

## SURVEY PAPER

## Symbol Emergence in Robotics: A Survey

Tadahiro Taniguchi<sup>a\*</sup>, Takayuki Nagai<sup>b</sup>, Tomoaki Nakamura<sup>b</sup>,  
Naoto Iwahashi<sup>c</sup>, Tetsuya Ogata<sup>d</sup>, and Hideki Asoh<sup>e</sup><sup>a</sup>*College of Information Science and Engineering,  
Ritsumeikan University,**1-1-1 Noji Higashi, Kusatsu, Shiga 525-8577, Japan*<sup>b</sup>*Department of Mechanical Engineering and Intelligent Systems,  
The University of Electro-Communications,**1-5-1 Chofugaoka, Chofu-shi, Tokyo 182-8585, Japan*<sup>c</sup>*Faculty of Computer Science and Systems Engineering,  
Okayama Prefectural University,**111 Kubogi, Soja-shi, Okayama 719-1197, Japan*<sup>d</sup>*Department of Intermedia Art and Science,**School of Fundamental Science and Engineering, Waseda University,**3-4-1 Ohkubo, Shinjuku, Tokyo 169-8555, Japan*<sup>e</sup>*Artificial Intelligence Research Center,**National Institute of Advanced Industrial Science and Technology (AIST),  
AIST Tsukuba Central 1, 1-1-1 Umezono, Tsukuba, Ibaraki 305-8560, Japan**(v1.0 released \*\*\*)*

Humans can learn the use of language through physical interaction with their environment and semiotic communication with other people. It is very important to obtain a computational understanding of how humans can form a symbol system and obtain semiotic skills through their autonomous mental development. Recently, many studies have been conducted on the construction of robotic systems and machine-learning methods that can learn the use of language through embodied multimodal interaction with their environment and other systems. Understanding human social interactions and developing a robot that can smoothly communicate with human users in the long term, requires an understanding of the dynamics of symbol systems and is crucially important. The embodied cognition and social interaction of participants gradually change a symbol system in a constructive manner. In this paper, we introduce a field of research called *symbol emergence in robotics (SER)*. SER is a constructive approach towards an emergent symbol system. The emergent symbol system is socially self-organized through both semiotic communications and physical interactions with autonomous cognitive developmental agents, i.e., humans and developmental robots. Specifically, we describe some state-of-art research topics concerning SER, e.g., multimodal categorization, word discovery, and a double articulation analysis, that enable a robot to obtain words and their embodied meanings from raw sensory-motor information, including visual information, haptic information, auditory information, and acoustic speech signals, in a totally unsupervised manner. Finally, we suggest future directions of research in SER.

**Keywords:** Developmental robotics, language acquisition, semiotics, symbol emergence, symbol grounding

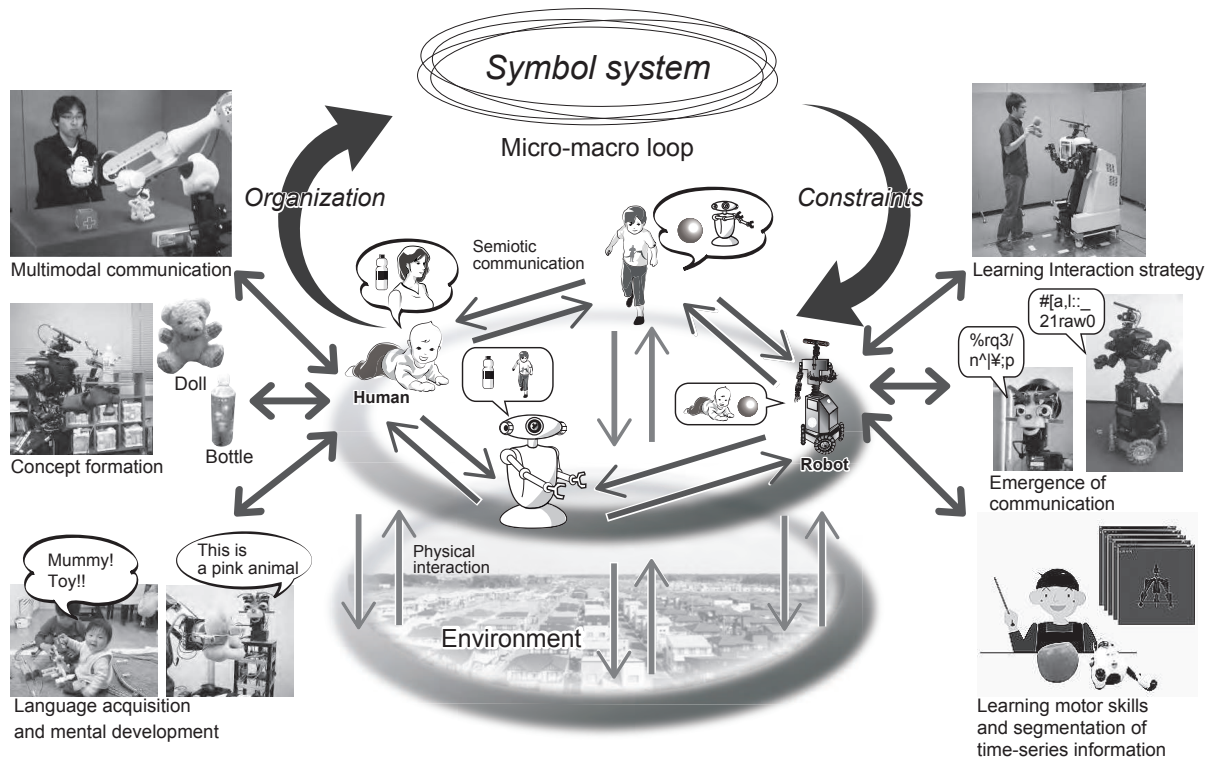


Figure 1. Emergent symbol system and research topics in symbol emergence in robotics.

## 1. Introduction

The development of an intelligent robot with which people would embrace a long-term interaction is one of the major challenges in the research field of robotics. Despite the rapid and remarkable progress in robotics, natural language processing, human–robot interaction, and related artificial intelligence technologies, we have not been able to develop such an autonomous robot. Even if an entertainment robot had sufficient capabilities for speech recognition, speech synthesis, and natural language processing, a user would find it unexciting if the robot behaved deterministically on the basis of finite hand-coded rules. To overcome this problem, a robot requires the ability of open-ended development through its physical interactions and semiotic communications. Moreover, such a robot must use a language system to communicate and collaborate with humans. To achieve actual long-term communication and collaboration between a human and a robot, the robot has to understand the meanings of utterances, estimate a speaker’s intention, learn new vocabularies, and promote a mutual understanding with people in the real world. To construct a mutual understanding between a human and a robot, both of them must be able to infer “what does he/she intend by the utterance?”, “what does the word he/she uttered represent?”, and “what should I say to make him/her understand what I want him/her to do?” by referring to many objects, events, contexts, situations, habits, and history. When we consider human–human communication and collaboration in the real world, it is easily understood that utterances, e.g., words, phrases, and sentences do not have one-to-one relationships with real-world phenomena. The modeling of such human semiotic communication for use in engineering applications is clearly important for developing a robot that can communicate and collaborate

---

\*Corresponding author. Email: taniguchi@ci.ritsumei.ac.jp

with others naturally like humans.

Generally, the language system is considered as a representative of the symbol systems in semiotics. A symbol system is an important philosophical and technical keyword, not only in semiotics, artificial intelligence, and cognitive science, but also in robotics. However, the adaptability and emergent properties of symbol systems have been underestimated and even ignored in the long history of research on robotic intelligence.

In contrast, studies that treat the adaptability of robots' internal representations and autonomous unsupervised learning processes for a language system have recently attracted attention [1–3]. One aspect of these studies is a constructive approach to human-embodied adaptive intelligence and human symbol systems. However, the majority of current approaches to symbol systems in artificial intelligence and robotics still cannot figure out the dynamics and emergent properties of human symbol systems, i.e., the symbol systems retained in our human society. A philosophical theory about the dynamics of human symbol systems should be established, and a more sophisticated understanding about human symbol systems should be obtained, in order to facilitate such studies and develop intelligent robots that enable long-term communication and collaboration with humans.

Based on this notion, the research field called “symbol emergence in robotics” (SER) has gradually emerged over the past decade, especially in Japan. We held the first organized session concerning SER in a domestic conference in Japan in 2011. SER is based on the concept of an “emergent symbol system,” which is introduced in Section 3. The emergent symbol system presumes that a human symbol system has emergent properties, and is self-organized through physical and semiotic interactions between cognitive agents, i.e., people and robots.

Figure 1 presents an abstract figure illustrating the emergent symbol system and research topics in SER. We assume that robots should be developed that semantically communicate with people and collaboratively interact with people in their environment, in order that they can be assimilated into the emergent symbol system. At the center of Figure 1, the dynamics of an emergent symbol system are schematically described. The background and concepts of emergent symbol systems are detailed in Section 2 and Section 3.

SER incorporates many research topics, such as multimodal communication, concept formation, language acquisition and mental development, learning interaction strategy, emergence of communication, and learning of motor skills and the segmentation of time-series information, as illustrated in Figure 1. In SER, robots are required to learn almost everything from their sensory-motor information flow, in a bottom-up manner. A top-down design of the intelligence of the robots would deprive them of the adaptability required for an emergent symbol system. In particular, concept formation that is grounded on a robot's multimodal sensory-motor information and autonomous language acquisition from raw speech signals are both fundamental topics in SER. SER aims to build computational models that can describe the overall dynamics and development of language acquisition and semiotic communication, on the basis of a robot's and a child's self-enclosed sensory-motor experiences. This would enable the establishment of a new theory of embodied semantics.

The remainder of this paper is organized as follows. Section 2 briefly reviews the history of “symbol systems” in artificial intelligence, cognitive science, and robotics. This forms the background of emergent symbol systems, which is a philosophical prerequisite of SER. Section 3 describes the concept of emergent symbol systems. In the subsequent sections, we provide a survey of works related to SER. In particular, we review previous studies on multimodal categorization (Section 4), word discovery (Section 5), and double articulation analysis (Section 6), which are important components of language acquisition, which is a fundamental challenge in SER. In Section 7, we describe further topics related to SER. Section 8 concludes this paper.

## 2. Background

### 2.1 *Physical symbol systems and robotics*

In the research field of robotics, the term “symbol” can be used in a variety of contexts, e.g., human–robot interaction, planning, reasoning, and communication. Historically, the physical symbol system hypothesis was proposed by Newell and Simon [4, 5]. This formed the starting point for a discussion about symbol systems in artificial intelligence and related fields. But, it was a problematic starting point. The philosophy is clearly inspired by early successes in computer science and programming languages. Many related works that followed the physical symbol system hypothesis and/or its way of thinking have placed an emphasis on the manipulation of symbols in research on artificial intelligence. This way of thinking was inherited from the tradition of “symbolic logic.”

In predicate logic, which is a representative of symbolic logic, predicates and variables that represent real-world phenomena are given as discrete representations in a top-down manner [6]. The fundamental assumption is that our world can be distinguished and segmented into a discrete “symbol” system, and that the system is deterministic and static. In other words, predicate logic can describe the world so far as such assumptions are satisfied. This represents a type of “approximation.” Almost all symbolic logic essentially shares the same assumptions. Physical symbol system hypothesis, proposed by Newel, is no exception [4, 5].

This convention has implicit effects on studies in robotics. In current robotics research, a “symbol” tends to be regarded as a “discrete” entity, having a “one-to-one” relationship with a word. A symbol is regarded as a manipulatable element in the mind, i.e., a robot’s memory system. For example, in [7] the authors call a type of trajectory of a humanoid’s entire bodily motion, modeled by a left-to-right hidden Markov model, a “proto-symbol.” Such notions as “a symbol is a discrete component of a memory system in a robot,” “a symbol is an internal representation in a robot,” and “a symbol system is a set or a network of such components,” have spread widely through the artificial intelligence and robotics communities.

