



Sensorimotor Representation Learning for an “Active Self” in Robots: A Model Survey

Phuong D. H. Nguyen¹ · Yasmin Kim Georgie² · Ezgi Kayhan^{3,4} · Manfred Eppe¹ · Verena Vanessa Hafner² · Stefan Wermter¹

Received: 12 June 2020 / Accepted: 13 January 2021
© The Author(s) 2021

Abstract

Safe human-robot interactions require robots to be able to learn how to behave appropriately in spaces populated by people and thus to cope with the challenges posed by our dynamic and unstructured environment, rather than being provided a rigid set of rules for operations. In humans, these capabilities are thought to be related to our ability to perceive our body in space, sensing the location of our limbs during movement, being aware of other objects and agents, and controlling our body parts to interact with them intentionally. Toward the next generation of robots with bio-inspired capacities, in this paper, we first review the developmental processes of underlying mechanisms of these abilities: The sensory representations of body schema, peripersonal space, and the active self in humans. Second, we provide a survey of robotics models of these sensory representations and robotics models of the self; and we compare these models with the human counterparts. Finally, we analyze what is missing from these robotics models and propose a theoretical computational framework, which aims to allow the emergence of the sense of self in artificial agents by developing sensory representations through self-exploration.

Keywords Developmental robotics · Body schema · Peripersonal space · Agency · Robot learning

1 Introduction

In order to bring robots to safely cooperate with humans within the same environment, the next generation of robots needs to be equipped with the abilities to learn, adapt, and act autonomously in unstructured and dynamic environments. In other words, we want robots to operate in the same situations and conditions as humans do, to use the same tools, to interact with, and to understand the same world in which humans’ daily lives take place. For achieving this, we want robots to be able to learn how to behave appropriately in our world and thus to find efficient ways to cope with

the challenges posed by our dynamic and unstructured environment, rather than providing a rigid set of rules to drive their actions. Robots should be able to deal efficiently with unexpected changes in the perceived environment as well as with modifications of their own physical structure in a scalable manner. So how can we build robots that possess such abilities?

We address this question by reviewing interdisciplinary research related to the developmental processes that form the representations of body schema and the peripersonal space (PPS), the sense of agency, and by discussing how these relate to the active self. We first review the development of body schema, sense of agency, and the PPS in humans in Sect. 2. Also, we highlight that the body schema and the PPS representations emerge by exploration and that they are critical for the development of agency and higher cognitive functions. Then, in Sect. 3, we discuss the behavioral function and properties of the body schema and PPS representations in humans and review the state of the art models in developmental robotics concerning these representations. Finally, in Sect. 4, we analyze these issues with respect to the so-called “active self” and conclude in Sect. 5 by proposing a general blueprint that builds on the *verification*

✉ Phuong D. H. Nguyen
p.nguyen@informatik.uni-hamburg.de

¹ Department of Informatics, University of Hamburg,
Vogt-Koelln-Str. 30, 22527 Hamburg, Germany

² Department of Computer Science, Humboldt-Universität zu
Berlin, Unter den Linden 6, 10019 Berlin, Germany

³ Department of Developmental Psychology, University
of Potsdam, Potsdam, Germany

⁴ Max Planck Institute for Human Cognitive and Brain
Sciences, Leipzig, Germany

principle by Stoytchev [167] to overcome the limitations of current robotic systems. In the remainder of this section, we provide a brief overall background and similar reviews in developmental robotics.

The main scope of this review paper focuses on the first phases of the development of the self in humans. We consider especially the physical mechanisms of the development and models inspired by these developmental processes. Therefore, our attention is on computational and robotics models proposed for humanoid robots or of equivalent configuration, i.e., having cameras as eyes, manipulators as arms, and possibly a tactile-covered system as artificial skin. We acknowledge that social interactions might have some influence on the development processes of multisensory representations and the self, which are discussed elsewhere, e.g., [34, Section 6], [161, Section 6], [63, 115, 116, 172, 176]. We leave the systematic review of social aspects for future work.

1.1 The Active Self and The Sense of Agency

The notion of the active self relates to the connections between perception, action, and prediction, and how these connections facilitate the emergence of a *minimal self*. For the term “active self”, we argue that the sensorimotor activities of an agent are a prerequisite for the emergence of the minimal self, in the sense that “the phenomenal, minimal self is empirically derived from sensorimotor experience and not a theoretical and empirical given” [182]. The minimal self, or “minimal phenomenal selfhood” [10], refers to the pre-reflective sense of being a self as being subject to immediate experience [68]. This minimal notion of the self involves the sense of agency—the sense of the self as the one causing or generating action, and the sense of ownership—the sense of the self as the one subjected to an experience [68].

The sense of agency and body ownership are emergent properties of a complex, embodied system that is situated in a dynamic environment that has a level of uncertainty. However, one can argue that the level of “complexity” of the environment is related to the system’s own sensory and motor capacities. Simply put, the more information a system can perceive from the interaction with the environment, and the “richness” with which the system can act upon the environment, the more “complex” the information from the environment will be to that system [138]. Thus, information is formed by the *interaction*, rather than being “provided” by the environment and decoded by the perceptual system of the agent. It follows then that the properties of the system, the body, along with the properties of the environment, govern the interaction between the embodied agent and the environment in which it is situated (the ecological niche), as well as the developmental process of the agent itself. Infants are born into a dynamic, uncertain environment with which the

interaction is complex. However, human infants (as well as other complex biological systems) are not born with complete pre-existing knowledge about their environment, nor their own body. Infants construct this knowledge over time, and progressively form a model of the body—a body representation, and a model of the environment through interactions ([128, 181]; also see [82, 90] for a review).

1.2 The Body Schema and The Peripersonal Space

In humans, the capabilities to deal with unexpected changes in the environment and modifications of our own physical structure (e.g., growth or by extending an arm with a tool) emerge from our ability to perceive our body in space, sensing the location of our limbs during movement, being aware of other objects and agents, and controlling our body parts to interact with them intentionally. These abilities are thought to be related to the presence of a body schema, peripersonal space (PPS), and the minimal self including the sense of body ownership and sense of agency.

As initially defined by Head and Holmes [79], the body schema is a sensorimotor representation of the structure and position of the human body, which is encoded in the brain and allows the agent to perform body movements. Also maintained by the brain, the PPS denotes the representation of the proximal space surrounding the agent’s body. This space is commonly defined as the reachable space but outside the body surface, differentiated from the extrapersonal space and the personal space. Specifically, PPS is the space where all motor activities of an agent such as object manipulations take place [161]. For example, consider grasping an external object in the reachable space of a robot. To execute this action, an agent requires two prerequisites: First, it needs to be aware of and monitor the position of its body part, e.g., a limb to execute the movement. Second, it needs to “compute” the dynamic position, dimension, etc., of the target with respect to the agent’s body. The brain provides awareness about the body schema and body configuration, and the computation of the target location is a result of the brain’s PPS representation. These two representations, the body schema and the PPS, emerge from the low-level integration of different sensory modalities available in a human body (see Table 1 for details). They are closely related and interact with each other.

Indeed, there are some overlapping functionalities of the body schema and the PPS representation, namely (1) they are both multisensory representations; (2) they convey frame of reference (FoR) transformations; (3) they have a strong link with actions within the reachable space and (4) the representations are plastic. According to [26, 27] these overlaps are potentially due to their causal relation (the extension of the body representation leads to the extension of PPS representation in tool-use cases) or their unique identity.

Table 1 Sensory inputs and functionality of body schema, body image, and the PPS representations

	Body schema	Body image	PPS
Sensory Sources	Proprioception	Proprioception	Vision
	Touch	Vision	Audio
	Vision		Touch Proprioception*
Functionality	Body structure	Body perception	Involuntary actions
	Body pose	Body conception	Voluntary actions

This table reproduces the results reported in [27, 45, 161]. Serino [161] suggests that the interaction of proprioceptive^(*) and visual signals about the body part is vital for frame transformations in the dynamic cases of the PPS

Nevertheless, the differences between the two representations still exist and stem from the involvement of external objects within reach in the environment¹ causing the non-bodily stimuli. Because of the requirement of body continuity, there are also certain tool-use cases, e.g., a remotely controlled tool like the computer mouse and its pointer, in which the body schema representation cannot include the tool. Hence, the representation is not affected. Instead, the spatial representation of PPS would be modulated due to the availability of visual-tactile correlation and action-effect association [27]. The recent behavioral study of D’Angelo et al. [52] suggests that there are separate mechanisms for the plastic changes of body schema and PPS representations.

As reviewed by de Vignemont [45], the representation of an agent’s body can be distinguished into body schema—the representation for actions², and body image—other body-related representations for perception, conception, and emotion (according to the *dyadic taxonomy*) (see also [50]). The body image can be further separated into two distinct representations, namely visuo-spatial body map—the structure description of body parts, and body semantics—the conceptual and linguistic level of body parts (according to the *triadic taxonomy*). However, with the perspective of the enactive approach [181], in which the sensorimotor exploration gives rise to perceptual experiences, the distinction between the bodily action-oriented and perceptual representation is quite blurry. For example, the visual appearance and boundary of a limb would have an effect on the agent’s perception of the length and position of the limb. Hence, it is reasonable to include the body structure description of the body image representation (and its sources of sensory information) when

¹ in tool-use cases, external out-of-reach objects also get involved

² Action-oriented representation is defined by the author as “it carries information about the bodily effector (and the bodily goal in reflective actions) that is used to guide bodily movements.”

considering the body schema in action, especially from the computational perspective. Indeed, most robotics models of the so-called body schema fall into this category (see Sect. 3.2 and Table 2). Furthermore, from the technical point of view, it is difficult to model the mental level of the body image when the definition is unclear. Therefore, we will use the term “body schema” in an extended meaning including both “body schema” and “body structure description”.

1.3 Developmental Processes of Emergent Selfhood and Body Awareness

Although the bodily senses exist in human adults as a result of multisensory integration processes, these abilities are not innate—newborns and infants develop these abilities over time [17]. Indeed, the senses of the bodily self, i.e., the sensation of the position of a body part, the surrounding space, and the feeling of owning and controlling one’s body, incrementally develop in newborns in the very first months of their lives, (e.g., [15, 17, 129, 146, 147, 111, 116]).

The sense of bodily self is the result of the gradual emergence of several abilities in infants. These abilities include perceiving multisensory spatial contingencies (i.e., visual-proprioceptive or visuotactile) soon after birth (e.g., [8]), spatial postural remapping [18], visual-elicited reaching movement [40], and goal-directed exploration behaviors [182].

Taken together, it is reasonable to argue that after birth, infants spend their first months of life undergoing many developmental milestones to incrementally develop the representation of their body. This body schema is related mainly to touch, proprioception, and vision (see Table 1) as these sensory modalities continue to develop from the fetal stage (see [2, 81] for reviews). Later on, the representation of the surrounding space of the body—the PPS—is aggregated from the proprioceptive and exteroceptive modalities (see Table 1). In addition, infants develop the capability to generate motor actions corresponding to desired outcomes, and the ability to distinguish between self and other, both related to the senses of body ownership and agency. At first, these developments may be triggered by self-exploration movements. However, then the enhanced perceptual capability may help infants in improving their motor control, from a reflexive manner to an intentionally goal-directed state during these processes (see Fig. 1). Insights from the developmental dynamics of these abilities may suggest important prerequisites for formulating developmental models of artificial intelligence.

1.4 Related Work

Several reviews relate to the topic of this paper. First, the review by Hoffmann et al. [82] on robotic models of body

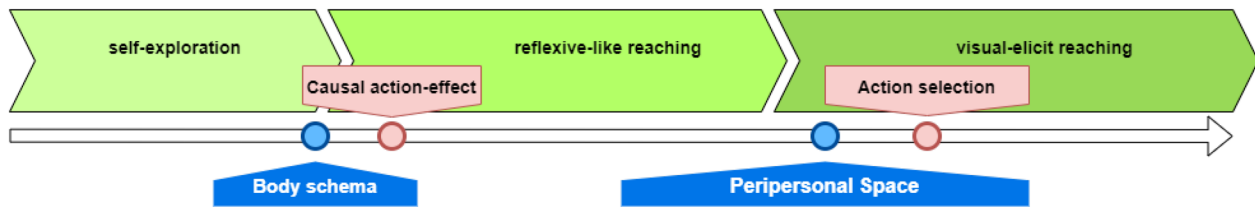


Fig. 1 Development path of different body-related representations and sensations in the first year in infants. The development of the body schema is reviewed and discussed in Sect. 2.1, from the fetal stage to the stage around 3 months of age after birth. The development of the PPS is discussed in Sect. 2.3, which is suggested to con-

schema surveyed the concept of body schema in biology, its properties, and its relation with the forward models used in the field of robotics. The review also provides a thorough overview of body schema-inspired robotic models. In this work, we will briefly review the body schema properties and further provide a complementary view on this sensory representation. Furthermore, we will provide an update on robotic models of the body schema representation.

Cangelosi and Schlesinger [24] presented both theoretical and experimental aspects of the developmental robotics approach. The approach promotes the idea of building artificial agents by receiving inspiration from human developmental science. The authors outline the theoretical principles of the approach including embodiment, enaction, cross-modality, and online, cumulative, open-ended learning. The experimental review provides an overview of developmental robotics models from intrinsic motivation, perceptual and motor developments, to social learning, language skills, and abstract knowledge developments.

Moreover, Schillaci et al. [159] reviewed and suggested the fundamental role of sensorimotor interaction in the development of both human and artificial agents. In this process, the agent's motor exploration in a situated environment serves as a means for gathering sensorimotor experiences, which facilitates the emergence of other cognitive functions. For example, sensorimotor experiences are used to learn a forward model, and a forward model can be the basis for learning high-level cognitive conceptual representations. In agreement with [159], we aim to delve into the role of multisensory information collected through exploration in the formation of an agent's body and peripersonal space representation, and how these sensorimotor representations affect the agent's sense of the active self, including the sense of agency and the sense of body ownership. Thus, motor explorations will be mentioned but not exhaustively discussed in this surveyed work. Instead, we focus on the body schema, the peripersonal space, and the emergence of the sense of agency.

Georgie et al. [71] discussed the development of body representations as a prerequisite for the emergence of the minimal self, which includes body ownership and agency.

tinue from 3 to 6–10 months of age. The development of the active self, which relates to causal action-effect and action selection, is discussed in Sect. 2.2. Literature suggests that the active self, as a process of self emergence, takes place from birth to 9 months of age.

They discuss some of the behavioral measures indicating the presence of body ownership and sense of agency in humans and survey some of the related robotics research that examined and developed these concepts. In their review, the authors suggested possible expansions to the robotics research for exploring the development of an artificial minimal self. Specifically, to focus on developing models that incorporate a whole developmental path in a real robot that would include, e.g., self-exploration and self-touch, where behavioral indices can be measured at different points along the developmental path.

Concurrently with this paper, Tani and White [171] review models of the sense of minimal and narrative self in cognitive neurorobotics, but mainly focus on models utilizing RNN architectures that follow the free-energy principle and active inference approach [65].

2 Development of the Body Schema, Peripersonal Space, and The Sense of Agency in Humans

Before reviewing robotic models of the body schema, PPS representations, and the sense of agency, we first consider the development of these representations and agency in humans. This involves the development of the body schema from gestation to infancy in Sect. 2.1, the PPS representation in infants in Sect. 2.3, and the emergence of the sense of agency in infants in Sect. 2.2.

