

Max-Planck-Institut
für Mathematik
in den Naturwissenschaften
Leipzig

A Theory of Cheap Control in Embodied
Systems

(revised version: November 2014)

by

Guido Montúfar, Keyan Ghazi-Zahedi, and Nihat Ay

Preprint no.: 70

2014



A Theory of Cheap Control in Embodied Systems

Guido Montúfar¹, Keyan Ghazi-Zahedi¹, and Nihat Ay^{1,2,3}

¹Max Planck Institute for Mathematics in the Sciences, Inselstraße 22, 04103 Leipzig, Germany
²Department of Mathematics and Computer Science, Leipzig University, 04009 Leipzig, Germany
³Santa Fe Institute, 1399 Hyde Park Road, Santa Fe, NM 87501, USA

Abstract

We present a framework for designing cheap control architectures for embodied agents. Our derivation is guided by the classical problem of universal approximation, whereby we explore the possibility of exploiting the agent’s embodiment for a new and more efficient universal approximation of behaviors generated by sensorimotor control. This embodied universal approximation is compared with the classical non-embodied universal approximation. To exemplify our approach, we present a detailed quantitative case study for policy models defined in terms of conditional restricted Boltzmann machines. In contrast to non-embodied universal approximation, which requires an exponential number of parameters, in the embodied setting we are able to generate all possible behaviors with a drastically smaller model, thus obtaining cheap universal approximation. We test and corroborate the theory experimentally with a six-legged walking machine. The experiments show that the sufficient controller complexity predicted by our theory is tight, which means that the theory has direct practical implications.

Keywords: cheap design, embodiment, sensorimotor loop, universal approximation, conditional restricted Boltzmann machine

1 Introduction

In artificial intelligence, learning is one of the central fields of interest. Crucial for the success of any learning method is the complexity of the underlying model, e.g. a neural network. If the model is chosen too complex, the learning algorithm will likely require too much time and get stuck in a suboptimal solution. If it is chosen too simple, it might not be able to solve the problem at all. It is known from biological systems, that the exploitation of the body and environment allows a reduction of the neural system’s complexity (Pfeifer and Bongard 2006).

The goal of this article is to provide a framework that allows to determine the complexity of a control architecture in accordance with the *cheap design* principle from embodied artificial intelligence (Pfeifer and Bongard 2006). Cheap design in this context refers to the relatively low complexity of the brain or controller in comparison with the complexity of an observed behavior. A classical example is given by the Braitenberg vehicles (Braitenberg 1984), which are *Gedankenexperiments* designed to show how a seemingly complex behavior can result from very simple control structures. Braitenberg discusses several artificial creatures with simple wirings between sensors and actuators. He then describes how these systems produce a behavior that an external observer would classify as complex if the internal wirings were not revealed. Most interestingly, he then relates the wiring of his vehicles to various neural structures in the human brain. The idea of a simple

wiring that leads to complex behaviors is also discussed by Pfeifer and Bongard (2006), who present the walking behavior of an ant as an example. Without taking the embodiment and, in particular, the sensorimotor loop into account, the complex behavior (of a complex morphology) seems to require a complex control structure (Pfeifer and Bongard 2006; p. 79). A strong indication that cheap design is a common principle in biological systems is given by the fact that the human brain accounts for only 2% of the body mass but is responsible for 20% of the entire energy consumption (Clark and Sokoloff 1999), which is also remarkably constant (Sokoloff et al. 1955). Further support for cheap design as a common principle is given by a recent study on the brain sizes of migrating birds. It is known that migrating birds have a reduced brain size compared with their resident relatives. Sol et al. (2010) have studied various species and the affected brain regions and point out that the reduced brain sizes could be a direct result from the need to reduce energetic, metabolic and cognitive costs for migrating birds.