Figuratively speaking, symbolic logic adopts the assumptions of equilibrium and determinism for modeling an actual human symbol system. It is assumed that the human symbol system is the same for everyone, and does not change over time. This approximation has been valid for solving many problems, in the same way that linear control theory has solved many problems, even though most real-world systems have nonlinear and stochastic properties, or that the theory of thermodynamics of equilibrium systems has provided fruitful results for engineering purposes.

However, these assumptions are crucially problematic, and have misled many researchers over the past four decades. This has resulted in people misunderstanding the human symbol system. The blind acceptance of the approximation has meant that people did not consider the following important characteristics of the human symbol system.

- C1** Grounded: A symbol does not have any meaning without being grounded or interpreted.
- C2** Dynamic: There does not exist an objectively true symbol system that can be determined in top-down manner in our human society.
- C3** Social: An individual representation system and the socially shared symbol system are not same.

The characteristics of symbols have been widely accepted in semiotics and in a broader context of humanities research [8].

### 2.2 *Physical grounding hypothesis*

There have been many criticisms of the physical symbol system hypothesis, and approaches to intelligent systems based on this hypothesis, in the field of artificial intelligence. Brooks representatively criticized and insisted that sensory-motor coupling with the environment is primarily important for robots to achieve everyday tasks in our daily environment [9, 10]. His famous paper

“intelligence without representation” provided a clear objection to the physical symbol system hypothesis. He proposed physical grounding hypothesis, of which the key observation is that “the world is its own best model.” He developed many robots based on a subsumption architecture, which is a reactive and decentralized robotic architecture. This behavior-based robotics places an emphasis on the primal sensory–motor interaction between a robot’s embodied system and its environment, and the emergence of behavior through interactions. Breazeal et al. even developed a “social” interactive robot using the subsumption architecture [11, 12]. However, the subsumption architecture is still a framework for “designing” a robot that behaves naturally in our daily environment. It is difficult for such an approach to build an autonomous system that gradually reaches an intelligent state such that it can communicate with people using a language system, i.e., a human symbol system.

The field of embodied cognitive science is closely related to Brooks’s approach [13]. This approach is also related to the research field of artificial life and complex systems. Metaphorically, Brooks’s approach is to develop an insect-like artifact. In contrast, the traditional approach to artificial intelligence can be thought of as an attempt to develop an obstinate mathematician who cannot behave appropriately in a real-world environment.

### 2.3 *Symbol grounding problem*

The other famous criticism in relation to the symbol system, was expressed by Harnad. He proposed the symbol grounding problem (SGP), which is one of the most famous problems in artificial intelligence [14]. The SGP focused on the relationship between a designed symbol system and real-world phenomena. The importance of the SGP has been widely recognized over the past two and a half decades. The design of a “symbol system” and application of it to an autonomous agent in a top-down manner inevitably leads to the SGP.

Advocates of physical symbol systems have insisted that the meaning of a symbol is syntactically determined in relation with other symbols. However, such a “relationship” cannot reach a conclusion on what anything means. A relationship between two signifiers can never provide the relationship between a signifier and a signified object. Harnad compared this phenomenon to a “merry-go-round”. To obtain any meaning, a word has to be grounded via sensory–motor information, or borrow meanings from other words using syntactic rules. Cangelosi et al. called these processes “sensorimotor toil” and “symbolic theft”, respectively [15].

In cognitive science, physical symbol systems have also been criticized. Barsalou proposed the concept of “perceptual symbol systems”, to place an emphasis on perceptual experiences for theories of knowledge [16]. He called the static symbolic system an “amodal symbol system,” and pointed out its drawbacks. Although the notion of a perceptual symbol system was not completely new, Barsalou mentioned that a perceptual theory of cognition may lead to a competitive, and perhaps superior, theory. The physical symbol system clearly relates to the SGP.

Many interdisciplinary studies have aimed to solve the SGP [2, 17–19]. Recently, Tellex et al. presented an approach to the SPG using probabilistic graphical models [20]. From a philosophical viewpoint, Taddeo et al. proposed the zero semantical commitment condition, which must be satisfied by any hypothesis seeking to solve the SGP [21].

However, despite the long history of the SGP, a clear solution has not been found. One of the reasons for this is that the SGP itself is naively defined. The SGP is based on the physical symbol system. Therefore, the SGP itself was misled by the physical symbol system hypothesis. That is, the SGP is an ill-posed problem. The SGP mainly considers **C1**, almost completely ignoring **C2** and **C3**. Steels pointed out the problem in an ambitious and impressive paper titled “The symbol grounding problem has been solved. So what’s next?” [22]. He described it as follows:

I propose to make a distinction between *c-symbols*, the symbols of computer science, and *m-symbols*, the meaning-oriented symbols in the tradition of the arts, humanities, and social and cognitive sciences.



This distinction is crucially important for the construction of theories concerning long-term human–robot interactions. The SGP starts with *c*-symbols, and attempts to make them grounded. To develop an intelligent robot that people would embrace long-term interactions with, we should clearly start from *m*-symbols, because from a human viewpoint, a human–robot interaction is composed of *m*-symbols. In our paper, we call a system of *c*-symbols and *m*-symbols as an internal representation system and a human symbol system, respectively, according to the conventions of robotics and semiotics. For example, Weng provided a critical survey on symbolic models and emergent models in artificial intelligence [23]. Those symbol models are about internal representation systems in our terminology.

In contrast, the emergent symbol system we introduce in the next section considers both types of symbol, in an integrative manner.

## 2.4 *Developmental robotics*

The epigenetic and/or developmental viewpoint is crucially important in creating artificial intelligent systems that can adapt to a dynamic real–world environment. The field of developmental robotics has emerged gradually over the past two decades [24]. Cangelosi et al. described developmental robotics as follows [24]:

Developmental robotics is an approach to the autonomous design of behavioral and cognitive capabilities in artificial agents (robots) that takes direct inspiration from the developmental principles and mechanisms observed in the natural cognitive systems of children.

The field is also referred to as “epigenetic robotics” [17]. Asada et al. used the term “cognitive developmental robotics” [25]. Asada et al. stated that cognitive developmental robotics places more emphasis on human/humanoid cognitive development than on related approaches. Our research field, SER, philosophically inherits many concepts and fundamental assumptions from the field of (cognitive) developmental robotics. In a manner of speaking, SER is a (crucially important) branch of developmental robotics.

Developmental robotics places an emphasis on an autonomous agent’s embodied interaction with the environment and the adaptive organization of the cognitive system, including cognitive capabilities relating to language and other symbol systems [24]. However, the scope of developmental robotics is tremendously large, because it involves almost all of human intelligence and its diachronic changes. Moreover, developmental robotics attaches importance to interdisciplinary communication between robotics and developmental psychology. Many efforts have been made to construct a fruitful interdisciplinary academic field.

However, these characteristics of developmental robotics have distracted its attention from a computational and constructive understanding of dynamic human symbol systems and the development of robots that achieve the overall dynamics and development of language acquisition and semiotic communication. We believe that these are central topics in robotics research for achieving long-term human–robot communication and collaboration. This is our motivation for introducing the field of SER.

## 2.5 *Symbol emergence in robotics*

The approach in SER places more emphasis on the computational understanding of emergent symbol systems. In addition to cognitive development, SER attempts to cover semiotic phenomena. The field of SER is an interdisciplinary field, which is not only related to robotics, artificial intelligence, development psychology, and cognitive science, but also to semiotics and linguistics as well.

To describe the diachronic changes in internal representation systems and human symbol systems that are caused by embodied interaction and social communication, we require mathematical models, such as generative models, neural networks and related statistical models, and robotic

models such as humanoids and mobile robots, for a productive discussion and development of the integrative theory. In cognitive science, the generative probabilistic model has recently been widely used to represent the human cognitive system [26]. In addition to such computational models, SER places an emphasis on embodied cognition. Therefore, researchers in the field of SER use robotic models to connect computational models to the real physical world. This involves the use of state-of-art machine-learning technology, including Bayesian nonparametrics and deep neural networks, to model the diachronic changes in the cognitive systems and symbol systems of human/robots.

### 3. Emergent Symbol Systems

The center of Figure 1 shows the schematic figure of an emergent symbol system that was originally introduced in [27]. SER is defined as a constructive approach towards emergent symbol systems [28]. In this section, we explain emergent symbol systems by referring to the figure.

#### 3.1 *Semiosis and umwelt*

We will start from a human symbol system, i.e., the m-symbols described by Steels [22]. The preexisting interdisciplinary research field that deals with human symbol systems is called semiotics. Semiotics is concerned with everything that can be interpreted as a sign, as explained by Eco [29]. Initially, semiology was introduced by Saussure, while Peirce independently introduced semiotics [30, 31]. Currently, the two fields have overlapped, and merged into the academic field called semiotics. From the viewpoint of semiotics, language is a representative of general symbol systems.

In Peircean semiotics, a symbol is defined as a process having three elements. The definition has a high affinity for the bottom-up approach to cognitive systems. The first is the *sign* (*representamen*), which describes the form that the sign takes, the second is the *object*, which is something that the sign refers to, and the third is the *interpretant*, which, rather than an interpreter, is the sense made of the sign. The important point of the Peircean definition of a symbol is that the sign, e.g., words, visual signs, or pointing, is not a symbol itself. The interpretant, the third element of a symbol, mediates between the sign and the object. This degree of freedom allows us to take a variety of interpretations and dynamics of a symbol system into consideration. In the Peircean definition, a symbol is not a static material, but a dynamic process of interpretation. Peirce calls this process “semiosis” [31].

The definition of a symbol is still abstract, but the definition clearly satisfies **C1**, **C2**, and **C3** from Section 2. The utterances of others are always interpreted on the basis of semiosis in **semiotic communication**.

Peircean semiotics places a thorough emphasis on the subjective viewpoint. Uexküll, who established biosemiotics, proposed the famous notion of *umwelt*. The *umwelt* represents an animal’s subjective world, which emerges on the basis of the animal’s sensory–motor system [32]. Brooks also cited Uexküll, when he attacked the physical symbol system in his famous paper [9]. We should start take the *umwelts* of robots and humans as a starting point. Their internal representation systems are initially formed through **physical interactions** with their **environment**, using their sensory–motor systems.

In this sense, *semiosis* is the key that connects Brooks’ physical grounding hypothesis, which eliminated internal representation systems, to semiotic communication, which is required for long-term human–robot communication.

### 3.2 *Arbitrariness and perspective of structuralists*

In contrast to Peirce, Saussure emphasized the synchronic structure of language. The defining notion of Saussurean semiotics is the *arbitrariness* of the sign (symbol). This embodies **C2**. The relationships between signs, such as labels and words, and categories are arbitrary, and the categories and segments of phenomena are also arbitrary. The arbitrarily determined categories, segments, and lexicons are retained in a language system to which many people, who speak the language, belong. Structuralists, the successors of Saussure, place an emphasis on arbitrariness. When we belong to a language system, i.e., a human symbol system, our cognition, interpretation, utterances, and even behaviors are affected by the symbol system. For example, we comprehend objects so as to classify them into preexisting categories that our language system retains. The symbol system provides **constraints** on our semiotic communication and physical interaction.

According to structuralists, “things” do not exist independently of the symbol system that we use; reality is the creation of the media that seems to simply represent it. The structuralist perspective tends to reverse the precedence of language and cognition. They stress that our language, which incorporates arbitrariness, determines the order of the world [33]. This places top-down constraints in an emergent symbol system.