2.1 Development of the Body Schema from Gestation to Infancy

The development of the body schema is inseparably linked with sensorimotor development, starting from as early as the fetal stage and continuing later on after birth. The body schema's neural foundation is formed by the neurological representations of the different anatomical divisions of the body. These are the cortical "homunculi" (see Fig. 2 for an illustration) in the primary sensory (S1) and motor (M1) cortices

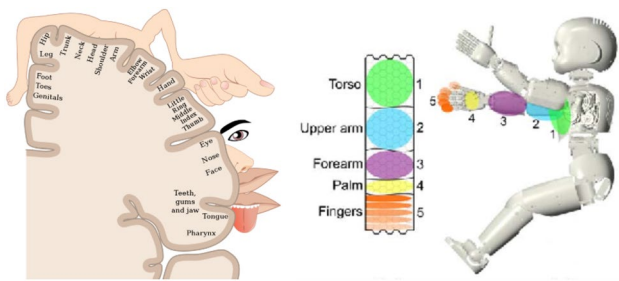


Fig. 2 The cortical homunculi in humans—left and an iCub robot—right (from [83, 127]).

[135]. The different anatomical divisions of the body are mapped onto brain areas in charge of sensory and motor processing along the S1 and M1 cortices. The organization of these specialized areas is realized in a somatotopic map, where adjacent body parts are represented closely together (for the most part—see [136], but also [47]). Moreover, the extent of the cortex dedicated to a body region is proportional to the density of innervation in that specific part (e.g., the mouth and palms) rather than to its size in the body. The establishment of the somatotopic organization in S1 and M1 is facilitated by genetic factors, and later refined through connectivity changes driven by embodied interactions both before and after birth [41].

In terms of motor development, fetuses in the first weeks of gestation typically display different types of motor patterns such as spontaneous startles that start at 7–8 weeks, general movements which start at 8 weeks, isolated movements that emerge soon after, and twitches which start at 10–12 weeks and are produced during active sleep [55]. These very early motor patterns seem to be spontaneous rather than responses to sensations. However, the first sense to develop in the fetus is the tactile sense [12], where fetuses are in a state of constantly being touched by their environment, the tactile sense develops at around the same time as motor movements. Once sensory receptors develop, the fetus' spontaneous movements inevitably lead to sensations, thus facilitating the formation of contingencies between movements and their sensory outcomes [55]. Also, fetuses engage in self-touch in the womb: They often touch body parts that are highly innervated and therefore most sensitive to touch such as the mouth and feet, and later on other parts of the body. The early tendency for movements and self-touch in parts of the body, which are more sensitive, points to a certain preference towards movements that induce more informative sensations (for a review on fetal sensorimotor development see [55]).

Positron emission tomography (PET) studies revealed that in infants under 5 weeks after birth, the dominant metabolic activity is in subcortical regions and the sensorimotor cortex, and by 3 months, metabolic activity increases in the

parietal, temporal, and dorsolateral occipital cortices [30]. It seems that at around 2 months after birth, behavioral control transitions from subcortical to cortical systems. Also, subcortical regions such as the superior colliculus have been investigated as a hub for multimodal integration in human and animal studies [7]. Specifically, the superior colliculus has been implicated as able to support social behavior in early infancy [139], due to its role in attentional behaviors [166, 179].

It seems that the ability to predict sensory consequences of actions, and subsequently to form sensorimotor contingencies begins to develop already in the uterus. There is evidence of fetus anticipation behavior of hand-to-mouth touch already at 19 weeks [121, 143], indicating the presence of a sort of sensorimotor mapping and inference. And from 22 weeks after gestation, movements seem to show an early form of goal-directedness, when the properties of a movement differ depending on the actions' target (more careful movement towards the eye than towards the mouth) [202]. In turn, these in-utero embodied interactions are thought to lay the foundation for the later integration of tactile-proprioceptive and visual information after birth. Using an embodied brain model of a human fetus in a simulated uterine environment, Yamada et al. [194] showed how these interactions promote the cortical learning of body representations by way of regularities in sensorimotor experiences and instantiate postnatal visual-somatosensory integration.

Right after birth, there is a certain regression in motor control, possibly due to the fundamental change in the environment: The newborn has to adapt to an aerial environment in which gravity is felt more strongly, and to the sudden change in brightness and is highly preoccupied with bodily functions such as feeding, sleeping, and crying [55]. Nonetheless, hand-mouth coordination continues to develop after birth. Infants seem to frequently explore their body at around 2 or 3 months, and from birth to 6 months, infants display self-touch progressively throughout their body, from frequently touching rostral parts such as the head and trunk to more caudal parts of the body such as the hips, legs, and feet [174]. From the evidence brought forth by Rochat and Morgan [118, 146, 147], it seems that infants develop the ability to perceive multisensory spatial contingencies (e.g., visual-proprioceptive or visuotactile) soon after birth (e.g., [8]; see also [17] for a review), and also form the perceptual body schema (via intermodal calibration) by 3 months old.

While evidence from neural development studies suggests that even before birth, the prenatal brain should be able to perceive information arising from the body—a rudimentary body schema involving tactile and proprioceptive information—the later maturation of cortical association areas constitutes higher level (multimodal) representations that are possibly formed during the first year after birth [81]. As Hoffmann [81] writes “However, the formation of more

holistic multimodal representations of the body in space occurs probably only after birth, in particular from about 2–3 months”.

Studies show that infants develop a body schema from early on in life allowing them to form expectations about how their bodies look and where they are located in space [146]. From 3 months of age on, when presented with a real-time display of their own legs, infants look longer at an unfamiliar, third-person perspective of their legs than at a familiar, first-person view [148]. Longer looking times of infants were interpreted such that infants expected the images to match their own body schema, thus, they were surprised when their expectations were violated in case of a mismatch between what they expected and what they observed on the display.

Others provided further evidence on infants’ body representations using an adapted version of the rubber-hand illusion paradigm [201]. In the first experiment, infants observed two adjacent displays of baby doll legs being stroked while their own leg was also stroked simultaneously. In the contingent display, the stroking of the infant’s own leg corresponded to the movements on the display whereas, in the non-contingent displays, there was a mismatch between the felt and observed stroking of the leg. Results showed that 10-month-old infants, but not 7-month-olds, looked longer at the contingent displays suggesting that at 10-months of age infants detected visual-tactile contingencies necessary for the identification of self-related stimuli. In this study, longer looking times were interpreted as indicating the early ability to detect visual-tactile contingencies

To find out whether morphological properties of the body facilitated the detection of visual–tactile contingencies, Zmyj et al. [201] ran a control experiment with 10-month-old infants in which infants observed wooden blocks instead of baby doll legs, which were stroked in synch or out of synch with their own leg. Data revealed that infants looked equally long at both contingent and non-contingent displays suggesting that they were able to detect visual-tactile contingencies only when the visual information was related to the body [201].

This preference for specifically body-related synchrony was also later found in newborns. Filippetti et al. [57] investigated the role of temporal synchrony in multisensory integration, to examine whether body-related temporal synchrony detection plays a role even from birth. In two experiments, Filippetti et al. presented newborns (from as early as 12 h after birth) with temporally synchronous and asynchronous visual-tactile stimulation. The visual information was either body-related (an upright newborn face in experiment 1) or non-body-related (an inverted newborn face in experiment 2). Preference or increased attention to the stimuli was measured by a longer looking time. Newborns showed a preference to the synchronous visual-tactile stimulus, only in the body-related condition, indicating that

this increased attention or preference was present only when the synchrony was related to their own body, rather than a general preference to synchrony. The results provide another piece of evidence to the notion that even right after birth, newborns are able to integrate multisensory information and detect synchronous multisensory stimulation, processes that are fundamental for body representations.

In another study, Filippetti et al. [59] presented newborns with videos of newborn faces being stroked with a paintbrush in either a spatially congruent or incongruent location of tactile stimulation. The newborns showed a preference towards the spatially congruent visual-tactile stimulation, suggesting that even shortly at birth, newborns are sensitive to visual-tactile multisensory information. These two studies showed that the ability for detecting temporal and spatial contingencies in multisensory information is present even shortly after birth, and it is present even without self-generated movement.

It is worth pointing out that besides methodological differences between studies (i.e., age groups, sample size), the different feedback modalities (i.e., visual-tactile vs. visual-proprioceptive) and task complexity might have played a role in different looking-time responses in infants. More research with different measures (e.g., pupillometry, EEG, etc.) is needed to clarify this point.

Following up on [57], Filippetti et al. [58] ran an fNIRS study to investigate the brain regions involved in visual-tactile contingency detection for body ownership in infants. 5-month-old infants observed either real-time or delayed videos of themselves while they received tactile stimulation on the cheek with a soft brush. Data revealed that infants showed bilateral activation over the superior temporal sulcus (STS), temporoparietal junction (TPJ), and inferior frontal gyrus (IFG) cortical regions in the contingent condition in response to visual-tactile (and visual-proprioceptive) contingencies. This finding shows that infants as young as 5 months of age show activation in brain regions similar to that of adults when they process information related to their own bodies.

Recently, employing neuroscience techniques, Marshall, Meltzoff, and colleagues conducted a set of experiments in infants’ representations of bodies at the neural level (see [111, 116] for reviews). Using EEG, Saby et al. [155] state that a group of 7-month-old infants shows some somatotopic patterns as the homunculi map: Tactile stimuli in infants’ feet corresponds to response in the midline area of the brain, whereas stimuli in their hands yield responses in lateral central areas. Even a younger group of infants (of 60-day-old) shows brain response when being touched in their hand, foot, and upper lip [117]. Especially, the magnitude of the response to lip touch is much higher than the responses to hand or foot touch, suggesting the tactile sensitivity of the lip area after birth.

2.2 Emergence of Sense of Agency

Developmental researchers have pointed towards two potential underlying mechanisms explaining how infants become agents over their bodies and the environment, namely (1) associative learning and (2) a causal representation of the world.

One line of research emphasized an associative learning mechanism that enables infants to detect the sensory contingencies in their environment. Although their focus in the paper was the memory functions of infants, the seminal work by Rovee and Rovee [153] has revealed some of the early findings on infants' sense of agency. In their mobile-paradigm experiments, infants at around 3 months of age laid in a crib above which a mobile was hanging. One of the limbs of the infant was connected to the mobile with a ribbon. In the connect phase, when the infant moved the connected limb, this resulted in the movement of the mobile. Infants moved their connected limb with increasing frequency when the limb was connected to the mobile, but not when the connected limb was switched or when there was a delay between the movement of the limb and the effect. Interestingly, infants showed increased kicking movement when the mobile was disconnected suggesting that they were trying to re-elicite the effect [154]. Using the mobile paradigm, Watanabe and Taga [186] have shown that whereas 2-month-old infants produced increased movement in all limbs as compared to a baseline period, by the age of 3–4 months, they showed increased movement only in the connected limb to activate the mobile [186]. These findings were interpreted such that at around 3 months of age infants learned the causal link between self-produced movements and their effects in the environment as an indication of “a sense of self-agency” [187]. Other researchers investigated infants' sense of agency in using different paradigms [149]. For example, they measured infants' sucking on a dummy pacifier to investigate whether 2-month-old infants showed differential oral activity based on auditory feedback. In the analog condition, each time infants sucked on the pacifier, they heard a pitch variation of the sound corresponding to the oral pressure applied on the pacifier. In the non-analog condition, each time infants applied pressure on the pacifier, they heard a random pitch variation. Data revealed that 2-month-old infants produced more frequent oral pressure on the pacifier when the auditory effect matched their sucking behavior suggesting that they detected the link between their sucking behavior and the sound effect.

Another line of research emphasized the causal representation of actions and their effects underlying the sense of agency. Researchers ascribing to this view argue that an associative learning mechanism would not be sufficient to account for infants' sense of agency because sense of agency requires a causal representation of the world [196]. Since the

behavioral patterns such as increased movement frequency when connected to a mobile can be explained by alternative mechanisms, these findings provide no evidence for infants' causal representations of their actions and the effects, i.e., sense of agency. Zaadnoordijk et al. [196] simulated the mobile paradigm with a “babybot” that functioned on operant conditioning, thus, it did not have a causal representation of itself and its environment to guide its actions. The simulation results showed that the non-representational babybot produced increased movement with the connected limb as compared to the baseline level of that limb as well as other unconnected limbs. That is, even in the absence of a causal model of the world, the babybot replicated the behavioral findings observed in infant experiments that have been interpreted as evidence for a sense of agency. However, unlike infants, the babybot did not increase its movement rate when the mobile was disconnected. In other words, in the absence of reinforcement, the babybot ceased its behavior. Based on these findings, the authors argued that a sense of agency requires representing the causal link between one's actions and an effect, which is observed in infants but not in non-representational agents.

In a follow-up EEG study, Zaadnoordijk et al. [195] tested whether 3- to 4.5-month infants showed neural markers of causal action-effect models that are required for a sense of agency. Infants' limbs were connected to a digital mobile on a computer screen with four accelerometers attached to each limb, one of which was functional to activate the mobile. In the connect phase, the image was animated when the infant moved their limb connected to the functional accelerometer. In the disconnected phase, the link between infants' movement and the effect was broken, that is, the image remained static even if infants moved their limb operating the mobile. Data showed that a group of infants who showed increased error response in their brains in the disconnected phase (i.e., when the action-effect link was broken) also showed an extinction burst in their behavior indicating that they had constructed a causal model of their actions and the effect. Moreover, the same group of infants moved their limb that operated the mobile more frequently than the other connected limbs. These findings show that causal action-effect models that are necessary for a sense of agency only begin to emerge between 3- and 4.5-month of age in infancy. It is worth noting that the causal relation of actions and sensory effects can be represented as computational forward models that map the current state of the system to the next state through actions.

Other evidence regarding infants' ability to detect sensory contingencies presented by Verschoor and Hommel [182] also supports the idea that the sense of agency would emerge through the agent's own sensorimotor experience at around the same time, rather than being innate. As the authors point out, however, the ability to anticipate the outcomes of actions,

realized by a forward model, is vital but not sufficient for the complete development of sense of agency in infants. Without the ability to control their own bodies to render actions to change the environment corresponding to expected sensory effects, it is hard to rule out the possibility that the increase of infants' activities (during and after the experiments) might be due to the entrainment effect. It is worth recalling that the infants' motor movement is highly reflexive-like during this early stage of development, rather than voluntary and controlled (see discussion in Sect. 2.3). No earlier than 9 months old, infants know to select which actions to perform to achieve an expected or desired outcome, which relates to the action selection process (some sort of inverse model) (see [182] for a review; also [53, 191–193]). This timeline corresponds with the development of other motor skills in infants, e.g., reaching, as we will discuss in Sect. 2.3. These processes are, of course, in coordination with the maturation of other skills in infants such as eye-head coordination, and postural control (see e.g., [2, 185] for a review). However, the ability to predict sensory outcomes of motor actions develops earlier and precedes the ability to predict motor actions that would produce a desired sensory state (see [90] for a discussion and review on the development of predictive abilities in humans). The bidirectional associations between actions and effects being refined through the forward and inverse models are hypothesized as a trigger for the sense of agency: while the forward model helps to predict outcomes of conducted actions, the inverse model maps expected effects to action to perform. The smaller the error between the predicted and the actual outcome of intentional action—the predictive-coding process [6, 64, 65], the stronger the agency experience (see [182] for a review; [28, 85, 177]).

2.3 Development of the Peripersonal Space

While there is a body of studies on the representation of peripersonal space (PPS) in adults (see Sect. 3.3 for a brief review), there is very little research on this representation in infants, especially in their first months after birth. In a recent study, Orioli et al. [129] present a modified version of the reaction times (RTs) measurement, developed by [25], to address the question of whether the boundaries of the PPS representation is available in newborns. Instead of measuring the participants' vocal response time to tactile stimuli during an audiotactile interaction task, they propose to measure the saccadic latency to visual targets (sRTs) as an indirect measure of infants' RTs. With the results of infants' sRTs showing a similar pattern as the adults' RTs, Orioli et al. [129] suggest that some sort of PPS boundaries exist already soon after birth and thus facilitate the simultaneous multisensory matching in newborns.