One way to achieve cheap design in this context is described as *compliance* in the embodied artificial intelligence community. A system is described as compliant, if it not only copes with the hard physical constraints it is subject to, but if it exploits them in order to minimize the required control effort. An illustrative example is the human walking behavior, which only needs to be actively controlled during the stance phase. The swing phase results mainly from the interaction of the physical properties of the leg with the environment (gravity). This is demonstrated by the Passive Dynamic Walker (McGeer 1990), which is a purely mechanical system that resembles the physical properties of human legs. The human walking behavior is emulated as a result of the interaction of the mechanical system with its environment (gravity and a slope). It is an impressive example of cheap design that requires no active control at all.

We are interested in quantifying to what extent a control structure can be reduced if the physical constraints are taken into account. Above, we referred to a system as cheaply designed, if it has a control structure of low complexity produces behaviors which an external observer would classify as complex. In this work, we are not concerned with the complexity of the behavior. Instead, we present an approach to determine the minimal complexity of a control structure that is able to produce a given set of desired behaviors (which can also be all theoretically possible behaviors) with a given morphology in a given environment. In other words, rather than comparing the complexities of the control structure and the behavior, we ask: what is the minimal brain complexity (or size) that can control all (desired) behaviors that are possible with the body and environment in which it is embedded?

There are various different complexity measures available in literature, of which the predictive information (Bialek and Tishby 1999), relevant information (Polani et al. 2006), and the Kolmogorov complexity (Schmidhuber 2009), are just a few examples. All these approaches have their specific strengths. However, they do not explicitly quantify how much the controller complexity can be reduced as a result of the agents embodiment, which is the focus of this work.

We follow a bottom-up, understanding by building approach (Brooks 1991a) to cognitive science, which is also known as behavior-based robotics (Brooks 1991b) and embodied artificial intelligence (Pfeifer and Bongard 2006). The core concept is that cognitive systems are considered as embedded and situated agents which cannot be understood if they are detached from the sensorimotor loop. This implicitly means that we assume sensor state sparsity and continuity of physical constraints. Consider the human retina as an example. We do not see random images but structured patterns and the sequence of these patterns is also highly dependent on our behavior. This behavior-dependent structuring of information is also known as *information self-structuring* and it has been

identified as one of the key principles of learning and development (Lungarella and Sporns 2005, Pfeifer et al. 2007). The second implication from the sensorimotor loop is continuity, e.g. natural systems are unable to teleport themselves from one place to another. Therefore, we can safely assume that the world around us will not be too different from the recent past and the recent future.

The sensorimotor loop (SML) (Klyubin et al. 2004, Ay and Zahedi 2014) is described by a type of partially observable Markov decision process (POMDP) where an embodied agent chooses actions based on noisy partial observations of its environment. An illustration of this causal structure is given in Figure 2. We aim at optimizing the design of policy models for controlling these processes. One aspect of the optimal design problem is addressed by working out the optimal complexity of the policy model. In particular, we are interested in the minimal number of units or parameters needed in order to obtain a network that can represent or approximate a desired set of behaviors within a given degree of accuracy. A first step towards resolving this problem is to address the minimal size of a universal approximator. In realistic scenarios, universal approximation is out of question, since it demands an enormous number of parameters – many more than actually needed. In this paper we reconsider the universal approximation problem by exploiting embodiment constraints and restrictions in the desired behavioral patterns.

We introduce the notions of *embodied behavior dimension* and *embodied universal approximation*, which quantify the effective dimension of a system that is subject to sensorimotor constraints (embodiment) and formalize the minimal control paradigm of cheap design in the context of the sensorimotor loop. We substantiate these ideas with theoretical results on the representational capabilities of conditional restricted Boltzmann machines (CRBMs) as policy models for embodied systems. CRBMs are artificial stochastic neural networks where the input and output units are connected bipartitely and undirected to a set of hidden units. Given the embodied behavior dimension, we derive bounds on the number of hidden units of CRBMs, that suffices to generate all possible behaviors by appropriate tuning of interaction weights and biases. In order to test our theory, we present an experimental study with a six-legged walking robot, and find a clear corroboration of our theorems. The experiments show that the sufficient controller complexity predicted by our theory is tight, which means that the theory has direct practical implications.