### 3.3 *Emergent symbol systems*

This structuralist exaggeration of unilateral determination is not accurate. Human symbol systems can be **organized** in a bottom manner. Genetic epistemology was proposed by Piaget [34], who is often called the father of cognitive development research. In genetic epistemology, the subjective world of humans is considered to be gradually “constructed” through interactions with their environment. Piaget introduced a schema system, which is a self-organized cognitive system that emerges through sensory–motor interaction, and is believed to be the basis of the language system [35].

An internal representation system is not a static system, but rather a dynamic system that is self-organized through physical interaction based on the sensory–motor system. Furthermore, a human symbol system is organized through semiotic communication on the basis of individual internal representation systems in a bottom-up manner (see *organization* in Figure 1). In semiotics, symbol systems are seldom treated as the static, closed, and stable systems that are inherited from preceding generations, but instead are regarded as constantly changing [8]. The bottom-up organization and the top-down constraints of the symbol system introduce an emergent property to the overall system.

Once a symbol system is generated in a society, people who use the symbol system must obey the rules of this system to communicate and collaborate with others. The symbol system includes phonetics, lexicons, syntax, and pragmatics as its constituents. If an agent belonging to the society does not follow the rules, i.e., the symbol system, then the agent cannot communicate its idea to others or collaborate with others. This means that the agent cannot make use of the powerful symbolic system for their further survival.

Such a bilateral relationship between an emerged symbol system at the macro level and a physical system consisting of communicating and collaborating agents at the micro level forms a **micro-macro loop**. Micro-macro loops are found in many complex systems, especially in living systems. This tells us that the entire system is an *emergent system*, i.e., a complex system having emergent properties. Polanyi, who introduced the notion of emergence, described it as follows [36]:

If each higher level is to control the boundary conditions left open by the operations of the next lower level, this implies that these boundary conditions are in fact left open by the operations going on at the lower level.

This stratification offered a framework for defining *emergence* as the action that produces the next higher level, first from the inanimate to the living, and then from each biotic level to the one above





#### 4.1 *Object category formation*

Before obtaining language, human children are considered to obtain object categories gradually through daily interactions with objects. Piaget insisted that the schema system self-organizes through the sensory–motor period, and that the system becomes the prerequisite for language [35]. An embodied multimodal sensory–motor experience must be a primal root for human category formation.

A category formation problem is different from a pattern recognition problem. In a pattern recognition problem, truth labels for the recognition results are provided in a supervised manner. A vast number of studies have been carried out regarding the development of an accurate pattern recognizer. Recently, deep learning methods have yielded excellent results [38, 39]. In contrast, category formation in a robot’s *umwelt* must be autonomously performed in an unsupervised manner. From the viewpoint of machine-learning, object categorization is regarded as a clustering problem, involving a type of unsupervised machine-learning tasks [40].

Historically, many studies have emphasized visual information in category formation by computational systems. However, the formation of object categories based solely on visual information is insufficient, because our human symbol system is organized on the basis of our multimodal sensory–motor experiences. The integration of multimodal information through category formation is important for a robot to predict future sensor information. By forming an object category on the basis of visual, auditory, and haptic information, a robot can infer auditory and haptic information from the recognized category, e.g., a bottle and a cymbal, from its visual information.

#### 4.2 *Computational models for multimodal object categorization*

Recently, various computational models for multimodal object categorization have been proposed [3, 18, 19, 41–51]. For example, Sinapov et al. proposed a graph-based multimodal categorization method, which allows a robot to recognize new objects on the basis of similarities to a set of familiar objects [42]. They also made a robot perform ten different behaviors; obtain visual, auditory, and haptic information; and explore 100 different objects, classifying them into 20 object categories [19]. However, their multimodal categorization is performed in a supervised manner. Celikkanat et al. modeled the context in terms of a set of concepts, allowing many-to-many relationships between objects and contexts using latent Dirichlet allocation (LDA), inspired by the notion of situated concepts introduced by Yeh and Barsalou [41, 52]. Mangin used a nonnegative matrix factorization algorithm to learn a dictionary of components from multimodal time series data [53]. Natale et al. have demonstrated that a robot can recognize objects with the help of a self-organizing map (SOM), using proprioceptive data extracted from the robot’s hand as it grasps an object [43]. Lallee et al. proposed multi-modal convergence maps on the basis of SOMs. This method can integrate visual, motor, and language modalities [54]. Invalid et al. proposed a cognitive architecture, and developed a child-like robot that can automatically learn object categories through active exploration [55].

Nakamura et al. have presented a series of studies on multimodal categorization using multimodal latent Dirichlet allocation (MLDA) and its extensions [3, 18, 44–48, 56]. They have extended LDA, which was first proposed by Blei et al. for document-word clustering, to a model that can treat multimodal information [3, 57, 58]. Concrete illustrations of the graphical models of LDA and MLDA and its extended models are presented in Figure 3. MLDA has several emission distributions for an object, i.e., a document in LDA. The object categories of objects in multimodal categorization in MLDA correspond to topics in document-word clustering in LDA.

Nakamura et al. developed a robotic system that can obtain visual, audio, and haptic information by interacting with objects. An overview of the robot is presented in Figure 4. The robot can grasp an object and observe it from various viewpoints. The robot has cameras, microphones, arms, and hands with pressure sensors. The robot obtains visual information by taking pictures of a target object from many directions, by rotating the object with its hand. The robot also obtains haptic information by grasping the target object several times, and audio information by

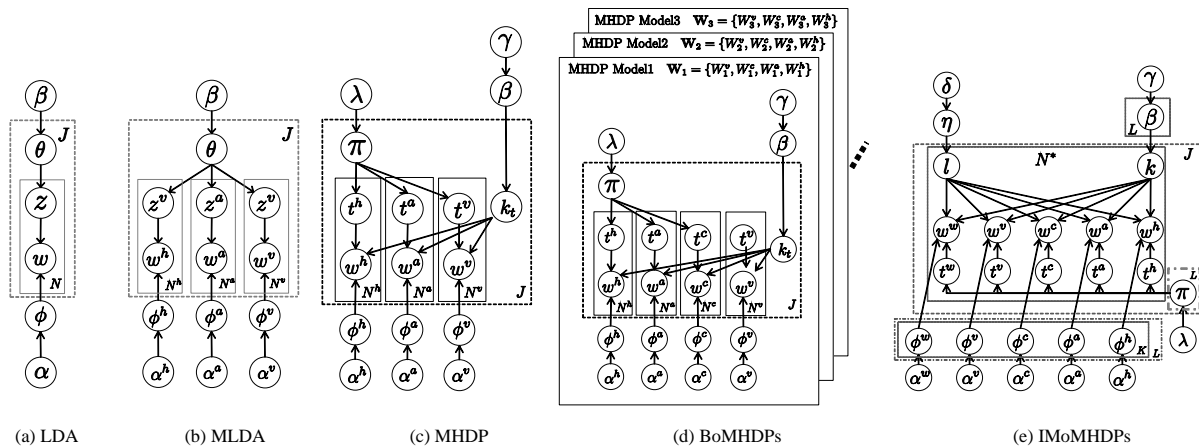


Figure 3. Graphical models for object categorization.

shaking the target object. Feature vectors are extracted from the observed information of each modality, and the feature vectors of each modality are transformed into bag-of-features representations using the K-means method, i.e., vector quantization. The bag-of-features representations are passed to MLDA, and the clustering procedure is performed. Nakamura et al. demonstrated that a robot can categorize a large number of objects in a home environment into categories that are similar to human categorization results [3]. Araki et al. developed online MLDA, and performed an experiment on multimodal category acquisition in a fully autonomous manner in a home environment [59]. The result indicates that a human symbol system is not necessarily required for a cognitive system to form human-like object categories. This suggests that the human symbol system is not completely arbitrary, but has a certain rationality brought by the latent structure embedded in sensory-motor information, as admitted by Saussure [30].

### 4.3 Estimating latent structure in multimodal categories

Although MLDA is able to form multimodal categories, it can not adaptively estimate the number of object categories. It is unlikely that the number of categories that a cognitive system has is determined in advance. The Bayesian nonparametric approach provides a reasonable solution to the problem. Nakamura et al. extended the hierarchical Dirichlet process (HDP), which was a nonparametric Bayesian clustering method proposed by Blei et al. for document-word clustering, to multimodal HDP (MHDP), which can treat multimodal information and automatically estimate the number of categories from the observed multimodal information [47, 60, 61].

Bayesian nonparametrics is a branch of the Bayesian approach. In general, the number of hidden variables in the Bayesian nonparametric approach can be automatically estimated by using the infinite dimensional prior distribution, for example using Dirichlet or Beta processes, and a feasible inference procedure [60, 62]. The mathematical framework is very important for constructing a computational model relating to emergent symbol systems. Nakamura et al. demonstrated that a robot can also estimate the number of object categories to give similar results to human categorization results [47].

MLDA and MHDP can easily be extended to treat “words.” LDA and HDP were originally applied to document and word clustering methods [57, 60]. By adding an observation variable for words to MLDA’s graphical model, MLDA is able to cluster multimodal information and words simultaneously. As a result, a robot can estimate the label of a category in an unsupervised manner.

More complex latent structures of multimodal categories can be estimated. Ando et al. proposed hierarchical MLDA (hMLDA), by extending hierarchical LDA (hLDA) for hierarchical multimodal categorization [44, 63]. This method enabled a robot to form a hierarchical struc-

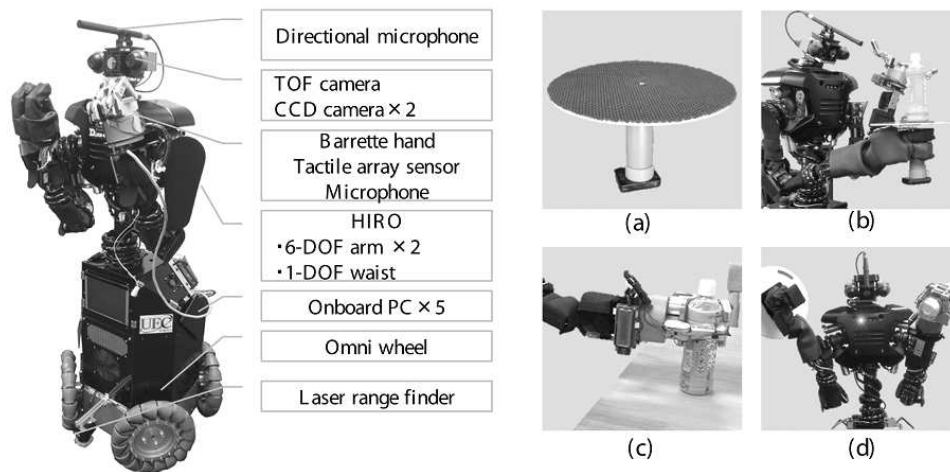


Figure 4. The robot used in multimodal categorization experiments.

ture of object concepts from multimodal sensory–motor information, e.g., “plastic bottle” is a subcategory of “water container.”

Nakamura et al. proposed the Bag of MLDA (BoMLDA), Bag of the MHDP (BoMHDP), and infinite mixture of MHDPs (IMoMHDPs) methods, which can perform various types of categorization with different perspectives [46, 56, 64]. The methods emphasize some modalities, and organize categories based on information of the modalities focused on. By adopting these methods, their robot was able to form concepts about not only “objects”, but also “attributes,” e.g., soft, hard, green, red, or yellow, from multimodal information.

Compared to related methods for multimodal categorization, the MHDP-based approach is sophisticated from the viewpoint of Bayesian modeling. Its mathematical soundness and theoretical consistency help us to build new methods that are based on it, e.g., active perception [65].