More systematically, Bremner et al. [15, 16] propose that the development of PPS representation relates to two main

mechanisms, namely the visual-spatial reliance and postural remapping. The former mechanism, which develops as early as 6 months of age, allows infants to statistically estimate the body and surroundings based on the statistical variability of sensory sources, and the canonical layout of their body. This seems to follow the ability to detect sensory contingencies, which contributes to constructing some sort of perceptual body schema (as discussed in Sect. 2.1). However, these sensory contingencies, at the early age, may not necessarily be encoded in a certain body part reference frame, which is an important functionality of PPS representation [16]. The latter mechanism, the postural remapping, takes into account the postural changes to dynamically mapped external stimuli and limb position. This mechanism develops (and works alongside) in infants at around 6.5–10 months. In their experiments, Bremner et al. [18] reveal that 6.5-month-old infants bias their crossmodal responses to the typical side of their hands, whereas 10-month-old infants can respond appropriately in both sides even in crossed-hand postures. That said the findings suggest that PPS representation emerges through the combination of the two mechanisms and is not yet fully-developed prior to 6.5 months. This stage-wise development is in line with a recent neuroscience finding on somatosensory processing in 6–7-month-old infants (using somatosensory mismatch negativity (sMMN)), which speculates that the somatotopic phase of tactile processing does exist at that age while the later phases involving the frame of reference shifting are still under development [163].

As we present later (in Sect. 3.3) the sensorimotor mapping of PPS representation takes part in the voluntary movements to nearby objects within the reachable space. The development of these motor movements in infants can be observed as a source of behavioral measures for PPS development [15]. Furthermore, these changes in properties of the motor movements, in turn, provide sensory experiences for the refinement and the alignment of different sensory modalities, underlying the PPS representation. In the first year after birth, reaching movements develop from discontinuous, reflexive-like movements, to more direct, organized, and visually-elicited reaching (see [40] for a review; also [173]). In the former phase, the movement appears to be in a trial-and-error manner [173] and monitored mainly by proprioceptive feedback [18, 160]. That is, the movements to the goal can be conducted without visual feedback of the infant's hand (e.g., [31–33]). During this pre-reaching phase, infants are also observed to accidentally touch their own bodies—double touch [146], or clothes during spontaneous movements, giving rise to the grounding of bodily perception by integrating proprioception and touch. At the reaching onset, infants prefer looking at the space in which the hand and object make contact [39]. This suggests that tactile feedback facilitates the emergence of hand-eye coordination when the perception of the body and the external

space intersect and are being calibrated [40]. These events are in agreement with results from [17], arguing that this development of reaching behaviors is due to the infants' improvement in using both familiar and unfamiliar postural information (e.g., crossed-hands) to competently align spatial information from different sensory sources. These observations and results approximate the emergence of PPS in infants at around 6–10 months of age.

3 Computational and Robotic Models of Body Schema and PPS Representations

In this section, we first discuss the behavioral functionalities and properties of the body schema and PPS representations in humans (Sects. 3.1 and 3.3). Second, we review computational and robotics models of the representations (Sects. 3.2 and 3.4). This structure may encourage readers in comparing models of those sensory representations constructed in artificial agents with the ones in humans directly.

3.1 Properties and Function of the Body Schema Representation

As discussed above, the representation of the body schema seems to develop at a very early stage in newborns (in the continuity of the development during the fetal stage) and is based upon multisensory integration, i.e., from proprioceptive, tactile, and possibly visual information (see Table 1 and e.g., [27, 69, 84]). Along with the maturation of the visual modality, the body schema representation would be grounded and extended with the perceptual representation.

Due to the integration of sensory information, the body schema representation can plastically be modulated to include other objects such as a tool. This is known as the body schema extension paradigm, where agents are trained to actively use a tool to conduct motor actions [27, 112, 113, 161]. It is worth noting that this plasticity property does not exist when the tool is passively held by the agents. This dynamic plasticity of the body schema enables humans (and primates) to use tools flexibly.

The role of body schema in the agent's actions has been suggested as related to the motor control process through two types of internal models of the agent, namely the forward and inverse models. These two models construct the bi-directional mapping between the sensory information with motor information. Taking into account the temporal properties of sensory information forming the body representation, there exists a short-term representation, updated constantly like the angle of a joint, and a long-term representation, such as the size of a limb, which is relatively stable over time. Jointly, these two representations provide a good initial estimate for the body schema. This is required

for the inverse computation (of the inverse model) for motor commands generation to achieve a desired state of the body. Concurrently, the forward model predicts the outcomes of the motor commands, resulting in the predicted body schema, and receives the feedback from the sensory system as the updated body schema [45, 82].

Another key function of the body schema is to allow the coordinate transformations between different sensory modalities conducted by the brain. The transformations are thought to be processed under the population-based encoding conducted by gain field neurons (see [82] for a review; also [4, 9, 11, 22, 140, 156]). In robotics, the frame of reference (FoR) transformation is normally computed by the chain of transformation matrices, each represented by Denavit–Hartenberg (D–H) parameterization [164, 165]. However, the D–H transformations do not directly allow the mapping between different sensory modalities like the gain-field neurons.

3.2 Computational and Robotic Models of the Body Schema

The problem of learning the robot's body schema is often broken down into two main problems: (1) kinematics models identification/calibration, and (2) visuomotor learning/mapping, depending on the type of input signals. Models of the former group mostly require only body-related sensors including proprioception and touch, e.g., [49, 83, 103, 151, 200]. The latter group additionally requires visual information and takes advantage of the relation between the internal and external sensory modalities to construct the robot's body, e.g., [97, 125, 158, 178, 183, 189]. As a result, the former category requires some sort of a priori knowledge of the robot's body in terms of parameterized functions, e.g., CAD model, Forward kinematic, Inverse Kinematic, etc. The approaches of the latter category can work completely model-free and without a priori knowledge.

In the following, we present a survey of models on robotic body schema in ascending order of the amount of a priori knowledge provided in the learning problem. By organizing reviewed models in this order, we aim to emphasize one important aspect of autonomous systems: The ability to learn and adapt to dynamic environments. Ideally, an autonomous system should be able to learn to complete different tasks with only minimal provided information. A summary of the reviewed models is presented in Table 2.

Inspired by infants' self-touch behaviors for "body calibration", Roncone et al. [151] present a strategy for a humanoid robot to self-calibrate its body schema by bringing an end-effector of an arm to touch various locations in the other arm (which are covered by artificial skin taxels). In this work, the body schema is represented in the form of

Table 2 Summary of models of body schema representations

Model	Sensory information	Type of representation	Means of representation	Agent's body	Learning method
Ronccone et al. [151]	P & T	Body schema	Kinematics chain	iCub humanoid robot	Model-based and self-touch (single chain reformulation)
Li et al. [103]	P and T	Body schema	Kinematics chain	Two KUKA arms	Model-based and self-touch (sliding)
Vicente et al. [183, 184]	P and V	Body schema	Kinematics and particle filter	iCub humanoid robot	Model-based and online adaptation
Zenha et al. [200]	P and T	Body schema	Kinematics and extended Kalman Filter	iCub simulator	Model-based and goal babbling
Diaz Ledezma and Haddadin [49]	P	Body schema	FOPnet—variation of Newton–Euler equations	ATLAS simulator and Franka Emika	Model-based and constrained movements
Hoffmann et al. [83]	P and T	Tactile homunculus—body surface topology	MRF-SOM	iCub humanoid robot	Model free and multitouch (human stimulation)
Gama and Hoffmann [70]	P and T	Proprioception homunculus	MRF-SOM	NAO humanoid	Model free and self-touch
Abrossimoff et al. [1]	V and P	Body schema	Gain-field networks	A 3DoF simulated robot	Model free
Ulbrich et al. [178]	V and P	Body schema	Kinematic Bézier Maps	ARMAR-IIIa robot	Model free and Visual marker and motor babbling
Lallee and Dominey [97]	V and P	Body schema	MMCMs (SOM-based map)	iCub simulator	Model free and motor babbling
Schillaci et al. [158]	V and P	Body schema	DSOMs	NAO humanoid	Model free and motor babbling + Hebbian learning
Wijesinghe et al. [190]	V and P	Body schema	GASSOM and FC NN	iCub simulator	Motor babbling
Nguyen et al. [125]	V and P	Body schema	CNN and FC NN	iCub humanoid robot	Model free and arm and head babbling
Lanillos and Cheng [99]	V and T and P	Body schema	Predictive coding with Gaussian process regression	TOMM robot	Model free and limited arm babbling

Sensory information is coded as: vision—V, proprioception—P, touch—T

kinematic chains. Positions of the end-effector computed from proprioceptive input (i.e., joint encoders) and estimated from the skin system are utilized for kinematic calibration by an optimization algorithm.

Similarly, Li et al. [103] consider the problem of learning the body schema as kinematic calibration, in which they can exploit the CAD model for initialization. In detail, the authors utilize continuous self-touch movements (sliding) to calibrate the closed kinematic chain formed by both KUKA LWR arms (i.e., the slave and master in a torso setup) touching each other. Hence, the calibration problem becomes computing the relative transformation matrix by least squares estimation, given pairs of measured contact locations in the two arms.

Vicente et al. [183, 184] cast the internal process of adapting the robot body schema into a hand-eye coordination problem: first, the hand pose and initially calibrated offsets are estimated with the particle filter method, using stereo-vision and encoder measurements; then the internal model is updated by reducing differences between the model prediction of the end-effector and its observed value. For this approach, it is vital to have prior knowledge about the kinematic structure of the robot, i.e., a kinematics model, transformation matrices, and the camera's intrinsic parameters. In contrast to [183], Zenha et al. [200] employ an Extended Kalman filter instead of the Monte Carlo Particle filter for incremental kinematics model calibration in iCub simulation. Besides, tactile input caused by touch events between

the robot's finger and known surfaces during the robot's random movements is employed instead of visual input. The prior knowledge of the robot model is also employed in a goal babbling strategy toward the desired contact surfaces.

Diaz Ledezma and Haddadin [49] present a versatile and dedicated framework using the First-Order-Principle (FOP), derived from Newton-Euler equations, for learning both the body schema, i.e., topology and morphology, and the inverse dynamics, i.e., the inertial properties, of a simulated ATLAS humanoid and a Franka Emika arm in a modular manner. Parameters of FOP are learned from only the proprioceptive signals, including Kinematics-related measurements \mathcal{K} and dynamics-related measurements \mathcal{D} , collected during random trajectories generated by a PD controller. Especially, in this approach, the authors propose to exploit knowledge regarding the physical system, i.e., physical laws and joints connectivity, as optimization constraints in facilitating the topology search problem.

Differently, Hoffmann et al. [83] present an approach to construct the representation for the iCub robot's whole body skin surface in a form of a 2D map—a robotic somatosensory homunculus—by employing the dot product based SOM (DP-SOM) with an additional mask vector as a way to impose the binding constraint between neurons and input layer, i.e., skin taxels, to steer the learning process of the network. Finally, the authors show that the new variety of SOM—maximum receptive field SOM (MRF-SOM)—allows handling multiple tactile contacts simultaneously and enables the robot to learn a topological representation similar to the primary somatosensory cortex of primates. In a later study, Gama and Hoffmann [70] extend the MRF-SOM in the proprioceptive domain, to preliminary results. They aim to enable a robot to learn a proprioceptive representation of its joint space to resemble the proprioceptive representations in the somatosensory cortex. The underlying hypothesis is that the body representations may arise as a consequence of the agent's self-touch.

Inspired by the gain-field mechanism in human brains for the spatial transformation, Abrossimoff et al. [1] propose a neural network model consisting of two gain-field networks, the sigma-pi networks of radial basis function, for sensorimotor transformation and multimodal integration. The former is a visuomotor network for inverse dynamic learning, and the latter is to learn a body-centered coordinate system of the robot's hand and the target. After being trained, the networks enable a three-link robot to complete the reaching visual targets in a simulated 2D environment.

Ulbrich et al. [178] propose a method to learn the forward kinematics (FK) mapping from a robot's joint configuration and visual position of the end-effector as body schema learning. Moreover, they represent the FK with kinematic Bézier maps (KB-Maps), a derived technique from computational geometry, and show that the model can be learned

more efficiently with linear least square optimization by constraining the KB-Map with some topology knowledge. The learning method is validated on noisy data collected from random joint movements of the ARMAR-IIIa humanoid robot in both simulation and hardware.

Lallee and Dominey [97] propose so-called multi-model convergence maps (MMCMs)—a SOM-based implementation of the convergence–divergence zones framework—for multiple sensory modalities integration to encode sensorimotor experiences of iCub robots. MMCMs contain the bi-directional connections from each sensory modality, through a hierarchical structure (i.e., unimodal-amodal). Thus after being trained, it allows predicting the activation of missing modalities given the other(s). Herein, the visuomotor mapping³ is constructed by training the MMCMs with proprioceptive data from the arm and head, and image data from the robot camera during gazing and reaching activities. The encoded map of the learned internal representation allows the robot to “mentally imagine” the appearance and position of its body parts.

Schillaci et al. [158] learn a visuomotor coordination task in a Nao humanoid robot with a model consisting of two Dynamic Self-organising maps (DSOMs) encoding the arm and head joint space input, associated by Hebbian links to simulate synaptic plasticity of the brain. Two learning processes, one for updating DSOMs and another for Hebbian learning, are employed to train the model in an online manner during the robot's motor babbling. As a result, the robot improves its ability to gradually track the movement of its arm during the exploration process by controlling the head with output from the DSOMs based model.

Widmaier et al. [189] propose an algorithm based on Random Forest to estimate the robot's arm pose by regressing directly the joint angles from the depth input images on the pixel-level. The model operates in a frame-by-frame manner, without the requirement of an initialization or segmentation step. Instead of the random forest, Nguyen et al. [125]'s model utilizes a deep neural network to regress the joint angles of the iCub humanoid robot, given a pair of stereo-vision images and 6-DoF joint configuration of the robot's head (and eyes). The model is trained by a self-generated dataset from the robot's motor babbling of its head and arms in a simulated environment and the real robot. Furthermore, a framework based on a GAN network is also designed for transferring the learned visuomotor mapping from the simulation to a real robot, which helps to overcome calibration errors that often occur in physical robots.

Based on the hypothesis about the slow dynamics of the agent's own body compared to the dynamics of the environment, Laflaquiere and Hafner [95] propose a deep neural

³ Considered as PPS representation by the authors.

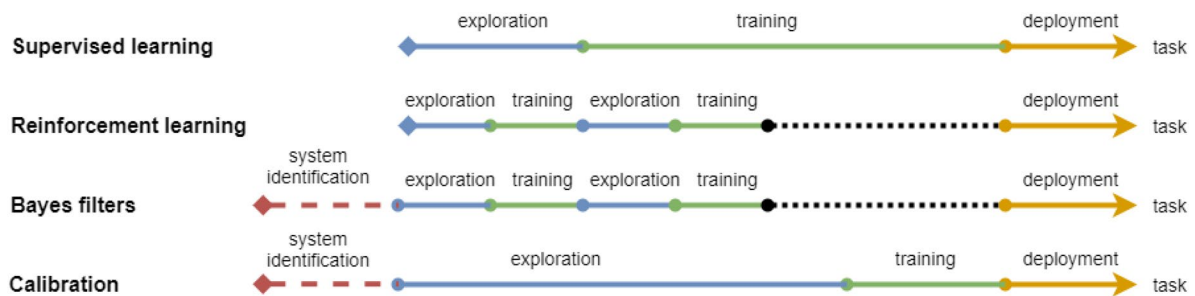


Fig. 3 Learning approaches employed for artificial agents and robots

network model for body representation estimation. The network is composed of two branches consisting of deconvolution and convolution layers. The former branch generates images of the robot's body with respect to the robot motor input, whereas the latter estimates the pixel-wise prediction error between the generated image from the former branch and the ground-truth. After training, the robot is able to predict the image of its own body in the environment, and to differentiate which part from the predicted image, i.e., a pixel, belongs to the agent's body or the environment based on its element-wise prediction error.

