis of consciousness," *Conscious. Cognit.*, vol. 8, pp. 529–565, 1999.
- [32] R. Sun and S. Franklin, "Computational models of consciousness," in *Cambridge Handbook of Consciousness*, P. Zelazo and M. Moscovitch, Eds., Cambridge, U.K.: Cambridge Univ. Press, 2007, pp. 151–174.
- [33] T. Kitamura, T. Tahara, and K. Asami, "How can a robot have consciousness?," *Adv. Robot.*, vol. 14, pp. 263–275, 2000.
- [34] A. Chella, "Towards robot conscious perception," in *Artificial Consciousness*, A. Chella and R. Manzotti, Eds., New York, NY, USA: Academic, 2007, pp. 124–140.
- [35] A. Chella, M. Frixione, and S. Gaglio, "A cognitive architecture for robot self-consciousness," *Artif. Intell. Med.*, vol. 44, pp. 147–154, 2008.
- [36] M. I. Posner and S. Dehaene, "Attentional networks," *Trends Neurosci.*, vol. 17, pp. 75–79, 1994.
- [37] M. Velmans, *The Science of Consciousness*. Evanston, IL, USA: Routledge, 1996.
- [38] P. M. Merikle and S. Joordens, "Parallels between perception without attention and perception without awareness," *Conscious. Cogn.*, vol. 6, pp. 219–236, 1997.
- [39] J. K. O'Regan and A. Noe, "A sensorimotor account of vision and visual consciousness," *Behav. Brain Sci.*, no. 24, pp. 939–973, 2001.
- [40] C. Koch and N. Tsuchiya, "Attention and consciousness: Two distinct brain processes," *Trends Cognit. Sci.*, vol. 11, pp. 16–22, Jan. 2007.
- [41] J. G. Taylor, "The creativity effect: Consciousness versus attention," in *Proc. Int. Joint Conf. Neur. Netw. (IJCNN)*, 2010, pp. 1–8.
- [42] C. Koch, *The Quest for Consciousness: A Neurobiological Approach*. Greenwood Village, CO, USA: Roberts and Publishers, 2004.
- [43] V. G. Hardcastle, "Attention versus consciousness: A distinction with a difference," *Cognit. Stud. Bull. Jpn. Cognit. Sci. Soc.*, vol. 4, pp. 56–66, 1997.
- [44] V. A. Lamme, "Why visual attention and awareness are different," *Trends Cognit. Sci.*, vol. 7, pp. 12–18, 2003.
- [45] G. F. Woodman and S. J. Luck, "Dissociations among attention, perception, and awareness during object-substitution masking," *Psychol. Sci.*, vol. 14, pp. 605–611, 2003.
- [46] C. Tallon-Baudry, "On the neural mechanisms subserving consciousness and attention," *Front. Psychol.*, vol. 2, pp. 397–397, Jan. 2012.
- [47] S. Grossberg, "Linking attention to learning, expectation, competition, and consciousness," Dep. of Cog. and Neur. Syst., Boston Univ. Center for Adap. Sys., Tech. Rep. CAS/CNS-2003-007, Boston, Apr. 2010.
- [48] C. J. Tinsley, "Using topographic networks to build a representation of consciousness," *Biosystems*, vol. 92, pp. 29–41, Jan. 2008.
- [49] P. Haikonen, *The cognitive approach to conscious machines*. New York, NY, USA: Academic, 2003.
- [50] P. O. A. Haikonen, *Robot Brains. Circuits and Systems for Conscious Machines*. Hoboken, NJ, USA: Wiley, 2007.
- [51] P. Haikonen, "Essential issues of conscious machines," *J. Conscious. Stud.*, vol. 14, pp. 72–84, 2007.
- [52] J. A. Starzyk and D. K. Prasad, "Machine consciousness: A computational model," *Brain-Inspired Cognit. Syst. (BICS)*, Madrid, Spain, Jul. 14–16, 2010, pp. 428–439.
- [53] D. Stachowicz and G. M. Kruijff, "Episodic-like memory for cognitive robots," *IEEE Trans. Auton. Ment. Develop.*, vol. 4, no. 1, pp. 1–16, Mar. 2012.
- [54] P. Ekman, "Basic emotions," in *Handbook of Cognition and Emotion*, D. G. Leish T. and Power M., Eds. Hoboken, NJ, USA: Wiley, 1999.
- [55] J. Panksepp, "Affective consciousness: Core emotional feelings in animals and humans," *Conscious. Cognit.*, vol. 14, pp. 30–80, 2005.
- [56] T. Ziemke and R. Lowe, "On the role of emotion in embodied cognitive architectures: From organisms to robots," *Cognit. Comput.*, vol. 1, pp. 104–117, 2009.
- [57] A. R. Damasio, *Looking for spinoza: Joy, sorrow and the feeling brain*. Orlando, FL, USA: Harcourt, 2003.
- [58] J. Wang, M. Lewis, and S. Hughes, "Validating USARSim for use in HRI research," in *Proc. Hum. Factors Ergon. Soc.*, 2005, pp. 457–461.



Alexandre da Silva Simões (M'95) received the B.E. (1998) degree in electrical engineering from UNESP, Bauru, Brazil, and the M.Sc. (2000) and Ph.D. degrees (2006) in electrical engineering from the University of São Paulo (USP), São Paulo, Brazil.

He is a Full Professor with the Department of Control and Automation Engineering, UNESP. He is currently the Vice Director of the Sorocaba Campus of UNESP. His research interests include robotics, machine consciousness, machine learning, arts, and education. Dr. Simões was the General Chair of RoboCup 2014 and the founding president of RoboCup in Brazil.



Esther Luna Colombini received the M.Sc. (2005) and Ph.D. degrees (2014) in computing engineering from Technological Institute of Aeronautics, São Paulo, Brazil.

She is an Associate Professor with the University of Campinas (UNICAMP). Her research interests include machine learning, autonomous robotics, cognitive modeling, and robotics in education.

Dr. Colombini is the President of RoboCup Brazil and the Organizer of the Brazilian Robotics Olympiad, a public initiative that involves 100 000

students supported by the National Council of Technological and Scientific Development.



Carlos Henrique Costa Ribeiro received the M.Sc. (1993) degree in electrical engineering from Federal University of Rio de Janeiro, Rio de Janeiro, Brazil, and the Ph.D. degree (1998) in electrical and electronic engineering from Imperial College London, London, U.K.

He is an Associate Professor with the Technological Institute of Aeronautics, São Paulo, Brazil. He has been a Researcher with the National Council of Technological and Scientific Development (CNPq), Brasília, Brazil, since 2001. His research interests

include machine learning (especially algorithms for autonomous learning), multiagent robotics, complex networks, and engineering education.