#### 4.4 *Multimodal representation learning using neural networks*

Another way to integrate multimodal information involves approaches that utilize neural networks. Ngiam et al. applied deep networks to learn features over multiple modalities. By integrating visual, i.e., lip motions, and auditory information, they developed a robust speech classification system. They also demonstrated that the system shows McGurk effect that is an audio-visual perception phenomenon in the similar way as human [66]. Noda et al. integrated auditory, visual, and motor information using a deep neural network [67]. In their experiment, a robot could recall upper bodily motion from visual and audio information, and retrieve image information from sound and joint angle inputs. Heinrich et al. extended the multiple timescale recurrent neural network (MTRNN), and obtained multi-modal MTRNN, to integrate visual, auditory, and motor information. Recently, neural networks with a deep architecture have attracted attention. Le et al. demonstrated that large-scale unsupervised learning using a deep neural network could build high-level features automatically from image data [68]. Bridging representation learning using neural networks and multimodal categorization using generative models represents important work in this field [69].

## 5. Word Discovery

### 5.1 *Word discovery by human children*

In language acquisition, word discovery, i.e., word segmentation, is an important task for children. A word is an elemental pattern of a linguistic sign. A phoneme is an element that is

acoustically but not semantically distinguishable. Discovering words from continuous speech signals is a fundamental task that children must solve to acquire language. Unlike an automatic speech recognition system, children must learn a language model, i.e., a word inventory and transitional information about the words, and in an acoustic model, i.e., an organized memory about phonemes, this must be done from speech signals in an unsupervised manner.

What types of cue can be used by children to discover words from continuous speech signals? Three representative cues for word segmentation are listed by Saffran et al. [70]. These are *distributional*, *co-occurrence*, and *prosodic* cues. Distributional cues concern the statistical relationships between neighboring speech sounds. These can be modeled as n-gram statistics to some extent, once each phoneme is recognized correctly. Co-occurrence cues concern entities detected in the environment by children. For example, when a child hears two sentences while he/she looks at an apple, the sentences are likely to contain overlapping words, such as “apple.” Prosodic cues relate more to superficial acoustic information, such as stressed syllables, post-utterance pauses, and acoustically distinctive final syllables.

All of the above cues are believed to contribute to word discovery in an integrative manner. Among these, Saffran emphasized distributional cues. She reported that word segmentation from fluent speech could be accomplished by eight-month-old infants using only distributional cues [71]. By the age of seven months, infants are reported to use distributional cues [72].

## 5.2 *Word discovery by robotic systems*

In SER, autonomous word discovery by robots is one of the first challenges that should be solved. In the past two decades, many types of unsupervised machine-learning methods for word discovery (segmentation) have been proposed [73–81]. Conventionally, Brent proposed model-based dynamic programming to find word boundaries in a natural-language text whose word boundaries are deleted [73]. Venkataraman proposed a statistical model to improve Brent’s algorithm [74].

In contrast with such text-based approaches, Roy et al. developed a computational model and a robotic system that autonomously discovers words from a raw multimodal sensory input [51]. The experimental results of Roy et al. demonstrated the development of a cognitive robot that can acquire a lexicon from raw sensor data without human transcription or labeling. Although imperfect, the results were encouraging. Their results showed that it is possible to develop cognitive models that can process raw sensor data and acquire a lexicon, without the need for human transcription or labeling.

Contemporaneously, Iwahashi et al. independently proposed a sophisticated probabilistic method that enables a robot to acquire linguistic knowledge, including speech units, lexicons, grammar, and interpretation through human–robot embodied communication, in an unsupervised manner [82]. This integrated speech, visual, and behavioral information in a probabilistic framework. This work was updated by Iwahashi et al. [83]. The learning process was carried out online, incrementally, actively, and in an unsupervised manner. On the basis of this work, Iwahashi et al. developed an integrated online machine-learning system called LCore, which combined speech, visual, and tactile information obtained through interactions, and enabled robots to learn beliefs regarding speech units, words, the concepts of objects, motions, grammar, and pragmatic and communicative capabilities [50]. These pioneering studies clearly demonstrated the possibility of the SER approach.

## 5.3 *Nonparametric Bayesian word segmentation*

In word discovery and segmentation tasks, the efficient management of a word inventory through the learning process is a fundamental computational problem. Although a robot can only memorize a finite number of words, there are potentially an infinite number of words in our society, i.e., in the human symbol system. The selection of an appropriate set of words constitutes a type of



model selection problem, and usually requires a very large computational cost. Recently, Bayesian nonparametrics have provided a sophisticated theoretical solution to this problem [26, 60, 62]. A nonparametric Bayesian language model, e.g., a hierarchical Pitman-Yor process language model (HPYLM), can assign an adequate probability to an infinite number of possible words using a fully Bayesian framework [84]. On the basis of this framework, word segmentation methods that assume that there is an infinite number of possible words can be developed. Goldwater proposed an HDP-based word segmentation method [75, 76]. Mochihashi et al. proposed a nested Pitman-Yor language model (NPYLM), in which a letter n-gram model based on a hierarchical Pitman-Yor language model is embedded in the word n-gram model [77]. An efficient blocked Gibbs sampler, employing the forward filtering backward sampling procedure, was also introduced in that study. These methods have made it possible to discover words from transcribed phoneme sequences or text data without any recognition errors.

However, in practice phoneme recognition errors are inevitable, especially during the language acquisition phase. In order to overcome this problem, several extensions have been proposed. Neubig et al. extended the word segmentation methods of Mochihashi et al., making it possible to analyze phoneme lattices that are a type of expression of noisy speech recognition results [85]. Heymann et al. modified the algorithm of Neubig et al., and proposed a suboptimal algorithm [86, 87]. Elsner et al. developed a learning method that jointly performs word segmentation and learns an explicit model of phonetic variation [88]. Hinoshita et al. solved a similar problem using MTRNN [89].

Recently, several advanced machine-learning methods have attempted to learn words from acoustic data without a preexisting language model and acoustic model [90–93]. However, the problem remains a challenging topic.

#### 5.4 *Integrative word discovery by robots*

It is difficult to discover words from uttered sentences using only distributed cues. Researchers have developed robotic systems and computational models in which co-occurrence cues help to discover words using distributional cues. Taguchi et al. proposed a method for the unsupervised learning of place-names, from information pairs that consist of spoken utterances and a mobile robot’s estimated current location, without any prior linguistic knowledge other than a phoneme acoustic model [94]. They optimized a word list using a model selection method based on a description length criterion. Araki et al. proposed an integrative computational model that involved MLDA and NPYLM [18]. Through MLDA, a robot can detect co-occurrence cues in the environment, and use the information to increase word segmentation performance. It was shown that the iterative learning process comprising MLDA and NPYLM increased the word segmentation accuracy. However, they reported that the accuracy decreases as the phoneme recognition error rate increases [18]. This implies that phoneme recognition errors and word segmentation errors should be mitigated simultaneously. Nakamura et al. developed an integrated statistical model for word segmentation, speech recognition, and multimodal categorization, in order to overcome this problem. The robot in the experiment simultaneously formed object categories and learned related words from continuous speech signals and continuous visual, auditory, and haptic information, i.e., sensory–motor information, through an iterative learning process [48].

Various word discovery methods which enable robots to obtain words and relationships between words and objects. However, situations in which robots can currently discover words are still limited. To simulate the robust word discovery process by human children in real-world environment, further studies will be required.

## 6. Double Articulation Analysis

### 6.1 *Double articulation structure*

Double articulation is an important property of human language systems. Chandler described double articulation as follows in a textbook on semiotics [8]:

One of the most powerful “design features” of language is called double articulation (or “duality of patterning”). Double articulation enables a semiotic code to form an infinite number of meaningful combinations using a small number of low-level units which by themselves are meaningless (eg., phonemes in speech or graphemes in writing). The infinite use of finite elements is a feature that about media, in general, has been referred to as “semiotic economy.”

Our speech signals and some semiotic time-series data are considered to have a double articulation structure. This means that a sentence can be decomposed into words, and a word can be decomposed into letters or phonemes. Automatic speech recognition systems usually presume double articulation in speech signals. A continuous speech signal is first segmented into phonemes, such as “a,” “e,” “i,” and “s.” Then, the phonemes are chunked into words, such as “can,” “dog,” and “pen.” A phoneme cannot usually act as a sign for an object in the sense of Peircean semiosis, but a word constitutes a sign, i.e., it has certain meaning. Usually, the relationship between speech signals and phonemes is stored in a phonetic model and/or acoustic model, and the relationship between phonemes and words is stored in a language model. Therefore, direct word discovery from speech signals can be regarded as the analyzing of a latent double articulation structure embedded in speech signals in an unsupervised manner.

In addition to speech signals, other time-series data generated by humans may have a double articulation structure. If such time-series data exists, then the analysis of such data would contribute to our understanding of emergent symbol systems. For this purpose, several researchers have developed a computational model that can automatically analyze double articulation structures [92, 93, 95–98].

### 6.2 *Segmentation of human bodily motion*

Human bodily motion is a candidate for doubly articulated time series data. Inamura et al. proposed the use of a left-to-right hidden Markov model (HMM) to recognize and reconstruct human bodily motion [7]. Essentially, left-to-right HMMs are often used to model words in automatic speech recognition systems. These computational models implicitly bridge speech recognition and human motion modeling.

There have been many previous studies on motion segmentation. However, the definition of a unit of motion has been unclear in many studies for a long time. Roughly speaking, some researchers have focused on physically elemental segments [99–103], and others on semantically elemental segments [104–106]. For example, when we semantically segment a baseball player’s motion, “pitching” definitely becomes a candidate for a unit motion. However, pitching consists of several low-level segmental motions from the viewpoint of physical dynamics. If we segment the pitching motion according to the criterion that an elemental motion has linear dynamics, then the pitching motion will be segmented into several elemental motions, e.g., “raising a knee” and “swinging an arm.” However, these elemental motions seem to be meaningless, and have no special names. The two-layer hierarchical structure is similar to that existing in speech signals. We call a short-term motion that corresponds to phoneme a *segment*, and a long-term motion that corresponds to word a *chunk*. A chunk corresponds to a sequence of segments.

Takano et al. developed a large-scale database of human whole-body motion, and modeled the motion data using a large number of HMMs [96, 107]. They roughly clustered the given motions, and constructed many left-to-right HMMs corresponding to meaningful motions. They hierarchically clustered the HMMs again, and obtained a large motion database. An online incremental learning method was also provided by Kulić [108]. They implicitly assumed double articulation

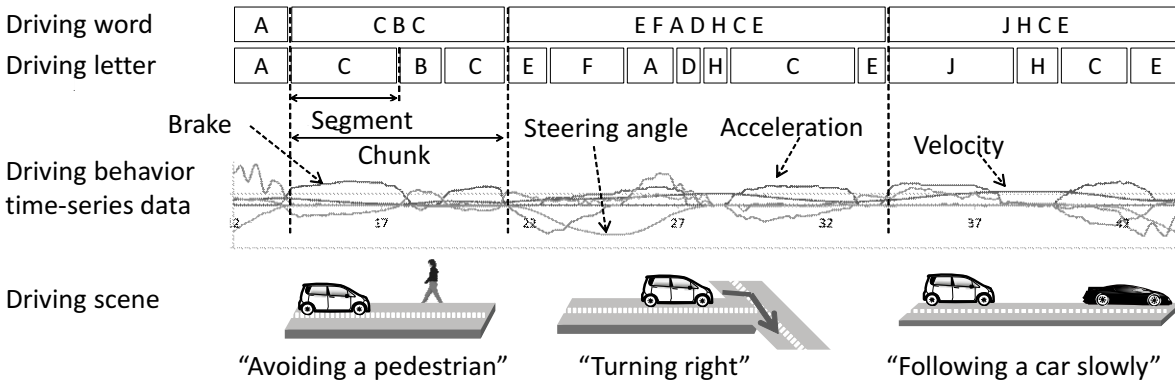


Figure 5. Double articulation in driving behavior.

in human bodily motion. Based on the database, they developed a machine translation system that can translate continuous human motion into a sentence, and vice versa [109, 110].