Wijesinghe et al. [190] present a bio-inspired predictive model for visuomotor mapping to track the robot's end-effector from the visual and proprioceptive inputs (i.e., from the position, velocity, and acceleration of four arm joints and position and velocity of two eye joints). The authors employ the generative adaptive subspace SOMs (GASSOMs) in their neural model for two purposes: (1) to encode the raw visual stimuli before combining with proprioception to generate a one-step prediction of the encoded visual stimuli; (2) to combine the encoded visual stimuli with its prediction. The output of the network is further used to control the robot's eye in tracking the arm movements.

Lanillos and Cheng [99], introduce a computational perceptual model based on the Gaussian additive noise model and free-energy minimization that enables a robot to learn, infer and update its body configuration from different sources of information, i.e., tactile, visual, and proprioceptive. The model is evaluated on a real multisensory robotic arm, showing the contributions of different sensory modalities in improving the body estimation, and the adaptability of the system against visuotactile perturbations.

So far, all models reviewed in this section share two common steps as shown in Fig. 3. The first step employs robots' movement as motor babbling for data generation. The second step constructs the relation between different sensory data by using analytical functions or machine learning techniques, e.g., artificial neural networks. While the performance of the analytically-based approaches depends mostly on the designers' choices of functions, the approaches using

machine learning techniques depend strongly on sensory data. Irrespective of the representation form employed as the body schema model, the main achievement of these approaches is the optimal estimation of the agents' body, i.e., joint configuration, end-effector position, or image of the hand/arm, with respect to the distribution of collected data from the babbling step. However, while these models demonstrate that they can (potentially) serve as a building block for more complex robotics behaviors, there are no possibilities for agents to continuously develop and learn these models outside the optimal estimation task they are meant to perform. We will discuss these points in detail in Sect. 5.

3.3 Peripersonal Space as a Brain's Representation of the Dynamic Interface Between the Body and the Environment

Similar to the body schema, the representation of the PPS representation is a result of various multisensory integration processes happening in the brain. The sources of sensory information include touch on the body, and vision and audio close to the body. Additionally, proprioception is also thought to take part in the process [161], especially in the arm-center PPS (see below text for more details of body-part centered PPS). This spatial representation helps to facilitate the manipulation of objects [72, 84] and to ease a variety of human actions such as reaching and locomotion with obstacle avoidance [84, 106]. Notably, this is not the case for space farther from the human body [56].

In terms of neuronal activation, the neuronal network of *parieto-premotor* areas of the cortex plays a vital role in PPS representation. PPS encoding neurons are found to be stimulated in several regions in primate brains, namely ventral intraparietal area (VIP), parietal area 7b, and premotor cortex (PMC), i.e., F4 and F5 areas (see [34] for a review). Neuroimaging studies in humans show similar results: Neurons in ventral PMC and inferior parietal sulcus (IPS⁴) relates to the hand-PPS; IPS neurons also relates to

⁴ IPS in human is homologous with VIP region in monkey.

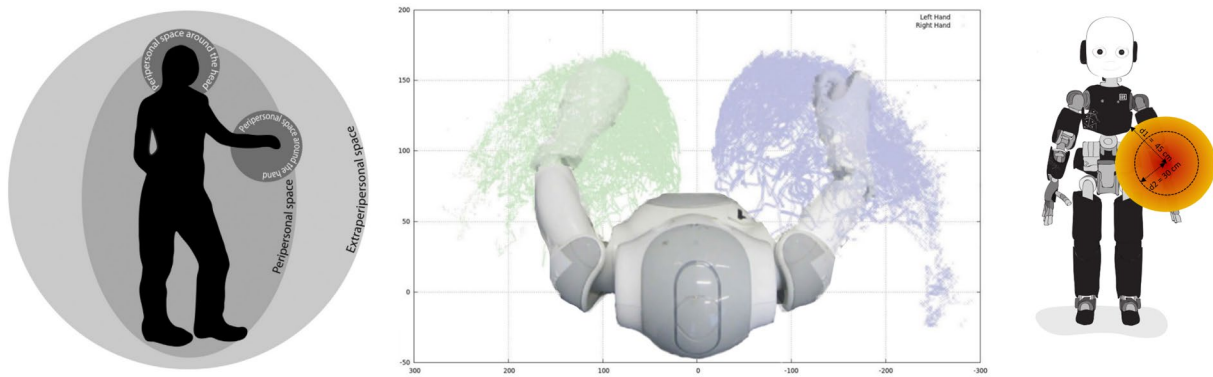


Fig. 4 PPS representations. Left: in humans (from [34]); center: active PPS as reachable regions in a Nao robot (from [159]); Right: defensive PPS as safety margin of a forearm in an iCub robot (from [123])

the face-PPS; many clusters of activation in parietal cortex and PMC correlate with PPS events (see [161] for a recent review; [76]). The activation of brain regions in the premotor cortex (during PPS events) also implies the link between the multisensory integration representation of PPS and motor activities.

The PPS representation serves as an interface between an agent's body and the environment through the multisensory neural network: it maps the sensory stimuli, e.g., objects via vision, directly to a body part frame of reference (FoR) to generate both voluntary and involuntary motor movements, e.g., reaching to grasp or avoidance reaction. The mapping is thought due to the multimodal receptive field (RF) of the activated PPS neurons anchored to this body part [61, 161]. Furthermore, the two types of PPS motor movements are not mutually exclusive [20, 21, 48, 161], and are thought to be related to two systems of PPS representation. First, the *active PPS* links with voluntary actions toward objects in the working reachable space. Second, the *defensive PPS* serves for involuntary defensive action [34, 46]. In the brains, there are specific networks for these two systems of PPS representation: The VIP-F4 network mainly process information for the defensive PPS [13, 14, 36, 37, 73–75]; the 7b-F5 network serves a core role of the active PPS [51, 60, 62, 114, 120, 144, 145] (see [34] for a review). Recent evidence from [132] suggests that the multisensory-motor mapping of PPS engages in both action plan and action execution processes.

The PPS representation is maintained (in the brain) by neurons with visuotactile RFs attached to different body parts, following the parts as they move (see e.g., [34, 84] for a recent survey). This forms a distributed and very dense coverage of the “safety margin” around the whole body. This defensive representation is not a unique space for the whole body but rather composed of many different sub-representations corresponding to different body parts. For example, the hands' PPS margin ranges around 30–45 cm from the

surface, the trunk 70–80 cm, and the face 50–60 cm [161] (see Fig. 4, left). Each sub-representation of a body part is closely coupled with that part even in movement, which is very useful for obstacle avoidance. When a body part moves, its PPS representation is modified independently from other body parts' representations, eliciting adaptive behaviors for only that specific body part.

That said, Cléry et al. [34] suggests that the separated PPS representations of body parts can interact and merge, depending on their relative positions. Besides, this protective safety zone is dynamically adapted to the action that the agent is performing, namely reaching vs. grasping [19]. It is also modulated by the state of the agent or by the identity and the “valence” (positive or negative) of the approaching object. For example, the safety zones are different in the cases of empty and full glasses of water [43], or in the cases of interacting with spiders and butterflies [42]. Furthermore, the social and emotional cues of interaction contexts also cause dynamic adjustment of the PPS representation [105, 172].

Moreover, the PPS representation is incrementally trained and adapted (i.e., expanded, shrunk, enhanced, etc.) through motor activities, as reported in, among others [34, 106, 162]. One of the motor actions being extensively studied is tool-use, where evidence from both primates and human studies reveal the enlargement of visuotactile RFs to include the tool [89, 110]⁵ or the increase of cross-modal extinction after actively using a tool to interact with far-space objects (see [112, 161] for reviews). Using short tools within the reachable space is not sufficient for this effect. More importantly, the degree of the extension of PPS representation

⁵ Though results from these works were originally interpreted as body schema extension, changes in the representation of the body itself after tool-use are not shown directly, but rather indirectly demonstrated through perceptual changes in PPS representations [26].

Table 3 Summary of models of PPS representation

Model	Sensory information	Type of representation	Means of representation	Agent's body	Learning method
De La Bourdonnaye et al. [44]	V & P	Body schema and PPS	CNN and autoencoder	Simulated stereo-camera and manipulator	Model free (reinforcement learning)
Pugach et al. [141]	V and T and P	Body schema and PPS	Gain-field network	Camera and manipulator	Model free and human touch
Nguyen et al. [123], Roncone et al. [150]	V and T and P	(defensive) PPS	Distributed visual RFs	iCub humanoid robot	Model-base and human touch/self-touch
Magosso et al. [108]	V and T	PPS	Unimodal and multimodal NN	No	No
Magosso et al. [107]	A and T	PPS	Unimodal and multimodal NN	No	Hebbian learning
Serino et al. [162]	A and T	PPS	Unimodal and multimodal NN	No	No
Straka and Hoffmann [168]	V and T	PPS	RBM and FC NN	No	Synthesized
Antonelli et al. [5] and Chinellato et al. [29]	V and P	PPS	RBF network	Tombato's humanoid robot	Model-free and gazing + reaching
Juett and Kuipers [92, 93]	V and P	PPS	PRM-like graph	Baxter robot	Model-free, motor babbling
Nguyen et al. [126]	V and T and P	PPS	CNN and FC NN	iCub simulator	Model free and arm babbling
Ramírez Contla [142]	V and P	PPS	FC NN	iCub simulator	Model-base vision and model-free robot's actions, body modification

Sensory information is coded as: vision—V, proprioception—P, touch—T, audio—A

depends on the way tools are used rather than the physical properties, e.g., the length, of the tools. In other words, bodily experiences are necessary for the plasticity of the PPS representation. The underlying reasons for this plasticity are temporally synchronous tactile and visual/audio stimulus during tool-use, which cause activation on the multisensory neurons integrated the corresponding unisensory tactile and visual/audio neurons. Thus these synapses between two sets of neurons are reinforced, according to the Hebbian learning principle.

The capabilities of PPS representation in updating the external stimuli to body parts (even in movements) imply the necessity of FoR transformations to align different sensory modalities coded in different FoRs. This is also the role of body schema (recall Sect. 3.1). However, in the PPS representation, the FoR transformations include both bodily and external stimuli (e.g., from vision, audio) [161]. To support this functionality, the proprioceptive stimuli may get involved with other sensory modalities, i.e., vision or audio,

especially in the case of the hand-centered PPS representations [161]. There is no clear evidence whether body schema representation takes part in the FoR transformation within the PPS representation. Cardinali et al. [27] suggest that the body schema may play as the “skeleton” for PPS but only it is not sufficient.

3.4 Computational and Robotic Models of PPS

Similar to Sect. 3.2, this section provides an overview of the research related to computational and robotics models of the PPS representation, organized in the increasing order of a priori information. The main differences between the approaches considered here are outlined in Table 3, which is constructed accounting for the following criteria: computation model for the PPS representation, sources of sensory information, agent's body, and learning approach (i.e., model-based or model-free, autonomous or not).

Roncone et al. [150, 152] propose a model of PPS representation as collision predictors distributed around an iCub robot's body, as a protective safety zone. The authors aim to investigate an integrated representation of the artificial visual and tactile sensors in the iCub humanoid robot. The multisensory information is integrated by probability associations between visual information, as the objects are seen approaching the body, and actual tactile information as the objects eventually physically contact the skin.

Nguyen et al. [123, 124] further extend this PPS model with the adaptability to the identity of approaching objects, e.g., neutral vs. dangerous, and interacting situation, e.g., hand-on interaction, to replicate the behavior of the protective PPS in humans [34, 42, 43]. Noticeably, the defensive behaviors of this PPS representation do not hinder the planned manipulating actions such as reaching, grasping an object. Instead, these two capabilities work harmoniously within the cognitive architectures through an optimal control algorithm. Hence the model facilitates the robot's activities alongside human partners in different human–robot interaction scenarios [119, 124].

Magosso et al. [108] propose and analyze a neural network model to integrate visuotactile stimuli for the PPS representation. This model is composed of two identical networks, corresponding to the left and right hemispheres of the brain. Each network is composed of unimodal neurons for visual and tactile stimuli input, and multimodal neurons for multisensory integration. Inhibitory connections also exist between the left and right hemisphere networks to model their mutually inhibiting relations: When one hemisphere activates, the other one will be to an equal extent inhibited. This brain-like construction allows modeling the behavior of the PPS at the physical level and to be compared with data collected from humans. Similar models are proposed for the case of audiotactile stimuli in Magosso et al. [107], Serino et al. [162]. The authors did not design a training procedure, except for the tool-use case presented in [107], where the Hebbian learning rule is employed.

Similarly, the PPS representation by Straka and Hoffmann [168]'s computational model associates visual and tactile stimuli in a simulated 2D scenario. The model is composed of Restricted Boltzmann Machine for object properties association (i.e. position and velocity), and a two-layer fully-connected artificial neural network for “temporal” prediction. After training, the model is capable of predicting the collision position, given the visual stimulus as in [150]. The designed scenario remains quite simple, however, since it boils down to simply a simulation in 2D space: The skin area is a line and there is no concept of the body, hence no transformation between sensory frames is taken into account.

Differently, Juett and Kuipers [92, 93] model the PPS representation as a graph of nodes in the robot's reachable space through a constrained motor babbling of a Baxter robot.

Each node in the graph is composed of inputs from joint encoder values and images). With the learned graph, search algorithms can be applied to find the shortest path connecting the current and the final state. In their most recent work, the final state search algorithm is extended to allow grasping objects. Although the graph model can be learned without a kinematics model, the authors utilize some image segmentation techniques to locate the robot's gripper during the learning phase, and the targets in the action phase from the input image(s). Requiring each node in the graph to store images is a memory-intensive solution.

Antonelli et al. [5] and Chinellato et al. [29] adopt radial basis function networks to construct the forward and inverse mappings between stereo visual data and proprioceptive data in a robot platform. This is conducted through the robot's gazing and reaching activities within the reachable space. Their mapping, however, requires visual markers to extract features with known disparity. Although authors aim to form a model of PPS representation, without the involvement of external objects and tactile sensing, there is not much difference between this model and visuomotor mapping models of the body schema.

Inspired by [5, 108], Nguyen et al. [126] present a model of the spatial representation by a visuo-tactile-proprioceptive integration neural network for reaching external object in reachable space on iCub robots. The model maps the visual input from 6-DoF stereo-vision system to the 10-DoF motor space including the torso and an arm. This is taken place under the supervision signal of touch events between objects and artificial skin taxels covering the robot body. After training, this model allows the robot to estimate the ability of reaching/colliding with visual stimulus within its reachable space, as similar as PPS representation.

De La Bourdonnaye et al. [44] present a stage-wise approach for a robotics agent learning to touch an object in the scene with a reinforcement learning algorithm. First, the robot learns to fixate the object by learning the configuration of the camera system to encode the object. Then it learns hand-eye coordination by constructing the mapping from the robot's motor space to the camera space. Finally, the previously learned information is used to shape the reward in learning to touch objects. While the first learning stage is equivalent to learning the PPS representation, the second phase is learning the body schema of the agent.

Pugach et al. [141] implement a gain-field network (recall [1] in Sect. 3.2) to construct the representations of a Jaco arm's body schema and PPS. Inputs for the network come from a fixed camera, a system of artificial skin covering the robot's forearm and its encoder, collected during one-degree-of-freedom movement of the arm. The tactile signal is employed to trigger the process of learning visual representation—the visual-tactile receptive field. Though

the approach requires some preprocessing steps, i.e., color-based object recognition, constraint movement of the robot, and denoised filters for outputs of the gain-field network, it presents some potential aspects of a defensive PPS representation as work by [150].