CRBMs are defined by clamping an *input* subset of the visible units of a Restricted Boltzmann machine (RBM) (Smolensky 1986, Freund and Haussler 1994). Conditional models of this kind have found a wide range of applications, e.g., in classification, collaborative filtering, and motion modeling (see Larochelle and Bengio 2008, Salakhutdinov et al. 2007, Sutskever and Hinton 2007, Taylor et al. 2007), and have proven useful as policy models in reinforcement learning settings (Sallans and Hinton 2004). These networks can be trained efficiently (Hinton 2002; 2012) and are well known in the context of learning representations and deep learning (see Bengio 2009). Although estimating the probability distributions represented by RBMs is hard (Long and Servedio 2010), approximate samples can be generated easily from a finite Gibbs sampling procedure. The theory and in particular the expressive power of RBM probability models has been studied in numerous papers (e.g., Le Roux and Bengio 2008, Montúfar and Ay 2011, Montúfar et al. 2011, Martens et al. 2013). Recently the representational power of CRBMs has been studied in detail (Montúfar et al. 2014). CRBMs can model non-trivial conditional distributions on high-dimensional input-output spaces using relatively few parameters, and their complexity can be adjusted by simply increasing or decreasing the number of hidden units. Hence we chose this model class for illustrating our discussion about the complexity of SML control problems.

This paper is organized as follows. Section 2 contains definitions around the SML. Section 3 presents the notions of embodied behavior dimension and embodied universal approximation, which we use to quantify and enforce dimensionality reduction. Section 4 contains our theoretical discussion on the representational power of CRBM models, comparing the non-embodied and the embodied settings, and pointing the role of the embodied behavior dimension. Section 5 puts the theory to the test in a robot control problem. Section 6 offers our conclusions and outlook. The Appendix contains technical proofs and details about possible generalizations of the discussion presented in the main part of the paper.

2 The Causal Structure of the Sensorimotor Loop

What is an embodied agent? In order to develop a theory of embodied agents that allows us to cast the core principles of the field of embodied intelligence into rigorous theoretical and quantitative statements, we need an appropriate formal model. Such a model should be general enough to be applicable to all kinds of embodied agents, including natural as well as artificial ones, and specific enough to capture the essential aspects of embodiment. How should such a model look like? First of all, obviously, an embodied agent has a body. This body is situated in an environment with which the agent can interact, thereby generating some behavior. In order to be useful, this behavior has to be guided or controlled by the agent’s brain or controller. Drawing the boundary between the brain on one side and the body, together with the environment, on the other side suggests a black box perspective of the brain. The brain receives sensor signals from and sends effector or actuator signals to the outside world. All it knows from the world is based on this closed loop of signal transmission. In other words, the world is a black box for the brain with which it interacts through sensing and acting. In particular, the boundary between the body and the environment is

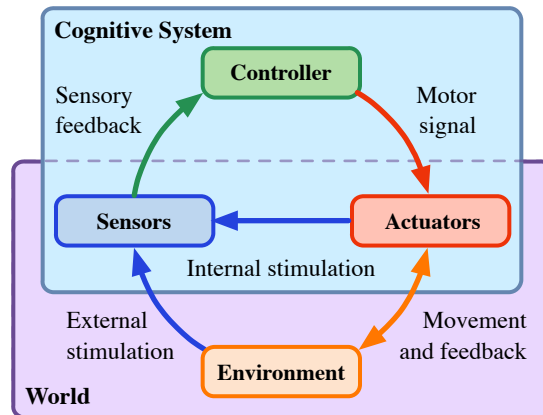


Figure 1: Sensorimotor loop

not directly “visible” for the brain. Both are parts of that black box and interact with the brain in an entangled way. Therefore, we consider them as being one entity, the outside world or simply the world. The brain is causally independent of the world, given the sensor signals, and the world is causally independent of the brain, given the actuator signals. This is the black box perspective.