Taniguchi et al. proposed a double articulation analyzer (DAA) by combining a sticky hierarchical Dirichlet process-HMM (sticky HDP-HMM) and NPYLM [95, 111]. The DAA explicitly assumes double articulation, and infers the latent letters, i.e., the segment or phoneme, and the latent words, i.e., the words or segments, in an unsupervised manner. These two nonparametric Bayesian models, sticky HDP-HMM and NPYLM, were sequentially applied to the target data, and a two-layered hierarchical structure was inferred in the DAA. They applied the DAA to human motion data, to extract unit motions from unsegmented human motion data.

### 6.3 Modeling driving behavior data

Data on driving behavior represents another candidate for doubly articulated time-series data. A meaningful chunk of driver behavior seems to consist of sequences of simple segments. Figure 5 presents an example. In this figure, when a driver “turns right at an intersection,” he/she “steps on the brake,” “turns the steering wheel to the right,” “steps on the accelerator,” “turns the steering wheel back to center,” and “steps on the accelerator” again as a sequence of physically elemental driving behaviors.

The DAA has been utilized for various application, such as segmentation [97], prediction [112, 113], data mining [114], topic modeling [115, 116], and video summarization [117]. Through experiments in a series of studies, it has become clear that driving behavior data has a double articulation structure. For example, a prediction method based on the DAA has outperformed conventional methods in a driving behavior prediction task [118]. This implies that the assumption of double articulation is appropriate. Taniguchi et al. called the latent states corresponding to semantically elemental driving behaviors *driving words*, and those corresponding to physically elemental driving behaviors *driving letters*. Recently, the DAA was applied to large-scale driving behavior data, and its effectiveness for driving behavior analysis has been verified [119].

### 6.4 Predicting long-term sensory-motor information

A third candidate for doubly articulated time-series data is given by sensory-motor information flow. Unlike speech signals, sensory-motor information flow data comprises sensor information as well as motor information. Modeling sensory-motor time-series data for a robot means estimating the forward dynamics of the system that the robot confronts. Many modular learning architectures, e.g., MOSAIC (modular selection and identification for control), HAMMER (hierarchical, attentive, multiple models for execution and recognition), and the dual schemata model, have been proposed to model switching forward dynamics [120–124]. Most of these involve forward-

inverse models in the learning system. Methods for imitation learning, reinforcement learning, and the emergence of communication have been proposed on the basis of such modular learning architectures [125–128].

In dynamic environments, the forward dynamics of a robot change intermittently. Such contextual information tends to have a certain structure. Tani et al. proposed a the use of hierarchical mixture of RNNs [98]. In their experiment, a longer context was coded into the activations of the context nodes of the RNNs at a higher level, and a shorter context was coded into those at a lower level. In other words, their model captured the double articulation structure in the environmental dynamics. Hierarchical MOSAIC and HAMMER architectures have also been proposed as computational models that could capture such hierarchical structures [121, 129]. MTRNN and recurrent neural networks with parametric bias (RNNPB) are candidates to model contextually changing forward dynamics [2, 130–132]

### 6.5 *Direct word discovery from speech signals*

As mentioned above, the double articulation analysis is deeply related to direct word discovery from an acoustic speech signal. Several studies have recently been carried out in relation to this difficult problem. The direct application of the DAA proposed by Taniguchi et al. [95] is one possible approach; however, poor results are expected to be obtained. In that approach, the DAA simply uses the two nonparametric Bayesian methods sequentially. These are not integrated into a single generative model. Therefore, if there are many recognition or categorization errors in the result of the first segmentation process by sticky HDP-HMM, then the performance of the subsequent unsupervised chunking by NPYLM deteriorates.

To overcome this problem, Taniguchi et al. proposed a nonparametric Bayesian double articulation analyzer (NPB-DAA), which is a two-layered generative model [93]. The generative model represents a complete data generation process that has a double articulation structure. An efficient blocked Gibbs sampler was also derived in the same study. They showed that the NPB-DAA could automatically find a word list from vowel speech signals directly and completely, in an unsupervised manner. Additional approaches have recently been proposed [90, 92].

It is interesting that superficially different time-series data generated from human behavior can be analyzed using almost the same computational model. In addition, the characteristic of double articulation is satisfied. That is, the elements in the first layer, i.e., the phonemes and segments, are meaningless, and the elements in the second layer, i.e., words and chunks, are meaningful. This suggests that different examples of such time-series data potentially share the same computational processes in our brain. In addition, we hypothesize that these are profoundly related to the nature of our emergent symbol system.

## 7. Further Topics

In order to construct a computational model that describes an entire emergent symbol system, and to develop a robot that can communicate and collaborate with humans in a long-term manner, there are many other challenges in the field of SER to consider and overcome. In this section, we describe some of them.

### 7.1 *Mutual belief system*

An utterance is not usually interpreted “as it is” by another person to whom it is told. A mutual belief system always affects a person’s interpretation. Roy provides a coffee scenario as an example in his survey paper [133]. Imagine a situation in which a cup of coffee is served to a customer by a waitperson, and the customer says, “This coffee is cold.” In this case, the referential meaning of the utterance is the fact that the temperature of the coffee is low in the

sense of thermodynamics, but the functional meaning is “please get me a hotter coffee.” The speech act conveys a meaning interpreted by referring to the physical situation shared by the communication partners. This means that the mutual belief system is important in generating natural sentences and interpreting uttered sentences. Roy emphasized this aspect of language for solving the symbol grounding problem.

A pioneering study that develops a constructive model involving a mutual belief system was presented by Iwahashi [82]. Iwahashi introduced a belief function that represents a mutual belief of a robot. The belief function contains several belief modules, i.e., speech, object images, motions, motion–object relationships, and behavioral contexts. The various external and internal contexts are taken into consideration to infer the speaker’s intention. The robotic system is truthfully an embodied natural language processing system that can take various contexts that an ordinal amodal natural processing system cannot access. Zuo et al. applied this model for detecting robot-directed speech [134]. Sugiura et al. proposed a method for estimating the ambiguity in commands by introducing an active learning scheme to the conversation system, based on the mutual belief [135].

A mutual belief is part of an emergent symbol system, and applies “constraints” not only to the interpretation of utterances, but also to the generation of our speech. Context is a crucial element of our natural dialog but is rarely taken into consideration in natural language processing. The SER approach is promising in this topic as well.

## 7.2 *Active perception and learning*

A robot that has the potential to be an element of an emergent symbol system, i.e., a member of our semiotic society, must be able to explore its environment, acquire knowledge, and communicate with people autonomously. Active perception and active learning are two of the most important capabilities of humans for achieving life-long development and communication [43, 133, 136–147].

Denzler et al. proposed an information theoretic action selection method to gather information that conveys the true state of a system through an active camera [136]. They used mutual information (MI) as a criterion for action selection. Krainin et al. developed an active perception method that made a mobile robot manipulate objects to build three-dimensional surface models of the objects [137]. Their method determines when and how a robot should grasp an object on the basis of the information gain (IG) criterion.

Modeling and recognizing a target object, as well as modeling a scene and segmenting objects from that scene, are important abilities for a robot in a realistic environment. An active perception planning method for scene modeling in a realistic environment was proposed by Eidenberger et al. [138]. A partially observable Markov decision process (POMDP) formulation was used to model the planning problem, and the differential entropy was introduced as part of the reward function. Hoof et al. proposed an active scene exploration method. An autonomous robot is able to segment a scene into its constituent objects by actively interacting with the objects using this method [139]. They used IG as a criterion for action selection. InfoMax control for acoustic exploration was proposed by Rebguns et al. [140]. In general, IG is mostly applied in active perception and learning [136, 137]. Taniguchi et al. developed an optimal active perception scheme for multimodal categorization using MLDA on the basis of the IG criterion. Localization, mapping, and navigation are also important targets of active perception [133, 141, 142]. In addition, various other studies on active perception have been conducted [43, 143–147].

In the context of reinforcement learning (RL), intrinsic motivation has been studied [148]. In reinforcement learning, an agent autonomously learns its policy, i.e., controller, to maximize the expected cumulative rewards. In related studies, internal reward systems are taken into consideration as well as external reward systems. Schmidhuber presented a survey on RL studies considering intrinsic motivation, and proposed a formal theory of fun, intrinsic motivation, and creativity [149].



When we develop an autonomous robot, the design of the intrinsic motivation and explanatory behavior, active perception, and in particular the IG criterion provide cues for the problem. The effectiveness of the IG criterion tells us that “curiosity” is properly treated computationally, in contrast with other emotions.

Sugiura et al. used active perception to determine the utterances of a robot [135]. When we talk to each other, we anticipate a reply, i.e., some information, from the other person. Therefore, we might generate a sentence so as to maximize the information gain or expected reward. Communication always involves some kind of decision-making problems. Active perception and learning will become more and more important in the wide range of decision-making problems concerning robots.

### 7.3 *Compositionality and semantics*

The hierarchical structure placed on segmented words that are extracted on the basis of double articulation analysis enables us to generate various meaningful sentences. A syntax is a rule that produces a meaningful sequence of words. From the viewpoint of an emergent symbol system, an important problem can be phrased as “how can combined words have adequate meanings for an embodied and situated agent?” An important research topic concerns computational models for grounding semantic composition [150]. Classically, the principle of compositionality, which is also called Frege’s principle, has been widely recognized. This is a principle stating that the meaning of a complex expression is determined by the meanings of its constituent expressions and the syntactic rules used for combining them [151]. A bottom-up approach to the principle of compositionality also represents an important topic in symbol emergence in robotics.

The notion of combinatoriality has been studied in a constructive manner. Sugita et al. and Ogata et al. developed a robotic system and neural network architecture that can simultaneously learn sentences and behaviors [130, 152, 153]. For this, they used an RNNPB. It was shown that compositionality emerges on the network in a distributed manner. Tuci et al. also introduced neurodynamic models that deal with compositionality problems in language and behavior association learning, and the learning of goal-directed actions [154, 155]. Hinaut et al. applied reservoir computing, which is a kind of RNN, to make a robot acquire and produce grammatical constructions [156, 157]. Tani and Cangelosi et al. have presented a comprehensive survey of related studies [1, 2].

Recently, distributional models of semantics, including word2vec, have been given attention [158–160]. For example, Mikolov showed that the relationship between a country and a capital city can automatically be extracted from an unlabeled text dataset only by a training predictor, on the basis of a skip-gram and recurrent neural net language model (RNNLM). Le et al. proposed an unsupervised machine-learning method, called paragraph vector, that can estimate fixed-length feature representations from variable-length sections of texts, e.g., sentences, paragraphs, and documents [161]. Many preliminary studies exist concerning compositionality and semantics. Incorporating both embodied cognition and formal language structure must be important to construct a robot that can understand uttered sentences in the real world. @ The connection of such learning methods to a series of studies in SER also represents an important topic.

## 8. Conclusion

In this paper, we have provided an introduction to emergent symbol systems and surveyed the research field of SER. Semioticians sometimes call humans “Homosignificans,” which means meaning-makers [8]. Comprehending signs from natural or artificial environments and applying semiosis in the mind are human characteristics. When developing an autonomous robot that can engage in long-term communication and collaboration with people, the robot must adapt to

the human symbol system. To provide a philosophical framework for the diversity and dynamics of a symbol system, we introduce a concept—the emergent symbol system—that constitutes a basic assumption in SER. People can acquire language through physical interactions with their environment and semiotic communication with other people (see Figure 1). This phenomenon is comprehended as a type of assimilation process, in which a personal symbol system that is supported by an internal representation system becomes coupled with the emergent symbol system. To achieve such assimilation, the person must have the capability to learn the language in an unsupervised manner. For robots, this must be same. To achieve long-term interaction with people, a robot has to have the capability to learn the language in an unsupervised manner, so that the robot’s symbol system becomes coupled with the emergent symbol system that the target society has. Therefore, it is crucially important to computationally understand how humans can learn a symbol system and obtain semiotic skills through their autonomous mental development.