On the other hand, Ramírez Contla [142] focuses on the plastic nature of PPS representation to account for the modification the body undergoes, and the impact of this plasticity on the confidence levels in respect to reaching activities. In their experiments, the author first assesses the contribution of visual and proprioceptive data to reaching performance, then measures the contribution of posture and arm-modification to reaching regions. The modifications applied to the arm, i.e., the changes in the arm's length, have similar effects as the extension of the PPS representation during tool-use.

As we discussed earlier, the main difference between models of PPS representation reviewed in this section and body-schema models is the involvement of external objects in the vicinity of the agents' body and thus the tactile sensing. Unsurprisingly, most approaches to modeling PPS representation also apply similar steps as the body schema models: (1) generating sensory data through the agent's movement for; (2) learning the model of PPS representation. The PPS representations are mostly constructed by artificial neural networks. The approaches are able to fulfill the main function of the PPS representation: correlating information from different sensory modalities including FoR transformations; and mapping the external objects within reach onto the agents' body parts. However, they also lack the ability to learn continuously outside the context of the designed learning tasks, as with the cases of body schema models.

4 The Active Self

4.1 The Self in Humans

The process of infants' development involves, among other things, the acquisition of "body knowledge". The body knowledge has been described within the context of infants' development as the formation of the body's sensorimotor map (the body schema) and the variety of actions that support motor and cognitive development [109]. The formation of the body schema—the sensorimotor representation of the body, begins with the genetic predisposition for the organization of body parts representation in the S1 and the M1. It is later elaborated through early (fetal stage) body-environment involuntary interactions such as the touch of the amniotic fluid with the skin (part of the development of tactile perception), and most importantly, body–body interactions (e.g., self-touch). In the first months of life, the infant is more focused on body-body

interactions, for example, acquiring body knowledge through self-touch behaviors. This goes alongside motor development, and as the body is the most accessible part of the environment, and also the most predictable, the body is the first part of the environment to be modeled [167]. At this time, the agent is learning the forward model—the causal relationship between motor actions and their sensory effects on the body. However, motor actions do not necessarily have to be voluntary, intentional, or goal-directed in order to construct the forward model and develop the causal representation of action-effect links. The *bidirectional* associations between actions and effects will develop with an inverse model that is involved with goal-directed movements—selecting actions that produce a predicted or desired sensory effect. This stage can be thought of as one that incorporates *verification* [167].

According to the basic principles of developmental robotics [167], artificial agents and robots need to be able to verify what they learn about the environment [169], in order to effectively interact with a complex and dynamic external environment. Verification requires the ability to act upon the environment, hence, the agent needs to be embodied [167]. In addition, the verification needs "grounding"—a process or its outcome that establishes what is valid verification. Because the environment is probabilistic, grounding requires the agent to construct action-effect links, and therefore to have a causal representation of actions and effects in a probabilistic manner. The process of grounding the verification in a probabilistic way requires the agent to repeat its actions to test and refine what it learns about the environment as a causal representation of action-effect, through, for example, detecting temporal contingencies [167]. In this view, it arises then that the developmental process advances from exploring the most predictable and verifiable parts of the environment (i.e., the body) to the least. Exploration is driven intrinsically: the agent is "drawn" to explore that part of the environment which has intermediate variability, until the variability is reduced, and the attention shifts to other parts (see also [159]).

The recent term "body know-how" focuses more on the practical aspects of body knowledge [91], and was defined as "the ability to sense and use the body parts in an organized and differentiated manner" [91, p. 109]. Body know-how and its acquisition are therefore interlinked with motor development. The more body know-how is accumulated, motor skills enhanced, and the forward model perfected, the more the agent can learn about its environment. This is because more body know-how leads to more informative and complex interactions. These are "informative" in the sense that the verification becomes more and more efficient as the agent learns about the morphological properties of the body, and about how to move the body. One can argue that the sort

Table 4 Summary of active self models based on multisensory sources

Model	Sensory information	Body representation	Means of representation	Active self ability	Learning method
Zambelli and Demiris [199]	V and P and T and A and M	Implicit	Ensembles of algorithms	Agency	Forward model learning by imitation
Copete et al. [38]	V and P and T	Implicit	Deep autoencoder	Agency	Imitation
Hwang et al. [87, 88]	V and P	Implicit	P-VMDNN	Agency	Forward model learning by imitation
Zambelli et al. [198]	V and P and T and A and M	Implicit	MVAE	Agency	Forward model learning by imitation
Saponaro et al. [157]	V and P	Explicit	Partical filter and PCA + Bayesian network	Agency	Affordance learning
Lang et al. [98]	V v P	Explicit	CNN	Agency and self-other distinction	Supervised learning
Lanillos et al. [100]	V and P and T	Explicit	Hierarchical Bayesian model	Body ownership—self-detection	Bayesian inference
Hinz et al. [80]	V and P and T	Explicit	Predictive coding with Gaussian Process regression	Body ownership—Rubber hand illusion	Limited arm babbling and Bayesian inference
Lanillos et al. [101]	V and P	Explicit	Mixed density networks	Agency	Limited arm babbling and Bayesian inference
			CNN	Body ownership—self-other distinction	Classifier

Sensory information is coded as: vision—V, proprioception—P, motor—M, touch—T, audio—A

of information that the agent learns from the interaction with the environment is statistical information: Spatiotemporal, sensorimotor contingencies, as well as causal links between actions and effects. Because the world is not deterministic, this information is therefore probabilistic.

Developing a representation of causal links between actions and effects on the environment is necessary, but not sufficient for the development of the sense of agency. This is because having a representation of associations between actions and effects is not informative with regards to who the author of the action was. In order to verify that the author of an action having led to an effect was oneself, the agent needs to perform goal-directed actions. In computational terms, the forward model represents the causal links between actions and effects and allows the agent to predict the sensory outcomes of actions. The agent makes use of the predictions brought by the forward model to produce goal-directed actions. The agent also needs to perform goal-directed actions to refine the inverse model, a representation of the links between a sensory effect and the action that will cause it, i.e., bidirectional action-effect links. Verschoor and Hommel [182] argue that goal-directed action is a prerequisite for the emergence of the minimal self, rather than an indication of its emergence.

Moreover, the developmental process is iterative: Acquiring knowledge about the body (“this movement led to this body sensation”—what body sensation does a certain movement elicit?) leads to acquiring knowledge about the environment (“this movement led to this perceived effect on the environment”—what is the perception that comes from this movement?), which leads back to knowledge about the body (“to get effect x on the environment, I need to move this way”—how to move to achieve a certain goal). The interface through which body know-how is acquired is the body schema representation, and the interface through which complex knowledge about the environment is acquired is the PPS representation.

The notion of verification reflects the active inference approach [67], which postulates that to reduce uncertainty (free energy), an embodied agent uses an internal generative model that samples sensory data through action. Sampling is done through approximate Bayesian inference to induce posterior beliefs, under the assumption that active sampling will update model priors. The uncertainty is resolved with actions that hold “epistemic value” to the agent, i.e., information-seeking behaviors [67]. The principles of active inference and free energy present the forward model as a mechanism to fulfill curiosity by minimizing the expected prediction error [66].

Table 5 Summary of active self models based on a single sensory input

Model	Sensory information	Body representation	Means of representation	Active self ability	Learning method
Watter et al. [188]	V	Implicit	VAE	Agency	Unsupervised representation learning
Van Hoof et al. [180]	T	Implicit	VAE	Agency	Reinforcement learning
Byravan et al. [23]	V (3D point cloud)	Explicit	SE3-POSE-NETS	Agency	Unsupervised representation learning
Agrawal et al. [3]	V	Implicit	CNN	Agency	Forward model learning by supervised learning
Pathak et al. [134]	V	Implicit	CNN	Agency	Reinforcement learning with intrinsic reward
Park et al. [131]	P	Implicit	RNNBP	Agency	Kinesthetic teaching and Imitation

Sensory information is coded as: vision—V, proprioception—P, touch—T

In this probabilistic framework, the agent gathers information about statistical regularities, through predictive processes—making predictions about sensory outcomes of generated actions, and resolving “prediction errors”—either in favor of updating the model or in favor of adopting the sensory information itself (see [104] for a review on the minimal self in this framework). One might think about the body model as explicitly distinct from the “world model”. However, the boundary between the body and the environment can also be thought of as a sort of statistically-dependent boundary: The body is the most predictable and consistent part of the environment, and therefore the most verifiable [167].

Lending this notion to the minimal self, the boundary between the (sensorimotor, minimal) self model and non-self model can also be thought of as statistically-dependent. For example, the notion of nested Markov blankets [94] postulates that biological systems tend to autonomously self-organize in a coherent way, through active inference, to separate their internal states from external ones, with nested hierarchical Markov blankets that define its boundaries in a statistical sense. Similarly, Hafner et al. [78] propose the notion of the self-manifold for an artificial agent, which is defined as a dynamic and adaptive outline for the boundaries of the self, and related to both body ownership and agency, as, in their view, they cannot be separated. They propose to formalize the self-manifold as a Markov blanket around the sensorimotor states of an agent.

4.2 Robotic Models of the Active Self

In this section, we review robotics models of the active self or models owning a common feature, which employs the predictive coding mechanism or the forward model. This focus stems from the idea that the feeling of agency can emerge in an agent with an ability to anticipate the effect of its own

action (see detailed discussions in Sects. 2.2 and 4.1). We first review models employing multisensory modalities (in Sect. 4.2.1 and Table 4), then continue with using a single sensory modality (mostly from visual input, in Sect. 4.2.2 and Table 5). For the latter cases, they allow to capture the dynamics of the whole system (including the agent and the interactive environment) via only a single input due to the special design of the input, i.e., the visual input is not taken from the first perspective viewer (as in humans and other animals).

4.2.1 Models with Multisensory Input

Zambelli and Demiris [199] introduce a learning architecture where forward and inverse models are coupled and updated as new data becomes available, without prior information about the robot’s kinematic structure. The ensemble learning process of the forward model combines different parametric and non-parametric online algorithms to build the sensorimotor representation models, while the inverse models are learned by interacting with a piano keyboard, thus engaging vision, touch, motor encoders, and sound. Zambelli et al. [198] extended this idea but trained a multimodal variational autoencoder (MVAE) model from motor babbling data that included combinations of complete and missing data from the joint position, vision, touch, sound, and motor command modalities. They tested the model in the same imitation task that involved predicting the sensory state of the robot arising from visual input alone when observing another agent’s actions.

The computational model by Copete et al. [38] allows a simulated robot to (1) acquire the ability to predict the intention of others’ actions, and (2) learn to produce the same actions. The main component of the model is a deep autoencoder-based predictor, whose aim is to integrate

visual, motor, and tactile signals (spatially and temporally). In the action learning mode, the autoencoder receives input from all sensory modalities to train the network, while in the action observation mode (of the other robot), the learned network receives only visual signals as input and is able to produce the missing sensory modalities, i.e., tactile and joint signals. Feeding the output signals back into the input of the network allows it to predict the future sensorimotor signals.

Hwang et al. [87] construct a multilayer predictive model (P-VMDNN) with two pathways for visual and proprioceptive inputs, in which pathways are only connected in the top-layer to simulate the link between perception and action. These pathways employ variations of RNN, namely predictive-multiple spatio-temporal scales and multiple timescales for processing visual and proprioceptive input, respectively. The model is trained end-to-end by backpropagation through time (BPTT) to minimize the (one-step ahead) prediction errors of the two inputs. As a result, a simulated iCub can imitate some primitive hand-waving gestures of another displayed on a screen, even in the case of missing one of the sensory inputs (similar to models using an autoencoder). Recently, this model was also employed for imitative interaction between an iCub robot and a human [88].

Saponaro et al. [157] further exploit the body schema and forward model (developed from visual and proprioceptive information by Vicente et al. [184]) in “mental” simulation of sensory outcomes in an affordance learning task. This is carried out by employing Principal Component Analysis (PCA) and an additional Bayesian Network to construct the relation between four pre-defined actions (in various directions) of robots with the known hand configurations or objects/tools.

Lang et al. [98] employ a deep convolutional neural network that integrates proprioception, vision, and motor commands to predict the visual outcomes of a Nao robot’s actions. This forward model was trained with self-generated data from the robot’s motor babbling and was employed in a self-other distinction task. It is expected that the prediction error of the forward model is lower when observed arm movements are performed by the robot itself than by other agents. The authors also showed how predictions can be used to attenuate self-generated movements, and thus create enhanced visual perceptions, where the sight of objects—originally occluded by the robot body—was still maintained.

Lanillos et al. [100] conceive a hierarchical Bayesian model, which aims to integrate movement and touch from an artificial skin system with vision from a camera. The hierarchical model consists of three layers: the first two deal with self-detection using inter-modal contingencies to avoid relying on visual assumptions like markers, whereas the last layer employs self-detection to enable conceptual interpretation such as object discovery. To validate the model, the authors design an experiment entailing object

discovery through interactions, in which the robot has to discern between its own body, usable objects, and illusion in the scene.

Hinz et al. [80] extend the model of body estimation by Lanillos et al. [99] (see discussion in Sect. 3.2) with an additional visual-tactile sensation, in the task of replicating the Rubber Hand Illusion in a humanoid robot. In this experiment, the authors consider the differences between the robot’s estimated end-effector position and the ground truth as the drift of the illusion, which shows similar patterns with the experiment in human participants.

Instead of the Gaussian process regression in previous models [99], Lanillos et al. [101] employ the mixture density network (MDN) to encode the visual generative model and follow the free energy minimization framework to estimate the robot’s body. The authors further utilize a deep learning-based classifier for contingency learning, i.e., the probability of association between the visual input from optical flow and the joint velocity of the robot. Finally, both the prediction error of the robot’s body estimation and the sensory contingency contribute to the tasks of self-recognition and self/other distinction at a sensorimotor level.

4.2.2 Models with Single Sensory Input

Watter et al. [188] employ a Variational Autoencoder (VAE) to probabilistically infer the visual depiction of the system state into a latent space, where the dynamic transition from the current latent state to the next state (under the untransformed action) is assumed to be linear. As a result, the problem of non-linear system identification and control from high-dimensional images becomes locally optimal control in linearized latent space. The learned feature allows locally optimal actions to be found in closed form stochastic optimal control algorithms. An additional constraint is also employed to enforce the similarity between samples from the state transition distribution and the inference distribution, thus guaranteeing a valid encoded representation for long-term prediction. Both autoencoder and transition networks are learned jointly.

Similarly, Van Hoof et al. [180] propose a variation of a VAE to encode low-dimensional features of the raw tactile input for more efficient reinforcement learning. The VAE is modified to take into account the transition dynamics by linearly combining the estimated latent state with action (through a linear neural network layer), and generating a prediction of the next latent state. The feature is learned by optimizing the marginal likelihood of sensory input concerning the prediction of the next latent state (instead of the latent state).

Borrowing some ideas from [188], Byravan et al. [23] develop a deep learning based predictive model to learn the

latent space from a pair of successive input images related by an action. The predictive model is formed as a U-net with an encoder of convolutional layers and a decoder of de-convolutional layers. Specifically, the network can (1) model the structure of the scene x_t in the form of segmented moving parts $k \in K$ (predefined) and their 6D pose; (2) predict the changes of each part k under the applied action; and (3) output the prediction of the scene dynamics, i.e., a predicted point cloud, as a result of the rigid rotation and translation of all point x^j that belong to the part k . The model is trained by the joint prediction losses at the point cloud and the pose level. After training, the model is employed for closed-loop control directly in latent space with a reactive controller using gradient-based methods.

Agrawal et al. [3] propose a method to jointly learn the forward model (for action-outcome prediction) and inverse model (for a greedy planner to generate the robot's discretized poking action) from the feature space of visual input in a supervised manner. The authors show that the forward model helps to regularize the inverse model and generalizes better than the case using only the inverse model (especially when the robot is tasked to poke the object in a long-distance).

Park et al. [131] deploy a computational model based on an RNNPB—recurrent neural network with parametric bias (PB)—on robots (i.e., a virtual 2 DoF arm and an NAO humanoid) and gradually allow them to imitate the movement shape of goal-directed motor behaviors. In order to do so, the network is trained by BPTT with the prediction error between the network output and the reference during the learning phase. During the imitation phase, with observed actions, the PB is first recognized by BPTT and then can be used to generate imitated actions as the output of the network.

Pathak et al. [134] propose to use an ensemble of forward dynamics functions within a policy-gradient-based deep reinforcement learning agent. The model also exploits the disagreement among prediction errors in the ensemble as the intrinsic motivation to drive the agent's exploration without external reward from the environment. Furthermore, the authors formulate the intrinsic reward as a differentiable function to perform policy optimization in a supervised learning manner instead of reinforcement. The authors show that a robotics manipulator can learn to touch a random object in the scene with only visual input.

5 Discussion

5.1 From Biological Agents to Artificial Agents

In humans, the senses of body ownership and agency develop through interaction with the environment, which is perceived and controlled with the available sensorimotor system. The underlying mechanisms are the associations

between different sensory modalities and sensorimotor contingencies. This learning process leads to the formation of representations of the body and the surrounding environment within reach, including other objects and agents.

Most of the research on learning multisensory representations that we review in Sects. 3.2 and 3.4 casts the development of multisensory representations in bio-agents into equivalent robotics learning tasks, namely *body calibration*, *pose estimation*, and *visuomotor mapping* for the body schema representation; or *reaching estimation* and *collision estimation* for the PPS representation (refer to Fig. 3 for different learning approaches). Tackling the development problems in this way and following two-step approaches, most approaches are able to find the optimal solution for the designed learning tasks and provide the learning outcome as a building block in a more complex architecture for robotics behaviors. This is, however, different from the development of sensorimotor representations in biology, which is a continuous iterative and interactive process. For example, the body schema representation in humans not only adapts during the motor babbling phase in infants but also continues to adapt during the tool-use context, where the agent's intention is to optimize the actions of grasping and manipulating the tool rather than optimizing the estimation of the position of the hand and arm. In other words, human sensorimotor representations develop in multiple settings: they are not only learned once through random actions and serve as input for more complex actions, but these representations are also continuously refined through feedback from the perceived outcomes of complex actions.

Similarly, models of the active self presented in Sect. 4.2.1 focus on learning to optimize the prediction loss of the forward models concerning the raw sensory input from multiple sources directly—without constructing explicit representations of the body and the environment. The prediction errors of the learned forward models are then employed to generate movements as similar as learned ones through imitation or babbling. By additionally constructing the explicit sensory representation of the agents' body—in forms of generative images or joint estimation—other models like [80, 98, 101] enable agents to distinguish between the agent's body and external objects. However, all of the existing approaches lack the ability to generalize beyond the learning tasks.

The predictive models with single sensory input that we review in Sect. 4.2.2 lack certain properties of bio-agents that are related to multisensory integration. However, their proposed architectures can efficiently enable agents to develop the ability to predict outcomes of their own actions in a latent representational space. In these models, the latent state abstraction serves as dimensionality reduction for the desired learning tasks. However, all existing models learn these two steps separately instead of simultaneously [133].

Table 6 Summary of sub-problems focused by reviewed models

	Biological agents	Artificial agents	
Multisensory representation	Body schema	Calibration Pose estimation Visuomotor mapping	Sect. 3.2
	Peripersonal space	Reaching estimation Collision estimation	Sect. 3.4
Minimal self sensation	Agency	Forward and inverse model	Sect. 4.2
	Body ownership		

In humans, the involvement of the body schema and PPS representations in various motor activities (as we review in Sects. 3.1, and 3.3) suggests that the brain might learn and use these representations as a process of dimensionality reduction or state abstraction, which then facilitate the ability to learn manipulation skills and transferring knowledge between different learned skills. Furthermore, the sense of touch plays a crucial role in the development of PPS and body schema representations, especially in the later development of manipulation skills when interacting with the external environment. Results from models taking into account the tactile sensing capability as one of the sensory modalities, e.g., [80, 99, 124, 126, 141, 150, 151] present similar behaviors to humans such as in the cases of self-touch and body-object interaction. Thus it is worth considering this sensory modality in the architecture for developmental agents.

5.2 A Conceptual Sketch for the Development of an Artificial Minimal Self

Our review on the state of the art in models of the active self and bodily-related representations suggests certain guidelines and principles that are important for modeling a self computationally. Here we propose a sketch of architecture to integrate these principles (see Fig. 5), aiming to enable artificial agents to develop the active self through self-exploration within an environment as discussed by Schillaci et al. [159].

Our review points out that agents require two critical components to develop a self: (1) a representation of multimodal sensorimotor contingencies, and (2) bidirectional associations of actions and effects. The former condition is addressed in our proposal with the *Multisensory integration* module. The latter condition is fulfilled by two modules, namely the *Predictor* and the *Action generator*.

The *Predictor* is a multimodal forward model that predicts a sensory effect $\hat{\phi}(s_{t+1})$ from a currently conducted action a_t and the currently perceived sensory state representation $\phi(s_t)$. The *Action generator* generates motor actions a_t under constraints exerted by the environment and under consideration of the prediction error e_{t+1} of the *Predictor*. Both the *Predictor* and the *Action generator* operate in the latent space of the multimodal sensory input, which is compressed by the *Multisensory*

integration process. We specify the operation of these modules as follows:

Multisensory representations:

$$\phi(s_t) = \phi^{PPS}(s_t^{eU_i}) \cup \phi^{body}(s_t^i)$$

Predictor:

$$\hat{\phi}(s_{t+1}) = f(\phi(s_t), a_t) \tag{1}$$

Predictor error:

$$e_{t+1} = \frac{1}{2} \left\| \hat{\phi}(s_{t+1}) - \phi(s_{t+1}) \right\|_2^2$$

Here, $\phi^{PPS}(s_t^{eU_i})$ denotes the representation of the PPS, s_t denotes the current sensory state and $\phi^{body}(s_t^i)$ denotes the body schema representation. In terms of the implementation, all these modules can be constructed by a multiple head neural network with each head corresponding to each module output. The large part of the network is shared between different modules. This artificial neural architecture reflects the hierarchical structure of multisensory integration processes to generate abstract, multimodal predictions at the high level from low-level unimodal sensory signals [64].

Importantly, all modules learn simultaneously through the agent’s own interactive experience in the environment. Their behavior is driven by sparse extrinsic feedback and the intrinsic motivation to minimize prediction errors of their intentional actions. In this setting, learning to minimize the prediction errors and integrate multisensory input are the auxiliary tasks alongside the main task of learning to generate skill-dependent actions. One possibility to model the *Action generator* is to combine motor babbling as was used by most of the reviewed approaches and sampled outputs of the reinforcement learning policy, which is known as *ϵ -greedy* exploration [170, Chapter 13]. Taking an example of a reinforcement learning agent, at every time step, the agent selects an action drawn from the policy π —an action generator—based on the current state s_t , exerts it on the environment and receives an extrinsic reward r_t^e depending on the next state s_{t+1} of the whole system. Moreover, the predictor also provides another internal reward r_t^i based on the prediction error. In turn, the total reward $r_t = r_t^e + r_t^i$ guides the

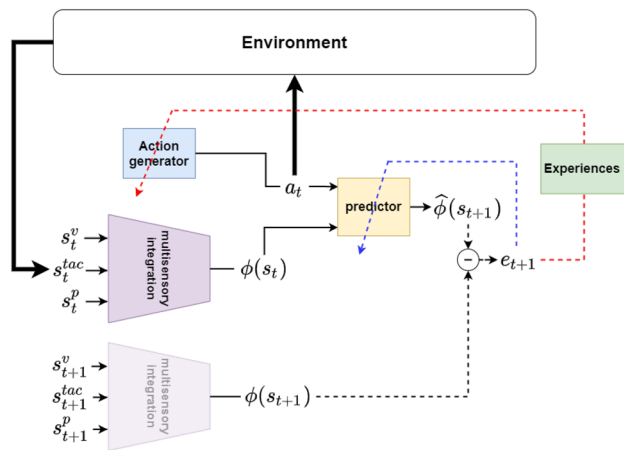


Fig. 5 Proposed model for developing the active self in artificial agents. s_t^v , s_t^{tac} , s_t^p denote raw visual, tactile, and proprioception input at time t respectively. The blue and red arrows denote the source of data affecting the learning of the target module: blue for the *predictor* and red for the *Action generator*.

improvement of the policy π through established algorithms such as policy gradient [170].

One problem, however, is that agents are prone to overfitting when learning only from a single task or in a single environment. As we point out in Sects. 3.2 and 3.4, irrespective of the chosen form for the models of the sensory representations, behaviors of trained agents are optimized with respect to the estimation task they are desired to perform. They lack the ability to learn these models continuously outside the context of the tasks. For example, an agent who is trained to perform a visuomotor tracking skill cannot easily adapt to completing the grasping skill without catastrophically forgetting the trained knowledge. To address this issue, we propose to use the sub-problems in the third column of Table 6 (i.e., calibration, pose estimation, visuomotor mapping, reaching estimation, and collision estimation) as benchmark tests instead of using them as objective functions for the learning task (e.g., object manipulation, tool use). Our main hypothesis is that since embodied agents are equipped with various sensory modalities such as vision, touch, and proprioception, the developed agents should pass the benchmark tests and show behaviors equivalent to humans, including sensory phenomena like the Rubber Hand Illusion. The general learning objective function is designed to maximize the agent's ability to learn skills while minimizing the prediction error of the agent's internal predictor. Furthermore, we propose to employ stage-wise or curriculum learning strategies for a set of different skills⁶, which are

⁶ Agents do not know about this whole set of skills to learn and choose from (e.g. as in [35])

gradually more difficult to achieve [130]. Since the sensory representations continuously mature during learning of one skill, e.g., object manipulation, the development implicitly facilitates transfer learning to other more sophisticated skills, e.g., grasping a tool and using a tool to manipulate objects, faster and easier than learning from scratch. During the learning process, while the skill-dependent objective function motivates the agent to generate actions to fulfill the skill requirement, the auxiliary objective function ensures multisensory representation learning to minimize the prediction error e_{t+1} (Eq. 1). The former learns with the stored long-term experiences, whereas the latter is trained with the short-term prediction error (as shown in right side of Fig. 5). The auxiliary task of learning multisensory representation plays as an intrinsic motivation for the transition from learning one skill to another skill.

The multitask learning process of the proposed architecture includes learning the multisensory representations and learning the predictive model for control tasks. This learning process is equivalent to state representation learning for control, as highlighted in a recent review by Lesort et al. [102]. Furthermore, our architecture shares some similarities with the proposal by Nagai [122], who focuses on modeling cognitive development by minimizing prediction errors of a forward model. However, we emphasize the importance of learning the sensory representations as a state abstraction from multiple sources simultaneously with learning the internal models in our proposal. In summary, we propose to combine several strategies to support the ability of continual learning, as highlighted in Parisi et al. [130], namely, multisensory learning and intrinsic motivation (of minimizing prediction error). This combination is supported by reviewed evidence from the development of biological agents and related computational and robotics models.

5.3 Towards Modeling a Self with Higher Cognitive Functions

The embodied conceptualization hypothesis by Lakoff and Johnson [96] entails that our body-specific sensorimotor apparatus and, therefore, our representations of body schema and PPS, determine how we conceptualize the world. Hence, these representations have strong influences on higher cognitive functions as they directly shape the way we think [137]. This becomes evident in natural language, where metaphorical expression involves basic body-related concepts [54, 175]. What remains open, though, is how we can model grounding of sensorimotor concepts computationally. Several approaches, including the theory of event-coding [86], and event segmentation theory [77, 197], exist. However, it is subject to future work to fully integrate these approaches within a unifying computational theory of high-level cognition. Research on the minimal active self fosters

the development of such a unifying theory as it allows one to investigate how basic body-related concepts emerge from sensorimotor interaction.

Acknowledgements EK was supported by the Max Planck Society and by a grant from the Deutsche Forschungsgemeinschaft (DFG) (grant numbers: KA 4926/1-1). YKG and VVH were funded by the DFG - “Prerequisites for the Development of an Artificial Self” (402790442). PDHN, ME, SW acknowledge funding by the DFG through the IDEAS (402776968) and LeCAREbot (433323019) projects.

Funding Open Access funding enabled and organized by Projekt DEAL.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article’s Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article’s Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

- Abrossimoff J, Pitti A, Gaussier P (2018) Visual learning for reaching and body-schema with gain-field networks. In: 2018 Joint IEEE 8th International Conference on Development and Learning and Epigenetic Robotics, ICDL-EpiRob 2018. IEEE, pp 197–203
- Adolph KE, Joh, AS (2007) Motor development: How infants get into the act. Introduction to infant development, 2nd edn., pp 63–80
- Agrawal P, Nair A, Abbeel P et al (2016) Learning to poke by poking: Experiential learning of intuitive physics. *Adv Neural Inf Process Syst*:5074–5082
- Ajemian R, Bullock D, Grossberg S (2001) A model of movement coordinates in the motor cortex: Posture-dependent changes in the gain and direction of single cell tuning curves. *Cereb Cortex* 11(12):1124–1135
- Antonelli M, Chinellato E, Del Pobil AP (2013) On-line learning of the visuomotor transformations on a humanoid robot. *Intelligent autonomous systems*, vol 12. Springer, New York, pp 853–861
- Apps MA, Tsakiris M (2014) The free-energy self: a predictive coding account of self-recognition. *Neurosci Biobehav Rev* 41:85–97
- Bahrick LE, Lickliter R (2002) Intersensory redundancy guides early perceptual and cognitive development. *Adv Child Dev Behav* 30:153–187
- Bahrick LE, Watson JS (1985) Detection of Intermodal Proprioceptive-Visual Contingency as a Potential Basis of Self-Perception in Infancy. *Dev Psychol* 21(6):963–973
- Baraduc P, Guigon E, Burnod Y (2001) Recoding arm position to learn visuomotor transformations. *Cereb Cortex* 11(10):906–917
- Blanke O, Metzinger T (2009) Full-body illusions and minimal phenomenal selfhood. *Trends Cogn Sci* 13(1):7–13
- Blohm G, Crawford JD (2009) Fields of gain in the Brain. *Neuron* 64(5):598–600
- Bradley RM, Mistretta CM (1975) Fetal sensory receptors. *Physiol Rev* 55(3):352–382
- Bremmer F, Duhamel JR, Ben Hamed S, Graf W (2002a) Heading encoding in the macaque ventral intraparietal area (VIP). *Eur J Neurosci* 16(8):1554–1568
- Bremmer F, Klam F, Duhamel JR et al (2002b) Visual-vestibular interactive responses in the macaque ventral intraparietal area (VIP). *Eur J Neurosci* 16(8):1569–1586
- Bremner AJ, Holmes NP, Spence C (2008a) Infants lost in (peripersonal) space? *Trends Cogn Sci* 12(8):298–305
- Bremner AJ, Holmes NP, Spence C (2012a) The development of multisensory representations of the body and of the space around the body. In: *Multisensory development*. Oxford University Press, Oxford
- Bremner AJ, Lewkowicz DJ, Spence C (2012b) *Multisensory development*. Oxford University Press, Oxford
- Bremner AJ, Mareschal D, Lloyd-Fox S, Spence C (2008b) Spatial localization of touch in the first year of life: early influence of a visual spatial code and the development of remapping across changes in limb position. *J Exp Psychol Gen* 137(1):149–162
- Brozzoli C, Cardinali L, Pavani F, Farné A (2010) Action-specific remapping of peripersonal space. *Neuropsychologia* 48(3):796–802
- Brozzoli C, Makin TR, Cardinali L (2011) Peripersonal space: a multisensory interface for body-object interactions. In: Murray MM, Wallace MT (eds) *The neural bases of multisensory processes*. Number, et al (May 2014) in *Frontiers in Neuroscience*, chapter 23. Taylor & Francis, London, pp 449–466
- Brozzoli C, Pavani F, Urquizar C et al (2009) Grasping actions remap peripersonal space. *NeuroReport* 20(10):913–917
- Bullock D, Grossberg S, Guenther FH (1993) A self-organizing neural model of motor equivalent reaching and tool use by a multijoint arm. *J Cogn Neurosci* 5(4):408–435
- Byravan A, Lceb F, Meier F, Fox D (2018) SE3-Pose-nets: structured deep dynamics models for visuomotor control. In: *IEEE International Conference on Robotics and Automation*, pp 3339–3346
- Cangelosi A, Schlesinger M (2014) *Developmental robotics: from babies to robots*. The MIT Press, Cambridge
- Canzoneri E, Magosso E, Serino A (2012) Dynamic sounds capture the boundaries of peripersonal space representation in humans. *PLoS One* 7(9):e44306
- Canzoneri E, Ubaldi S, Rastelli V et al (2013) Tool-use reshapes the boundaries of body and peripersonal space representations. *Exp Brain Res* 228(1):25–42
- Cardinali L, Brozzoli C, Farné A (2009) Peripersonal space and body schema: two labels for the same concept? *Brain Topogr* 21(3–4):252–260
- Chambon V, Haggard P (2013) Premotor or ideomotor: how does the experience of action come about? *Action Sci*:358–380
- Chinellato E, Antonelli M, Grzyb BJ, del Pobil AP (2011) Implicit sensorimotor mapping of the peripersonal space by gazing and reaching. *IEEE Trans Auton Ment Dev* 3(1):43–53
- Chugani HT (1994) Development of regional brain glucose metabolism in relation to behavior and plasticity. In: *Human behavior and the developing brain..* The Guilford Press, pp 153–175
- Clifton RK, Muir DW, Ashmead D, Clarkson M (1993) Is visually guided reaching in early infancy a myth?? *Child Dev* 64(4):1099–1110
- Clifton RK, Rochat P, Litovsky RY, Perris EE (1991) Object representation guides infants’ reaching in the dark. *J Exp Psychol Human Percept Perform* 17(2):323–329