Many challenges have been investigated in relation to the construction of robotic systems and machine-learning methods that can obtain some parts of language through the embodied multimodal interaction with the environment. In order to understand human social interactions, and develop a robot that can smoothly communicate with human users, it is fundamentally important to understand symbol systems that change dynamically on the basis of the embodied cognition of participants in a constructive manner.

In this paper, we introduced the research field called SER. This represents a constructive approach towards an emergent symbol system. The emergent symbol system is socially self-organized through both semiotic and physical interactions with autonomous cognitive developmental agents, i.e., humans and developmental robots. Among the numerous fields connected with SER, we have described some specific topics in this paper, such as multimodal categorization, word discovery, and double articulation analysis. SER presents various future challenges involving acquiring lexicons, learning syntax, obtaining skills using metaphors, learning pragmatics, and being able to generate appropriate sentences given a particular context.

The majority of studies relating natural language processing and linguistics have only treated documents. However, we have to communicate and collaborate with other agents, including people and robots, in a real-world environment. The appropriateness of emergent symbol systems and robotic systems must be evaluated in relation to embodied cognition, context, and collaboration. Real-world collaborative tasks should be considered in order to evaluate them. In this context, some competitions including real-world human–robot interaction must be effective for facilitating researchers to study embodied cognitive systems and emergent symbol systems, and for evaluate appropriateness of the developed systems. RoboCup@Home is an obvious candidate [162–165].

However, the learning of all of the knowledge required for communication and collaboration through situated and embodied interactions requires very large costs and time. For further research, some pseudo-real-world environment will be required to increase the speed of our research. For example, Inamura et al. have developed the SIGVerse, a SocioIntelliGenesis simulator that enables human–robot interactions in a virtual world [166]. Cloud-based semiotic and physical interactions will also be important components of further studies in SER.

SER is still an emerging research field, but one that shows promise. Further research on SER will push long-term human–robot interaction forward, and provide a new understanding of human intelligence.

## Acknowledgements

This research was partially supported by a Grant-in-Aid for Young Scientists (B) 2012-2014 (24700233) and a Grant-in-Aid for Young Scientists (A) 2015-2019 (15H05319) funded by the Ministry of Education, Culture, Sports, Science, and Technology, Japan. This research was sup-

ported by CREST, JST.

## References

- [1] Cangelosi A, Metta G, Sagerer G, Nolfi S, Nehaniv C, Fischer K, Tani J, Belpaeme T, Sandini G, Nori F, Fadiga L, Wrede B, Rohlfing K, Tuci E, Dautenhahn K, Saunders J, Zeschel A. Integration of action and language knowledge: A roadmap for developmental robotics. *IEEE Transactions on Autonomous Mental Development*. 2010;2(3):167–195.
- [2] Tani J. Self-organization and compositionality in cognitive brains: A neurorobotics study. *Proceedings of the IEEE*. 2014;102(4):586–605.
- [3] Nakamura T, Nagai T, Iwahashi N. Grounding of word meanings in multimodal concepts using LDA. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. 2009. p. 3943–3948.
- [4] Newell A, Simon Ha. *Completer Science as Empirical Inquiry: Symbols and Search*. *Communications of the ACM*. 1976;19(3):113–126.
- [5] Newell A. Physical symbol systems. *Cognitive Science*. 1980;4:135–183.
- [6] Russell S, Norvig P. *Artificial intelligence: A modern approach*. 3rd ed. Pearson. 2009.
- [7] Inamura T, Toshima I, Tanie H, Nakamura Y. Embodied Symbol Emergence Based on Mimesis Theory. *The International Journal of Robotics Research*. 2004;23(4):363–377.
- [8] Chandler D. *Semiotics the Basics*. Routledge. 2002.
- [9] Brooks R. Elephants Don’t Play Chess. *Robotics and Autonomous Systems*. 1990;6:3–15.
- [10] Brooks R. Intelligence without representation. *Artificial Intelligence*. 1991;47(1-3):139–159.
- [11] Breazeal C. Emotion and sociable humanoid robots. *International Journal of Human Computer Studies*. 2003;59(1-2):119–155.
- [12] Breazeal CL. *Designing sociable robots*. MIT press. 2004.
- [13] Pfeifer R, Scheier C. *Understanding intelligence*. A Bradford Book. 2001.
- [14] Harnad S. The symbol grounding problem. *Physica D: Nonlinear Phenomena*. 1990;42(1):335–346.
- [15] Cangelosi A, Greco A, Harnad S. From robotic toil to symbolic theft: Grounding transfer from entry-level to higher-level categories. *Connection Science*. 2000;12(2):143–162.
- [16] Barsalou LW. Perceptual symbol systems. *Behavioral and Brain Sciences*. 1999;22(04):1–16.
- [17] Cangelosi A, Riga T. An embodied model for sensorimotor grounding and grounding transfer: experiments with epigenetic robots. *Cognitive science*. 2006;30(4):673–689.
- [18] Araki T, Nakamura T, Nagai T, Nagasaka S, Taniguchi T, Iwahashi N. Online learning of concepts and words using multimodal LDA and hierarchical Pitman-Yor Language Model. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. 2012. p. 1623–1630.
- [19] Sinapov J, Schenck C, Staley K, Sukhoy V, Stoytchev A. Grounding semantic categories in behavioral interactions: Experiments with 100 objects. *Robotics and Autonomous Systems*. 2014; 62(5):632–645.
- [20] Tellex S, Kollar T, Dickerson S. Approaching the Symbol Grounding Problem with Probabilistic Graphical Models. *AI magazine*. 2011;32(4):64–76.
- [21] Taddeo M, Floridi L. Solving the Symbol Grounding Problem: a Critical Review of Fifteen Years of Research. *Journal of Experimental and Theoretical Artificial Intelligence*. 2005;17:419–445.
- [22] Steels L. *The symbol grounding problem has been solved, so what’s next ? Symbols, Embodiment and Meaning* Oxford University Press, Oxford, UK. 2008;(2005):223–244.
- [23] Weng J. Symbolic Models and Emergent Models: A Review. *IEEE Transactions on Autonomous Mental Development*. 2012;4(1):29–53.
- [24] Cangelosi A, Schlesinger M. *Developmental Robotics*. The MIT press. 2015.
- [25] Asada M, Hosoda K, Kuniyoshi Y, Ishiguro H, Inui T, Yoshikawa Y, Ogino M, Yoshida C. Cognitive Developmental Robotics: A Survey. *IEEE Transactions on Autonomous Mental Development*. 2009; 1(1):12–34.
- [26] Tenenbaum JB, Kemp C, Griffiths TL, Goodman ND. How to grow a mind: statistics, structure, and abstraction. *Science*. 2011;331(6022):1279–1285.
- [27] Tagniguchi T. *Can we create a robot that communicate with human? –constructive approach towards symbol emergence system–*. NTT publishing Co. Ltd.. 2010. (in Japanese).
- [28] Tagniguchi T. *Symbol emergence in robotics –introduction to mechanism of intelligence–*. Kodansha.

2014. (in Japanese).
- [29] Eco U. *A Theory of Semiotics*. London: Indiana University Press. 1976.
  - [30] de Saussure F. *Course in General Linguistics* (trans. Roy Harris). London. 1983.
  - [31] Peirce CS. *Collected Writings*. Cambridge: Harvard University Press. 1931-1958.
  - [32] Von Uexküll J. A stroll through the worlds of animals and men: A picture book of invisible worlds. *Semiotica*. 1992;89(4):319–391.
  - [33] Sturrock J. *Structuralism*. London: Paladin. 1986.
  - [34] Piaget J. *Genetic epistemology*. W W Norton & Co Inc. 1971.
  - [35] Flavell JH. *The developmental psychology of Jean Piaget*. Literary Licensing, LLC. 2011.
  - [36] Polanyi M. *The Tacit Dimension*. The University of Chicago Press. 1966.
  - [37] Maturana HR, Varela FJ. *The tree of knowledge: The biological roots of human understanding*. revised ed. Shambhala. 1992.
  - [38] Krizhevsky A, Sutskever I, Hinton GE. ImageNet Classification with Deep Convolutional Neural Networks. In: *Advances in neural information processing systems (NIPS)*. 2012. p. 1–9. 1102.0183.
  - [39] Dahl GE, Yu D, Deng L, Acero A. Context-Dependent Pre-Trained Deep Neural Networks for Large-Vocabulary Speech Recognition. *IEEE Transactions on Audio, Speech, and Language Processing*. 2012;20(1):30–42.
  - [40] Bishop C. *Pattern recognition and machine learning (information science and statistics)*. Springer. 2010.
  - [41] Celikkanat H, Orhan G, Pugeault N, Guerin F, Erol S, Kalkan S. Learning and Using Context on a Humanoid Robot Using Latent Dirichlet Allocation. In: *Joint IEEE International Conferences on Development and Learning and Epigenetic Robotics (ICDL-Epirob)*. 2014. p. 201–207.
  - [42] Sinapov J, Stoytchev A. Object Category Recognition by a Humanoid Robot Using Behavior-Grounded Relational Learning. In: *IEEE International Conference on Robotics and Automation (ICRA)*. 2011. p. 184 – 190.
  - [43] Natale L, Metta G, Sandini G. Learning haptic representation of objects. In: *International Conference of Intelligent Manipulation and Grasping*. 2004.
  - [44] Ando Y, Nakamura T, Araki T, Nagai T. Formation of hierarchical object concept using hierarchical latent Dirichlet allocation. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. 2013. p. 2272–2279.
  - [45] Nakamura T, Nagai T, Iwahashi N. Multimodal object categorization by a robot. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. 2007. p. 2415–2420.
  - [46] Nakamura T, Nagai T, Iwahashi N. Bag of multimodal LDA models for concept formation. In: *IEEE International Conference on Robotics and Automation (ICRA)*. 2011. p. 6233–6238.
  - [47] Nakamura T, Nagai T, Iwahashi N. Multimodal categorization by hierarchical dirichlet process. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. 2011. p. 1520–1525.
  - [48] Nakamura T, Nagai T, Funakoshi K, Nagasaka S, Taniguchi T, Iwahashi N. Mutual Learning of an Object Concept and Language Model Based on MLDA and NPYLM. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. 2014. p. 600 – 607.
  - [49] Griffith S, Sinapov J, Sukhoy V, Stoytchev A. A behavior-grounded approach to forming object categories: Separating containers from noncontainers. *IEEE Transactions on Autonomous Mental Development*. 2012;4(1):54–69.
  - [50] Iwahashi N, Sugiura K, Taguchi R, Nagai T, Taniguchi T. Robots that learn to communicate: A developmental approach to personally and physically situated human-robot conversations. In: *Dialog with robots papers from the aaai fall symposium*. 2010. p. 38–43.
  - [51] Roy DK, Pentland AP. Learning words from sights and sounds: a computational model. *Cognitive Science*. 2002;26(1):113–146.
  - [52] Wenchi YEH, Barsalou LW. The situated nature of concepts. *American Journal of Psychology*. 2006;119(3):349–384.
  - [53] Mangin O, Oudeyer PY. Learning semantic components from subsymbolic multimodal perception. In: *Ieee 3rd joint international conference on development and learning and epigenetic robotics (icdl)*. 2013. p. 1 – 7.
  - [54] Lalle S, Ford Dominey P. Multi-modal convergence maps : From body schema and self-representation to mental imagery. In: *Adaptive behavior*. 2013.
  - [55] Ivaldi S, Nguyen SM, Lyubova N, Droniou A, Padois V, Filliat D, Oudeyer PY, Sigaud O. Object Learning Through Active Exploration. *IEEE Transactions on Autonomous Mental Development*. 2014;6(1):56–72.

- [56] Nakamura T, Ando Y, Nagai T, Kaneko M. Concept formation by robots using an infinite mixture of models. In: IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). 2015.
- [57] Blei DM, Ng AY, Jordan MI. Latent Dirichlet allocation. *The Journal of Machine Learning Research*. 2003;3(1):993–1022.
- [58] Griffiths TL, Steyvers M. Finding scientific topics. *Proceedings of the National Academy of Sciences of the United States of America (PNAS)*. 2004;101 (Suppl 1):5228–5235.
- [59] Araki T, Nakamura T, Nagai T, Funakoshi K, Nakano M, Iwahashi N. Autonomous acquisition of multimodal information for online object concept formation by robots. In: 2011 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). 2011. p. 1540–1547.
- [60] Teh Y, Jordan M, Beal M, Blei D. Hierarchical dirichlet processes. *Journal of the American Statistical Association*. 2006;101(476):1566–1581.
- [61] Sudderth EB, Torralba a, Freeman W, Willsky AS. Describing Visual Scenes using Transformed Dirichlet Processes. In: *Advances in Neural Information Processing Systems (NIPS)*. Vol. 18. 2005. p. 1297–1304.
- [62] Teh Y, Jordan M. Hierarchical Bayesian nonparametric models with applications. *Bayesian Nonparametrics*. 2009;:158.
- [63] Blei DM, Griffiths TL, Jordan MI. The nested Chinese restaurant process and Bayesian nonparametric inference of topic hierarchies. *Journal of the ACM (JACM)*. 2007;57(2):1–30. 0710.0845.
- [64] Nakamura T, Nagai T, Iwahashi N. Bag of multimodal hierarchical dirichlet processes: Model of complex conceptual structure for intelligent robots. In: IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). 2012. p. 3818–3823.
- [65] Tagniguchi T, Takano T, Yoshino R. Active perception for multimodal object category recognition using information gain. *IEEE Transactions on Autonomous Mental Development* (submitted). 2015; .
- [66] Ngiam J, Khosla A, Kim M, Nam J, Lee H, Ng AY. Multimodal Deep Learning. In: *Proceedings of The 28th International Conference on Machine Learning (ICML)*. 2011. p. 689–696.
- [67] Noda K, Arie H, Suga Y, Ogata T. Intersensory Causality Modeling Using Deep Neural Networks. In: *IEEE International Conference on Systems, Man, and Cybernetics (SMC)*. 2013. p. 1995–2000.
- [68] Le QV, Ranzato M, Monga R, Devin M, Chen K, Corrado GS, Dean J, Ng AY. Building high-level features using large scale unsupervised learning. In: *International conference in machine learning (ICML)*. 2011.
- [69] Bengio Y, Courville A, Vincent P. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. 2013;35(8):1798–1828.
- [70] Saffran JR, Newport EL, Aslin RN. Word Segmentation: The Role of Distributional Cues. *Journal of Memory and Language*. 1996;35(4):606–621.
- [71] Saffran JR, Aslin RN, Newport EL. Statistical learning by 8-month-old infants. *Science*. 1996; 274(5294):1926–1928.
- [72] Thiessen ED, Saffran JR. When cues collide: use of stress and statistical cues to word boundaries by 7- to 9-month-old infants. *Developmental psychology*. 2003;39(4):706–716.
- [73] Brent MR. An efficient, probabilistically sound algorithm for segmentation and word discovery. *Machine Learning*. 1999;34:71–105.
- [74] Venkataraman A. A statistical model for word discovery in transcribed speech. *Computational Linguistics*. 2001;27(3):351–372.
- [75] Goldwater S, Griffiths TL, Johnson M, Griffiths T. Contextual dependencies in unsupervised word segmentation. In: *Proceedings of the 21st International Conference on Computational Linguistics and the 44th annual meeting of the Association for Computational Linguistics*. 2006. p. 673–680.
- [76] Goldwater S, Griffiths TL, Johnson M. A Bayesian framework for word segmentation: exploring the effects of context. *Cognition*. 2009;112(1):21–54.
- [77] Mochihashi D, Yamada T, Ueda N. Bayesian unsupervised word segmentation with nested Pitman-Yor language modeling. In: *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP (ACL-IJCNLP)*. 2009. p. 100–108.
- [78] Johnson M, Goldwater S. Improving nonparameteric Bayesian inference: experiments on unsupervised word segmentation with adaptor grammars. In: *Annual conference of the north american chapter of the association for computational linguistics*. 2009. p. 317–325.
- [79] Chen M, Chang B, Pei W. A joint model for unsupervised Chinese word segmentation. In: *Conference on empirical methods in natural language processing (EMNLP)*. 2014. p. 854–863.



- [80] Magistry P. Unsupervised word segmentation : the case for Mandarin Chinese. In: Annual meeting of the association for computational linguistics. Vol. 2. 2012. p. 383–387.
- [81] Sakti S, Finch A, Isotani R, Kawai H, Nakamura S. Unsupervised determination of efficient Korean LVCSR units using a Bayesian Dirichlet process model. In: Ieee international conference on acoustics, speech and signal processing (ICASSP). 2011. p. 4664–4667.
- [82] Iwahashi N. Language acquisition through a human-robot interface by combining speech, visual, and behavioral information. *Information Sciences*. 2003;156:109–121.
- [83] Iwahashi N. Interactive learning of spoken words and their meanings through an audio-visual interface. *IEICE Transactions on Information and Systems*. 2008;(2):312–321.
- [84] Teh YW. A hierarchical Bayesian language model based on Pitman-Yor processes. In: International conference on computational linguistics and the annual meeting of the association for computational linguistics. 2006. p. 985–992.
- [85] Neubig G, Mimura M, Mori S, Kawahara T. Bayesian learning of a language model from continuous speech. *IEICE Transactions on Information and Systems*. 2012;E95-D(2):614–625.
- [86] Heymann J, Walter O. Unsupervised word segmentation from noisy input. In: Ieee workshop on automatic speech recognition and understanding (asru). 2013. p. 458–463.
- [87] Heymann J, Walter O, Haeb-umbach R, Raj B. Iterative bayesian word segmentation for unsupervised vocabulary discovery from phoneme lattices. In: IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). 2014. p. 4085–4089.
- [88] Elsner M, Goldwater S, Feldman N, Wood F. A joint learning model of word segmentation, lexical acquisition, and phonetic variability. In: Proceedings of the 2013 conference on empirical methods in natural language processing. Seattle, Washington, USA. 2013. p. 42–54.
- [89] Hinoshita W, Arie H, Tani J, Okuno HG, Ogata T. Emergence of hierarchical structure mirroring linguistic composition in a recurrent neural network. *Neural Networks*. 2011;24(4):311–320.
- [90] Kamper H, Jansen A, Goldwater S. Fully Unsupervised Small-Vocabulary Speech Recognition Using a Segmental Bayesian Model. In: Interspeech. 2015.
- [91] Brandl H, Wrede B, Joubin F, Goerick C. A self-referential childlike model to acquire phones, syllables and words from acoustic speech. In: IEEE International Conference on Development and Learning (ICDL). 2008. p. 31–36.
- [92] Lee Cy, Donnell TJO, Glass J. Unsupervised Lexicon Discovery from Acoustic Input. *Transactions of the Association for Computational Linguistics*. 2015;3:389–403.
- [93] Taniguchi T, Nakashima R, Nagasaka S. Nonparametric Bayesian double articulation analyzer for direct language acquisition from continuous speech signals. 2015. arXiv:1506.06646.
- [94] Taguchi R, Yamada Y, Hattori K, Umezaki T, Hoguro M. Learning place-names from spoken utterances and localization results by mobile robot. In: Interspeech. 2011. p. 1325–1328.
- [95] Taniguchi T, Nagasaka S. Double articulation analyzer for unsegmented human motion using Pitman-Yor language model and infinite hidden Markov model. In: IEEE SICE international symposium on system integration (SII). 2011. p. 250–255.
- [96] Takano W, Imagawa H, Kulić D, Nakamura Y. What do you expect from a robot that tells your future? The crystal ball. In: IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). IEEE. 2010. p. 1780–1785.
- [97] Takenaka K, Bando T, Nagasaka S, Taniguchi T, Hitomi K. Contextual scene segmentation of driving behavior based on double articulation analyzer. In: IEEE/RSJ international conference on intelligent robots and systems (IROS). 2012. p. 4847–4852.
- [98] Tani J, Nol S. Learning to perceive the world as articulated: An approach for hierarchical learning in sensory-motor systems. *Neural Networks*. 1999;12:1131–1141.
- [99] Rubin J, Richards. Boundaries of visual motion. AI Memo 835, Massachusetts Institute of Technology, Artificial Intelligence Laboratory. 1985;.
- [100] Fod A, Matarić M, Jenkins O. Automated derivation of primitives for movement classification. *Autonomous robots*. 2002;12(1):39–54.
- [101] Kawashima H, Matsuyama T. Multiphase learning for an interval-based hybrid dynamical system. *IEICE transactions on fundamentals of electronics, communications and computer sciences*. 2005; E88-A(11):3022–3035.
- [102] Barbič J, Safonova A, Pan J, Faloutsos C, Hodgins J, Pollard N. Segmenting motion capture data into distinct behaviors. In: Proceedings of graphics interface. 2004. p. 185–194.
- [103] Li Y, Wang T, Shum H. Motion texture: a two-level statistical model for character motion synthesis. In: Proceedings of the 29th annual conference on computer graphics and interactive techniques.