33. Clifton RK, Rochat P, Robin DJ, Bertheir NE (1994) Multimodal perception in the control of infant reaching. *J Exp Psychol Human Percept Perform* 20(4):876–886
34. Cléry J, Guipponi O, Wardak C, Ben Hamed S (2015) Neuronal bases of peripersonal and extrapersonal spaces, their plasticity and their dynamics: Knowns and unknowns. *Neuropsychologia* 70:313–326
35. Colas C, Fournier P, Sigaud O et al (2019) CURIOUS: intrinsically motivated modular multi-goal reinforcement learning. In: *International conference on machine learning (ICML)*
36. Cooke DF, Graziano MS (2003) Defensive movements evoked by air puff in Monkeys. *J Neurophysiol* 90(5):3317–3329
37. Cooke DF, Graziano MS (2004) Sensorimotor Integration in the Precentral Gyrus: polysensory neurons and defensive movements. *J Neurophysiol* 91(4):1648–1660
38. Copete JL, Nagai Y, Asada M (2017) Motor development facilitates the prediction of others' actions through sensorimotor predictive learning. In: *ICDL-EpiRob 2016*. IEEE, pp 223–229
39. Corbetta D, Thurman SL, Wiener RF et al (2014) Mapping the feel of the arm with the sight of the object: on the embodied origins of infant reaching. *Front Psychol* 5
40. Corbetta D, Wiener RF, Thurman SL, McMahon E (2018) The embodied origins of infant reaching: implications for the emergence of eye-hand coordination. *Kinesiol Rev* 7(1):10–17
41. Dall'Orso S, Steinweg J, Allievi A et al (2018) Somatotopic mapping of the developing sensorimotor cortex in the preterm human brain. *Cereb Cortex* 28(7):2507–2515
42. de Haan AM, Smit M, Van der Stigchel S, Dijkerman HC (2016) Approaching threat modulates visuotactile interactions in peripersonal space. *Exp Brain Res* 234(7):1875–1884
43. de Haan AM, Van der Stigchel S, Nijens C, Dijkerman HC (2014) The influence of object identity on obstacle avoidance reaching behaviour. *Acta Psychol* 150:94–99
44. De La Bourdonnaye F, Teuliere C, Triesch J, Chateau T (2018) Learning to touch objects through stage-wise deep reinforcement learning. In: *IEEE International Conference on Intelligent Robots and Systems*, pp 7789–7794
45. de Vignemont F (2010) Body schema and body image-Pros and cons. *Neuropsychologia* 48(3):669–680
46. de Vignemont F, Iannetti GD (2015) How many peripersonal spaces? *Neuropsychologia* 70:327–334
47. Di Noto PM, Newman L, Wall S, Einstein G (2013) The hermunculus: what is known about the representation of the female body in the brain? *Cereb Cortex* 23(5):1005–1013
48. di Pellegrino G, Ládavas E (2015) Peripersonal space in the brain. *Neuropsychologia* 66:126–133
49. Diaz Ledezma F, Haddadin S (2019) FOP networks for learning humanoid body schema and dynamics. In: *IEEE-RAS International Conference on Humanoid Robots*, pp 1121–1127
50. Dijkerman HC, de Haan EHF (2007) Somatosensory processes subserving perception and action. *Behav Brain Sci* 30(2):189–201
51. Durand JB, Nelissen K, Joly O et al (2007) Anterior regions of monkey parietal cortex process visual 3D shape. *Neuron* 55(3):493–505
52. D'Angelo M, di Pellegrino G, Seriani S et al (2018) The sense of agency shapes body schema and peripersonal space. *Sci Rep* 8(1):1–11
53. Elsner B (2007) Infants' imitation of goal-directed actions: The role of movements and action effects. *Acta Psychol* 124(1):44–59
54. Eppe M, Trott S, Raghuram V et al (2016) Application-independent and integration-friendly natural language understanding. In: *Global Conference on Artificial Intelligence (GCAI)*, pp 340–352
55. Fagard J, Esseily R, Jacquey L et al (2018) Fetal origin of sensorimotor behavior. *Front Neurobotics* 12
56. Farné A, Dematté ML, Ládavas E (2005) Neuropsychological evidence of modular organization of the near peripersonal space. *Neurology* 65(11):1754–1758
57. Filippetti ML, Johnson MH, Lloyd-Fox S et al (2013) Body perception in newborns. *Curr Biol* 23(23):2413–2416
58. Filippetti ML, Lloyd-Fox S, Longo MR et al (2015a) Neural mechanisms of body awareness in infants. *Cereb Cortex* 25(10):3779–3787
59. Filippetti ML, Orioli G, Johnson MH, Farroni T (2015b) Newborn body perception: sensitivity to spatial congruency. *Infancy* 20(4):455–465
60. Fogassi L, Ferrari PF, Gesierich B et al (2005) Neuroscience: parietal lobe: from action organization to intention understanding. *Science* 308(5722):662–667
61. Fogassi L, Gallese V, Fadiga L et al (1996) Coding of peripersonal space in inferior premotor cortex (area F4). *J Neurophysiol* 76(1):141–157
62. Fogassi L, Luppino G (2005) Motor functions of the parietal lobe. *Curr Opin Neurobiol* 15(6):626–631
63. Fotopoulou A, Tsakiris M (2017) Mentalizing homeostasis: the social origins of interoceptive inference. *Neuropsychanalysis* 19(1):3–28
64. Friston K (2012) Prediction, perception and agency. *Int J Psychophysiol* 83(2):248–252
65. Friston K, Kiebel S (2009) Predictive coding under the free-energy principle. *Philos Trans R Soc B Biol Sci* 364(1521):1211–1221
66. Friston K, Mattout J, Kilner J (2011) Action understanding and active inference. *Biol Cybern* 104(1–2):137–160
67. Friston K, Rigoli F, Ognibene D et al (2015) Active inference and epistemic value. *Cogn Neurosci* 6(4):187–214
68. Gallagher S (2000) Philosophical conceptions of the self: implications for cognitive science. *Trends Cogn Sci* 4(1):14–21
69. Gallese V, Sinigaglia C (2010) The bodily self as power for action. *Neuropsychologia* 48(3):746–755
70. Gama F, Hoffmann M (2019) The hermunculus for proprioception: toward learning the representation of a humanoid robot's joint space using self-organizing maps. [arXiv:1909.02295](https://arxiv.org/abs/1909.02295)
71. Georgie YK, Schillaci G, Hafner VV (2019) An interdisciplinary overview of developmental indices and behavioral measures of the minimal self. In: *2019 Joint IEEE 9th international conference on development and learning and epigenetic robotics, ICDL-EpiRob 2019*, pp 129–136
72. Goerick C, Wersing H, Mikhailova I, Dunn M (2005) Peripersonal space and object recognition for humanoids. In: *Humanoid robots, 2005 5th IEEE-RAS International Conference on*. IEEE, pp 387–392
73. Graziano MS, Cooke DF (2006) Parieto-frontal interactions, personal space, and defensive behavior. *Neuropsychologia* 44(13):2621–2635
74. Graziano MS, Taylor CS, Moore T (2002) Complex movements evoked by microstimulation of precentral cortex. *Neuron* 34(5):841–851
75. Graziano MSA, Hu XINT, Gross CG et al (1997) Visuospatial properties of ventral premotor cortex. *J Neurophysiol* 77(5):2268–2292
76. Grivaz P, Blanke O, Serino A (2017) Common and distinct brain regions processing multisensory bodily signals for peripersonal space and body ownership. *NeuroImage* 147:602–618
77. Gumbsch C, Butz MV, Martius G (2019) Autonomous identification and goal-directed invocation of event-predictive behavioral primitives. *IEEE Trans Cogn Dev Syst* ([online](https://doi.org/10.1109/TCDS.2019.2918448))
78. Hafner VV, Loviken P, Villalpando AP, Schillaci G (2020) Prerequisites for an artificial self. *Front Neurobotics* 14
79. Head H, Holmes G (1911) Sensory disturbances from cerebral lesions. *Brain* 34(2–3):102–254

80. Hinz NA, Lanillos P, Mueller H, Cheng G (2018) Drifting perceptual patterns suggest prediction errors fusion rather than hypothesis selection: Replicating the rubber-hand illusion on a robot. In: ICDL-EpiRob 2018. IEEE, pp 125–132
81. Hoffmann M (2017) The role of self-touch experience in the formation of the self. [arXiv:1712.07843](https://arxiv.org/abs/1712.07843)
82. Hoffmann M, Marques H, Arieta A et al (2010) Body schema in robotics: a review. *IEEE Trans Auton Mental Dev* 2(4):304–324
83. Hoffmann M, Straka Z, Farkaš I et al (2018) Robotic homunculus: learning of artificial skin representation in a humanoid robot motivated by primary somatosensory cortex. *IEEE Trans Cogn Dev Syst* 10(2):163–176
84. Holmes NP, Spence C (2004) The body schema and multi-sensory representation(s) of peripersonal space. *Cogn Process* 5(2):94–105
85. Hommel B (2015a) Action control and the sense of agency. Oxford University Press, New York, In *The Sense of Agency*
86. Hommel B (2015b) The theory of event coding (TEC) as embodied-cognition framework. *Front Psychol* 6:1318
87. Hwang J, Kim J, Ahmadi A et al (2018) Predictive coding-based deep dynamic neural network for visuomotor learning. In: 7th Joint IEEE International Conference on Development and Learning and on Epigenetic Robotics, ICDL-EpiRob 2017, vol 2018, pp 132–139
88. Hwang J, Kim J, Ahmadi A et al (2020) Dealing with large-scale spatio-temporal patterns in imitative interaction between a robot and a human by using the predictive coding framework. *IEEE Trans Syst Man Cybern Syst* 50(5):1918–1931
89. Iriki A, Tanaka M, Iwamura Y (1996) Coding of modified body schema during tool use by macaque postcentral neurones. *Neuroreport* 7(14):2325–2330
90. Jacquey L, Baldassarre G, Santucci VG, O'Regan JK (2019) Sensorimotor contingencies as a key drive of development: from babies to robots. *Front Neurobotics* 13
91. Jacquey L, Popescu ST, Vergne J et al (2020) Development of body knowledge as measured by arm differentiation in infants: From global to local? *Br J Dev Psychol* 38(1):108–124
92. Juett J, Kuipers B (2016) Learning to reach by building a representation of peri-personal space. In: 2016 IEEE-RAS 16th International Conference on Humanoid Robots (Humanoids). IEEE, pp 1141–1148
93. Juett J, Kuipers B (2018) Learning to grasp by extending the peripersonal space graph. In: 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems. IEEE, pp 8695–8700
94. Kirchhoff M, Parr T, Palacios E et al (2018) The Markov blankets of life: autonomy, active inference and the free energy principle. *J R Soc Interface* 15(138):20170792
95. Laflaquiere A, Hafner VV (2019) Self-supervised body image acquisition using a deep neural network for sensorimotor prediction. In: 2019 Joint IEEE 9th International Conference on Development and Learning and Epigenetic Robotics, (ICDL-EpiRob), pp 117–122
96. Lakoff G, Johnson M (1999) *Philosophy in the flesh*. Basic Books
97. Lalle S, Dominey PF (2013) Multi-modal convergence maps: from body schema and self-representation to mental imagery. *Adapt Behav* 21(4):274–285
98. Lang C, Schillaci G, Hafner VV (2018) A deep convolutional neural network model for sense of agency and object permanence in robots. In: 2018 Joint IEEE 8th International Conference on Development and Learning and Epigenetic Robotics (ICDL-EpiRob). IEEE, pp 257–262
99. Lanillos P, Cheng G (2018) Adaptive robot body learning and estimation through predictive coding. In: IEEE International Conference on Intelligent Robots and Systems, pp 4083–4090
100. Lanillos P, Dean-Leon E, Cheng G (2017) Yielding self-perception in robots through sensorimotor contingencies. *IEEE Trans Cogn Dev Syst* 9(2):100–112
101. Lanillos P, Pages J, Cheng G (2020) Robot self/other distinction: active inference meets neural networks learning in a mirror. In: 24th European Conference on Artificial Intelligence (ECAI 2020)
102. Lesort T, Díaz-Rodríguez N, Goudou JF, Filliat D (2018) State representation learning for control: an overview. *Neural Netw* 108:379–392
103. Li Q, Haschke R, Ritter H (2015) Towards body schema learning using training data acquired by continuous self-touch. In: IEEE-RAS International Conference on Humanoid Robots (Humanoids). IEEE, pp 1109–1114
104. Limanowski J, Blankenburg F (2013) Minimal self-models and the free energy principle. *Front Human Neurosci* 7:547
105. Lourenco SF, Longo MR, Pathman T (2011) Near space and its relation to claustrophobic fear. *Cognition* 119(3):448–453
106. Lådavas E, Serino A (2008) Action-dependent plasticity in peripersonal space representations. *Cogn Neuropsychol* 25(7–8):1099–1113
107. Magosso E, Ursino M, di Pellegrino G et al (2010a) Neural bases of peri-hand space plasticity through tool-use: Insights from a combined computational-experimental approach. *Neuropsychologia* 48(3):812–830
108. Magosso E, Zavaglia M, Serino A et al (2010b) Visuotactile representation of peripersonal space: a neural network study. *Neural Comput* 22(1):190–243
109. Mannella F, Santucci VG, Somogyi E et al (2018) Know your body through intrinsic goals. *Front Neurobotics* 12:30
110. Maravita A, Iriki A (2004) Tools for the body (schema). *Trends in Cogn Sci* 8(2):79–86
111. Marshall PJ, Meltzoff AN (2015) Body maps in the infant brain
112. Martel M, Cardinali L, Roy AC, Farné A (2016) Tool-use: an open window into body representation and its plasticity. *Cogn Neuropsychol* 33(1–2):82–101
113. Martin B, Wittmann M, Franck N et al (2014) Temporal structure of consciousness and minimal self in schizophrenia. *Front Psychol* 5:1175
114. Matelli M, Luppino G (2001) Parietofrontal circuits for action and space perception in the macaque monkey. *NeuroImage* 14(1 II):27–32
115. Meltzoff AN (2007) Like me: a foundation for social cognition. *Dev Sci* 10(1):126–134
116. Meltzoff AN, Marshall PJ (2020) Importance of body representations in social-cognitive development: new insights from infant brain science, 1st edn, vol 254. Elsevier B.V, Amsterdam
117. Meltzoff AN, Saby JN, Marshall PJ (2019) Neural representations of the body in 60-day-old human infants. *Dev Sci* 22(1):1–8
118. Morgan R, Rochat P (1997) Intermodal Calibration of the body in early infancy. *Ecol Psychol* 9(1):1–23
119. Moulin-Frier C, Fischer T, Petit M et al (2018) DAC-h3: a proactive robot cognitive architecture to acquire and express knowledge about the world and the self. *IEEE Trans Cogn Dev Syst* 10(4):1005–1022
120. Murata A, Fadiga L, Fogassi L et al (1997) Object representation in the ventral premotor cortex (Area F5) of the monkey. *J Neurophysiol* 78(4):2226–2230
121. Myowa-Yamakoshi M, Takeshita H (2006) Do human fetuses anticipate self-oriented actions? a study by four-dimensional (4d) ultrasonography. *Infancy* 10(3):289–301
122. Nagai Y (2019) Predictive learning: its key role in early cognitive development. *Philos Trans R Soc B Biol Sci* 374(1771)
123. Nguyen D, Hoffmann M, Roncone A et al (2018a) Compact real-time avoidance on a humanoid robot for human-robot interaction.