2002. p. 465–472.
- [104] Okada M, Nakamura D, Nakamura Y. Selforganizing Symbol Acquisition and Motion Generation based on Dynamics-based Information Processing System. In: Proc. of the second international workshop on man-machine symbiotic systems. 2004. p. 219–229.
- [105] Kadone H, Nakamura Y. Segmentation, memorization, recognition and abstraction of humanoid motions based on correlations and associative memory. In: IEEE-RAS International Conference on Humanoid Robotics. 2006. p. 1–6.
- [106] Chiappa S, Peters J. Movement extraction by detecting dynamics switches and repetitions. In: Advances in Neural Information Processing Systems (NIPS). 2010. p. 388–396.
- [107] Takano W, Imagawa H, Kulić D, Nakamura Y. Organization of behavioral knowledge from extraction of temporal-spatial features of human whole body motions. In: 3rd IEEE RAS and EMBS International Conference on Biomedical Robotics and Biomechatronics (BioRob). 2010. p. 52–57.
- [108] Kulić D, Ott C, Lee D, Ishikawa J, Nakamura Y. Incremental learning of full body motion primitives and their sequencing through human motion observation. *The International Journal of Robotics Research*. 2012;31(3):330–345.
- [109] Takano W, Nakamura Y. Integrating whole body motion primitives and natural language for humanoid robots. In: Ieee-ras international conference on humanoid robots (humanoids). 2008. p. 708–713.
- [110] Takano W, Nakamura Y. Incremental learning of integrated semiotics based on linguistic and behavioral symbols. In: Ieee/rsj international conference on intelligent robots and systems (iros). 2009. p. 2545–2550.
- [111] Fox EB, Sudderth EB, Jordan MI, Willsky AS. A sticky HDP-HMM with application to speaker diarization. *The Annals of Applied Statistics*. 2009;5(2A):1020–1056.
- [112] Taniguchi T, Nagasaka S, Hitomi K, Chandrasiri NP, Bando T. Semiotic prediction of driving behavior using unsupervised double articulation analyzer. In: IEEE intelligent vehicles symposium (IV). 2012. p. 849–854.
- [113] Taniguchia T, Nagasaka S, Hitomi K, Takenaka K, Bando T. Unsupervised hierarchical modeling of driving behavior and prediction of contextual changing points. *IEEE Transactions on Intelligent Transportation Systems*. 2014;16(4):1746 – 1760. in press.
- [114] Nagasaka S, Taniguchi T, Yamashita G, Hitomi K, Bando T. Finding meaningful robust chunks from driving behavior based on double articulation analyzer. In: IEEE/SICE intl symposium on system integration (SII). 2012. p. 535–540.
- [115] Bando T, Takenaka K, Nagasaka S, Taniguchi T. Unsupervised drive topic finding from driving behavioral data. In: IEEE intelligent vehicles symposium (IV). 2013. p. 177–182.
- [116] Bando T, Takenaka K, Nagasaka S, Taniguchi T. Automatic drive annotation via multimodal latent topic model. In: IEEE/RSJ international conference on intelligent robots and systems (IROS). 2013. p. 2744–2749.
- [117] Takenaka K, Bando T, Nagasaka S, Taniguchi T. Drive video summarization based on double articulation structure of driving behavior. In: ACM Multimedia. 2012. p. 1169–1172.
- [118] Taniguchi T, Nagasaka S, Hitomi K, Chandrasiri N, Bando T, Takenaka K. Sequence prediction of driving behavior using double articulation analyzer. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*. 2015;PP(99):1–1.
- [119] Mori M, Takenaka K, Bando T, Taniguchi T, Miyajima C, Takeda K. Automatic lane change extraction based on temporal patterns of symbolized driving behavioral data. In: IEEE Intelligent Vehicles Symposium (IV). 2015.
- [120] Wolpert DM, Kawato M. Multiple paired forward and inverse models for motor control. *Neural networks*. 1998;11(7):1317–1329.
- [121] Demiris Y, Khadhour B. Hierarchical attentive multiple models for execution and recognition of actions. *Robotics and Autonomous Systems*. 2006;54(5):361–369.
- [122] Taniguchi T, Sawaragi T. Assimilation and accommodation for self-organizational learning of autonomous robots: Proposal of dual-schemata model. In: IEEE International Symposium on Computational Intelligence in Robotics and Automation (CIRA). Vol. 1. 2003. p. 277–282.
- [123] Taniguchi T, Sawaragi T. Self-organization of inner symbols for chase: symbol organization and embodiment. In: IEEE International Conference on Systems, Man and Cybernetics (SMC). Vol. 2. 2004. p. 2073–2078.
- [124] Tagniguchi T, Sawaragi T. Incremental acquisition of multiple nonlinear forward models based on differentiation process of schema model. *Neural Networks*. 2008;21(1):13027.

- [125] Samejima K, Katagiri K, Doya K, Kawato M. Symbolization and imitation learning of motion sequence using competitive modules. *IEICE Transactions on Information and Systems*. 2002;85(1):90–100.
- [126] Wolpert D, Doya K, Kawato M. A unifying computational framework for motor control and social interaction. *Philosophical Transactions of the Royal Society B: Biological Sciences*. 2003;358:593–602.
- [127] Doya K, Samejima K, Katagiri Ki, Kawato M. Multiple model-based reinforcement learning. *Neural computation*. 2002;14(6):1347–1369.
- [128] Taniguchi T, Sawaragi T. Incremental acquisition of behaviors and signs based on a reinforcement learning schemata model and a spike timing-dependent plasticity network. *Advanced Robotics*. 2007;21(10):1177–1199.
- [129] Haruno M, Wolpert D, Kawato M. Hierarchical MOSAIC for movement generation. *International Congress Series*. 2003;1250:575–590.
- [130] Tani J, Ito M, Sugita Y. Self-organization of distributedly represented multiple behavior schemata in a mirror system: Reviews of robot experiments using RNNPB. *Neural Networks*. 2004;17(8-9):1273–1289.
- [131] Heinrich S, Magg S, Wermter S. Analysing the multiple timescale recurrent neural network for embodied language understanding. In: Koprinkova-Hristova P, Mladenov V, Kasabov NK, editors. *Artificial neural networks*. Vol. 4. Springer International Publishing. 2015. p. 149–174.
- [132] Murata S, Yamashita Y, Arie H, Ogata T, Tani J, Sugano S, Overview A. Generation of Sensory Reflex Behavior versus Intentional Proactive Behavior in Robot Learning of Cooperative Interactions with Others. In: *Joint IEEE International Conferences on Development and Learning and Epigenetic Robotics (ICDL-EpiRob)*. 2014. p. 242–248.
- [133] Roy N, Thrun S. Coastal Navigation with Mobile Robots. In: *Advances in Neural Processing Systems (NIPS)*. 1999.
- [134] Zuo X, Iwahashi N, Taguchi R, Funakoshi K, Nakano M, Matsuda S, Sugiura K, Oka N. Detecting robot-directed speech by situated understanding in object manipulation tasks. In: *IEEE International Workshop on Robot and Human Interactive Communication*. 2010. p. 608–613.
- [135] Sugiura K, Iwahashi N, Kawai H, Nakamura S. Situated Spoken Dialogue with Robots Using. *Advanced Robotics*. 2011;25:2207–2232.
- [136] Denzler J, Brown CM. Information Theoretic Sensor Data Selection for Active Object Recognition and State Estimation. *IEEE Transactions on pattern analysis and machine intelligence*. 2002; 24(2):1–13.
- [137] Krainin M, Curless B, Fox D. Autonomous generation of complete 3D object models using next best view manipulation planning. In: *IEEE International Conference on Robotics and Automation (ICRA)*. 2011. p. 5031–5037.
- [138] Eidenberger R, Scharinger J. Active perception and scene modeling by planning with probabilistic 6D object poses. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. 2010. p. 1036–1043.
- [139] van Hoof H, Kroemer O, Ben Amor H, Peters J. Maximally informative interaction learning for scene exploration. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. 2012. p. 5152–5158.
- [140] Rebguns A, Ford D, Fasel I. InfoMax Control for Acoustic Exploration of Objects by a Mobile Robot. In: *AAAI11 Workshop on Lifelong Learning*. 2011. p. 22–28.
- [141] Burgard W, Fox D, Thrun S. Active Mobile Robot Localization. In: *Proceedings of the Fourteenth International Joint Conference on Artificial Intelligence (IJCAI)*. 1997. p. 1346–1352.
- [142] Stachniss C, Grisetti G, Burgard W. Information Gain-based Exploration Using Rao-Blackwellized Particle Filters. In: *Robotics Science and Systems (RSS)*. 2005.
- [143] Gouko M, Kobayashi Y, Kim CH. Online Exploratory Behavior Acquisition of Mobile Robot Based on Reinforcement Learning. In: *Recent trends in applied artificial intelligence*. 2013. p. 272–281.
- [144] Saegusa R, Natale L, Metta G, Sandini G. Cognitive Robotics - Active Perception of the Self and Others -. In: *International Conference on Human System Interactions (HSI)*. 2011. p. 419–426.
- [145] Ji S, Carin L. Cost-Sensitive Feature Acquisition and Classification. *Pattern Recognition*. 2006; 40(5):1474–1485.
- [146] Tuci E, Massera G, Nolfi S. Active Categorical Perception of Object Shapes in a Simulated Anthropomorphic Robotic Arm. *IEEE Transactions on Evolutionary Computation*. 2010;14(6):885–899.
- [147] Schneider A, Sturm J, Stachniss C, Reisert M, Burkhardt H, Burgard W. Object identification with

- tactile sensors using bag-of-features. In: IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). 2009. p. 243–248.
- [148] Singh S, Chentanez N, Barto AG. Intrinsically Motivated Reinforcement Learning. In: Advances in Neural Information Processing Systems (NIPS). 2005. p. 1281–1288.
- [149] Schmidhuber J. Formal theory of creativity, fun, and intrinsic motivation (1990-2010). *IEEE Transactions on Autonomous Mental Development*. 2010;2(3):230–247.
- [150] Daoutis M, Mavridis N. Towards a Model for Grounding Semantic Composition. In: The 50th annual convention of the aisb. 2014.
- [151] Morris JF. The Principle of Semantic Compositionality. *Topoi*. 1994;13(1):11–24.
- [152] Sugita Y, Tani J. Learning Semantic Combinatoriality from the Interaction between Linguistic and Behavioral Processes. *Adaptive Behavior*. 2005;13(1):33–52.
- [153] Ogata T, Murase M, Tani J, Komatani K, Okuno HG. Two-way translation of compound sentences and arm motions by recurrent neural networks. In: IEEE International Conference on Intelligent Robots and Systems (IROS). 2007. p. 1858–1863.
- [154] Tuci E, Ferrauto T, Zeschel A, Massera G, Nolfi S. An experiment on behavior generalization and the emergence of linguistic compositionality in evolving robots. *IEEE Transactions on Autonomous Mental Development*. 2011;3:176–189.
- [155] Sandamirskaya Y, Zibner SKU, Schneegans S, Schöner G. Using Dynamic Field Theory to extend the embodiment stance toward higher cognition. *New Ideas in Psychology*. 2013;31(3):322–339.
- [156] Hinaut X, Petit M, Poiteau G, Dominey PF. Exploring the acquisition and production of grammatical constructions through human-robot interaction with echo state networks. *Frontiers in Neurobotics*. 2014;8:1–17.
- [157] Hinaut X, Lance F, Droin C, Petit M, Poiteau G, Dominey PF. Corticostriatal response selection in sentence production: Insights from neural network simulation with reservoir computing. *Brain and Language*. 2015;150:54–68.
- [158] Mikolov T, Sutskever I, Chen K, Corrado GS, Dean J. Distributed representations of words and phrases and their compositionality. In: Burges C, Bottou L, Welling M, Ghahramani Z, Weinberger K, editors. *Advances in Neural Information Processing Systems (NIPS)*. 2013. p. 3111–3119. Available from: <http://papers.nips.cc/paper/5021-distributed-representations-of-words-and-phrases-and-their-compos>
- [159] Mikolov T, Corrado G, Chen K, Dean J. Efficient Estimation of Word Representations in Vector Space. In: *International Conference on Learning Representations (ICLR)*. 2013. p. 1–12.
- [160] Mikolov T, Yih Wt, Zweig G. Linguistic regularities in continuous space word representations. In: *HLT-NAACL*. 2013. p. 746–751.
- [161] Le Q, Mikolov T. Distributed Representations of Sentences and Documents. In: *International Conference on Machine Learning (ICML)*. Vol. 32. 2014. p. 1188–1196. 1405.4053.
- [162] Stückler J, Behnke S. Integrating indoor mobility, object manipulation, and intuitive interaction for domestic service tasks. In: *IEEE-RAS International Conference on Humanoid Robots (HUMANOIDS)*. 2009. p. 506–513.
- [163] Iocchi L, Holz D, Ruiz-del Solar J, Sugiura K, van der Zant T. RoboCup@Home: Analysis and results of evolving competitions for domestic and service robots. *Artificial Intelligence*. 2015;1:1–24.
- [164] Nakamura T, Attamimi M, Sugiura K, Nagai T, Iwahashi N, Toda T, Okada H, Omori T. An Extended Mobile Manipulation Robot Learning Novel Objects. *Journal of Intelligent & Robotic Systems*. 2012;66(1):187–204.
- [165] Stückler J, Droschel D, Gräve K, Holz D, KläßJ, Schreiber M, Steffens R, Behnke S. Towards robust mobility, flexible object manipulation, and intuitive multimodal interaction for domestic service Robots. *Lecture Notes in Computer Science*. 2012;7416 LNCS:51–62.
- [166] Inamura T, Shibata T, Sena H, Hashimoto T, Kawai N, Miyashita T, Sakurai Y, Shimizu M, Otake M, Hosoda K, Umeda S, Inui K, Yoshikawa Y. Simulator platform that enables social interaction simulation - SIGVerse: SocioIntelliGenesis simulator. In: *IEEE/SICE International Symposium on System Integration (SII)*. 2010. p. 212–217.