- In: ACM/IEEE International Conference on Human-Robot Interaction
124. Nguyen PD, Bottarel F, Pattacini U et al (2018b) Merging physical and social interaction for effective human-robot collaboration. In: 2018 IEEE-RAS 18th International Conference on Humanoid Robots (Humanoids), pp 1–9
 125. Nguyen PD, Fischer T, Chang HJ et al (2018c) Transferring visuomotor learning from simulation to the real world for robotics manipulation tasks. In: 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). IEEE, pp 6667–6674
 126. Nguyen PD, Hoffmann M, Pattacini U, Metta G (2019) Reaching development through visuo-proprioceptive-tactile integration on a humanoid robot—a deep learning approach. In: 2019 Joint IEEE 9th International Conference on Development and Learning and Epigenetic Robotics (ICDL-EpiRob). IEEE, pp 163–170
 127. OpenStax College (2020) Central processing
 128. O'Regan JK, Noë A (2001) A sensorimotor account of vision and visual consciousness. *Behav Brain Sci* 24(5):939–973
 129. Orioli G, Santoni A, Dragovic D, Farroni T (2019) Identifying peripersonal space boundaries in newborns. *Sci Rep* 9(1):1–11
 130. Parisi GI, Kemker R, Part JL et al (2019) Continual lifelong learning with neural networks: a review. *Neural Netw* 113(5):54–71
 131. Park JC, Kim DS, Nagai Y (2018) Learning for goal-directed actions using RNNPB: developmental change of “What to Imitate”. *IEEE Trans Cogn Dev Syst* 10(3):545–556
 132. Patané I, Cardinali L, Salemme R et al (2018) Action planning modulates peripersonal space. *J Cogn Neurosci* 31(8):1141–1154
 133. Pathak D, Agrawal P, Efros AA, Darrell T (2017) Curiosity-driven exploration by self-supervised prediction. In: International Conference on Machine Learning (ICML)
 134. Pathak D, Gandhi D, Gupta A (2019) Self-supervised exploration via disagreement. In: International Conference on Machine Learning (ICML), pp 5062–5071
 135. Penfield W, Boldrey E (1937) Somatic motor and sensory representation in the cerebral cortex of man as studied by electrical stimulation. *Brain* 60(4):389–443
 136. Penfield W, Rasmussen T (1950) The cerebral cortex of man; a clinical study of localization of function. Macmillan, New York
 137. Pfeifer R, Bongard JC (2006) How the body shapes the way we think: a new view of intelligence (Bradford Books). The MIT Press, Cambridge
 138. Pfeifer R, Lungarella M, Iida F (2007) Self-organization, embodiment, and biologically inspired robotics. *Science* 318(5853):1088–1093
 139. Pitti A, Kuniyoshi Y, Quoy M, Gaussier P (2013) Modeling the minimal newborn's intersubjective mind: the visuotopic-somatotopic alignment hypothesis in the superior colliculus. *PLoS One* 8(7):e69474
 140. Pouget A, Deneve S, Duhamel JR (2002) A computational perspective on the neural basis of multisensory spatial representations. *Nat Rev Neurosci* 3(9):741–747
 141. Pugach G, Pitti A, Tolochko O, Gaussier P (2019) Brain-inspired coding of robot body schema through visuo-motor integration of touched events. *Front Neurobotics* 13:5
 142. Ramírez Contla S (2014) Peripersonal space in the humanoid robot iCub. Ph.D. thesis, University of Plymouth
 143. Reissland N, Francis B, Aydin E et al (2014) The development of anticipation in the fetus: a longitudinal account of human fetal mouth movements in reaction to and anticipation of touch. *Dev Psychobiol* 56(5):955–963
 144. Rizzolatti G, Luppino G (2001) The cortical motor system. *Neuron* 31(6):889–901
 145. Rizzolatti G, Matelli M (2003) Two different streams form the dorsal visual system: anatomy and functions. *Exp Brain Res* 153(2):146–157
 146. Rochat P (1998) Self-perception and action in infancy. *Exp Brain Res* 123(1–2):102–109
 147. Rochat P, Morgan R (1995a) Spatial determinants in the perception of self-produced leg movements by 3- to 5-month-old infants. *Dev Psychol* 31(4):626–636
 148. Rochat P, Morgan R (1995b) Spatial determinants in the perception of self-produced leg movements in 3-to 5-month-old infants. *Dev Psychol* 31(4):626
 149. Rochat P, Striano T (1999) Social-cognitive development in the first year. Understanding others in the first months of life, early social cognition, pp 3–34
 150. Roncone A, Hoffmann M, Pattacini U et al (2016) Peripersonal space and margin of safety around the body: learning visuo-tactile associations in a humanoid robot with artificial skin. *PLoS One* 11(10):e0163713
 151. Roncone A, Hoffmann M, Pattacini U, Metta G (2014) Automatic kinematic chain calibration using artificial skin: self-touch in the icub humanoid robot. In: Robotics and automation (ICRA), 2014 IEEE International Conference on. IEEE, pp 2305–2312
 152. Roncone A, Hoffmann M, Pattacini U, Metta G (2015) Learning peripersonal space representation through artificial skin for avoidance and reaching with whole body surface. In: Intelligent robots and systems (IROS), 2015 IEEE/RSJ International Conference on. IEEE, pp 3366–3373
 153. Rovee CK, Rovee DT (1969) Conjugate reinforcement of infant exploratory behavior. *J Exp Child Psychol* 8(1):33–39
 154. Rovee-Collier CK, Morrongiello BA, Aron M, Kupersmidt J (1978) Topographical response differentiation and reversal in 3-month-old infants. *Infant Behav Dev* 1:323–333
 155. Saby JN, Meltzoff AN, Marshall PJ (2015) Neural body maps in human infants: Somatotopic responses to tactile stimulation in 7-month-olds. *NeuroImage* 118:74–78
 156. Salinas E, Abbott LF (2001) Coordinate transformations in the visual system: How to generate gain fields and what to compute with them. *Prog Brain Res* 130:175–190
 157. Saponaro G, Vicente P, Dehban A et al (2018) Learning at the ends: from hand to tool affordances in humanoid robots. In: 7th Joint IEEE International Conference on Development and Learning and on Epigenetic Robotics, ICDL-EpiRob 2017, pp 321–337
 158. Schillaci G, Hafner VV, Lara B (2014) Online learning of visuomotor coordination in a humanoid robot. A biologically inspired model. In: ICDL-EPIROB 2014. IEEE, pp 130–136
 159. Schillaci G, Hafner VV, Lara B (2016) Exploration behaviors, body representations, and simulation processes for the development of cognition in artificial agents. *Front Robot AI* 3:1–18
 160. Schlesinger M, Parisi D (2001) Multimodal control of reaching - Simulating the role of tactile feedback. *IEEE Trans Evol Comput* 5(2):122–128
 161. Serino A (2019) Peripersonal space (PPS) as a multisensory interface between the individual and the environment, defining the space of the self. *Neurosci Biobehav Rev* 99(August 2018):138–159
 162. Serino A, Canzoneri E, Marzolla M et al (2015) Extending peripersonal space representation without tool-use: evidence from a combined behavioral-computational approach. *Front Behav Neurosci* 9
 163. Shen G, Weiss SM, Meltzoff AN, Marshall PJ (2018) The somatosensory mismatch negativity as a window into body representations in infancy. *Int J Psychophysiol* 134(October):144–150
 164. Siciliano B, Khatib O (eds) (2016) Springer handbook of robotics. Springer International Publishing, Cham

165. Siciliano B, Sciavicco L, Villani L, Oriolo G (2009) Robotics: modelling, planning and control, advanced textbooks in control and signal processing. Springer, OCLC, New York
166. Stein BE, Meredith MA (1993) The merging of the senses. The MIT Press, Cambridge
167. Stoytchev A (2009) Some basic principles of developmental robotics. *IEEE Trans Auton Mental Dev* 1(2):122–130
168. Straka Z, Hoffmann M (2017) Learning a peripersonal space representation as a visuo-tactile prediction task. International conference on artificial neural networks. Springer, New York, pp 101–109
169. Sutton RS (2001) Verification, the key to AI. In: On-line essay (online). <http://www.cs.ualberta.ca/sutton/IncIdeas/KeytoAI.html>
170. Sutton RS, Barto AG (2017) Reinforcement learning: an introduction. MIT Press, Cambridge
171. Tani J, White J (2020) Cognitive neurorobotics and self in the shared world, a focused review of ongoing research. *Adapt Behav*
172. Teneggi C, Canzoneri E, di Pellegrino G, Serino A (2013) Social modulation of peripersonal space boundaries. *Curr Biol* 23(5):406–411
173. Thelen E, Corbetta D, Kamm K et al (1993) The transition to reaching: mapping intention and intrinsic dynamics. *Child Dev* 64(4):1058–1098
174. Thomas BL, Karl JM, Whishaw IQ (2015) Independent development of the reach and the grasp in spontaneous self-touching by human infants in the first 6 months. *Front Psychol* 5:1526
175. Trott S, Eppe M, Feldman J (2016) Recognizing intention from natural language: clarification dialog and construction grammar. In: Workshop on Communicating Intentions in Human-Robot Interaction @ IEEE International Symposium on Human and Robot Interactive Communication
176. Tsakiris M (2017) The multisensory basis of the self: from body to identity to others. *Q J Exp Psychol* 70(4):597–609
177. Tsakiris M, Hesse MD, Boy C et al (2007) Neural signatures of body ownership: a sensory network for bodily self-consciousness. *Cereb Cortex* 17(10):2235–2244
178. Ulbrich S, De Angulot VR, Asfour T et al (2009) Rapid learning of humanoid body schemas with kinematic bézier maps. In: 9th IEEE-RAS International Conference on Humanoid Robots (Humanoids), pp 431–438
179. Valenza E, Simion F, Cassia VM, Umiltá C (1996) Face preference at birth. *J Exp Psychol Human Percept Perform* 22(4):892
180. Van Hoof H, Chen N, Karl M et al (2016) Stable reinforcement learning with autoencoders for tactile and visual data. In: IEEE International Conference on Intelligent Robots and Systems, pp 3928–3934
181. Varela FJ, Thompson E, Rosch E (1991) The embodied mind: cognitive science and human experience. The MIT Press, Cambridge
182. Verschoor SA, Hommel B (2017) Self-by-doing: the role of action for self-acquisition. *Soc Cogn* 35(2):127–145
183. Vicente P, Jamone L, Bernardino A (2016a) Online body schema adaptation based on internal mental simulation and multisensory feedback. *Front Robotics AI* 3:7
184. Vicente P, Jamone L, Bernardino A (2016b) Robotic hand pose estimation based on stereo vision and GPU-enabled internal graphical simulation. *J Intell Robotic Syst Theory Appl* 83(3–4):339–358
185. Von Hofsten C (2004) An action perspective on motor development. *Trends Cogn Sci* 8(6):266–272
186. Watanabe H, Taga G (2006) General to specific development of movement patterns and memory for contingency between actions and events in young infants. *Infant Behav Dev* 29(3):402–422
187. Watanabe H, Taga G (2011) Initial-state dependency of learning in young infants. *Human Mov Sci* 30(1):125–142
188. Watter M, Springenberg JT, Boedecker J, Riedmiller MA (2015) Embed to control: a locally linear latent dynamics model for control from raw images. In: Advances in neural information processing systems, pp 2746–2754
189. Widmaier F, Kappler D, Schaal S, Bohg J (2016) Robot arm pose estimation by pixel-wise regression of joint angles. *IEEE International Conference on Robotics and Automation 2016*:616–623
190. Wijesinghe LP, Triesch J, Shi BE (2018) Robot end effector tracking using predictive multisensory integration. *Front Neurobotics* 12(October):1–16
191. Willatts P (1999) Development of means-end behavior in young infants: pulling a support to retrieve a distant object. *Dev Psychol* 35(3):651–667
192. Woodward AL, Sommerville JA (2000) Twelve-month-old infants interpret action in context. *Psychol Sci* 11(1):73–77
193. Woodward AL, Sommerville JA, Gerson S et al (2009) Chapter 6 the emergence of intention attribution in infancy
194. Yamada Y, Kanazawa H, Iwasaki S et al (2016) An embodied brain model of the human foetus. *Sci Rep* 6:27893
195. Zaadnoordijk L, Meyer M, Zaharieva M et al (2020) From movement to action: an eeg study into the emerging sense of agency in early infancy. *Dev Cogn Neurosci* 42:100760
196. Zaadnoordijk L, Otworowska M, Kwisthout J, Hunnius S (2018) Can infants' sense of agency be found in their behavior? Insights from babybot simulations of the mobile-paradigm. *Cognition* 181:58–64
197. Zacks JM, Speer NK, Swallow KM et al (2007) Event perception: a mind-brain perspective. *Psychol Bull* 133(2):273–293
198. Zambelli M, Cully A, Demiris Y (2020) Multimodal representation models for prediction and control from partial information. *Robot Auton Syst* 123:103312
199. Zambelli M, Demiris Y (2017) online multimodal ensemble learning using self-learned sensorimotor representations. *IEEE Trans Cogn Dev Syst* 9(2):113–126
200. Zenha R, Vicente P, Jamone L, Bernardino A (2018) Incremental adaptation of a robot body schema based on touch events. In: 2018 Joint IEEE 8th International Conference on Development and Learning and Epigenetic Robotics, ICDL-EpiRob 2018, pp 119–124
201. Zmyj N, Jank J, Schütz-Bosbach S, Daum MM (2011) Detection of visual-tactile contingency in the first year after birth. *Cognition* 120(1):82–89
202. Zoia S, Blason L, D'Ottavio G et al (2007) Evidence of early development of action planning in the human foetus: a kinematic study. *Exp Brain Res* 176(2):217–226