Let us now develop a formal description of this sensorimotor loop. We denote the set of world states by \mathcal{W} . This set can be, for instance, the position of a robot in a static 3D environment. Information from the world is transmitted to the brain through sensors. Denoting the set of sensor states by \mathcal{S} , we can consider the sensor to be an information transmission channel from \mathcal{W} to \mathcal{S} as it is defined within information theory. Given a world state $w \in \mathcal{W}$, the response of the sensor can be characterized by a probability distribution of possible sensor states $s \in \mathcal{S}$ as result of w . For instance, if the sensor is noisy, then its response will not be uniquely determined. If the sensor is noiseless, that is, deterministic, then there will be only one sensor state as possible response to the world state w . In any case, the response of the sensor given w can be described in the following way: for a set S of sensor states we simply say how likely it is that the sensor will respond with a sensor state s that is contained in S . Formally, we can express this likelihood by a number $\beta(w; S)$ between zero and one, which is the probability of S given w . Collecting these numbers for all world states $w \in \mathcal{W}$ and sets S leads to the mathematical definition of a channel, also called a Markov kernel. It can be summarized as a map

$$\beta : \mathcal{W} \longrightarrow \Delta_{\mathcal{S}},$$

where $\Delta_{\mathcal{S}}$ denotes the set of probability distributions on the set \mathcal{S} of sensor states. The set of all such sensor channels is denoted by $\Delta_{\mathcal{W}}^{\mathcal{S}}$. Given a sensor channel β , there is another way to represent the probability distribution that is assigned to a world state w . Instead of providing a list $\beta(w; S)$ for all sets S of interest, we can restrict attention to infinitesimally small sets ds , leading to the notation $\beta(w; ds)$. In order to represent $\beta(w; S)$ in this notation, we have to integrate over all the infinitesimal ds in a set S , that is

$$\beta(w; S) = \int_S \beta(w; ds).$$

Whenever the base set \mathcal{S} is discrete, we simply replace ds by s and use $\beta(w; s)$ instead of $\beta(s; ds)$. Note that, as a Markov kernel, β has to satisfy various conditions. In order to provide a mathematically rigorous treatment, we assume that these conditions are satisfied. However, in order to improve the readability of the paper, we will not be very explicit with this (for the technical definitions see, e.g., Bauer 1996).

After having described in detail the mathematical model of a sensor, it is now straightforward to consider corresponding formalization of the other components of the sensorimotor loop. We continue with the notion of a policy. The agent can generate an effect in the world in terms of its actuators. Since we consider the body as part of the world, this can lead, for instance, to some body movement of the agent. In order guide this movement, it is beneficial for the agent to choose its actuator state based on the information about the world received through its sensors. Denoting the state set of the actuators by \mathcal{A} , we can again consider a channel from \mathcal{S} to \mathcal{A} as formal model of a policy, which we denote by π . Being more precise, with a sensor state s and a subset A of actuator states, $\pi(s; A)$ denotes the probability that the agent chooses an actuator state in A , given that its sensor state is s . Again, we have a Markov kernel

$$\pi : \mathcal{S} \longrightarrow \Delta_{\mathcal{A}},$$

where $\Delta_{\mathcal{A}}$ denotes the set of probability distributions on \mathcal{A} . We also use the notation $\pi(s; da)$ for an infinitesimal set da as we already introduced above in the context of β . Note that this definition

of a policy allows us to also consider a random choice of actions, so-called non-deterministic policies. The set of policies is denoted by $\Delta_{\mathcal{A}}^{\mathcal{S}}$.

Finally, we consider the change of the world state from w to w' in the context of an actuator state a as a channel which we denote by α . More precisely, given a world state w , an actuator state a , and a set W' of world states, $\alpha(w, a; W')$ denotes the probability that the actuator state a will generate a transition from w to a new world state that is in W' . As for the other channels, we use also in this case the notation $\alpha(w, a; dw')$ for infinitesimally small sets dw' . With the set $\Delta_{\mathcal{W}}$ of probability distributions on \mathcal{W} , we have

$$\alpha : \mathcal{W} \times \mathcal{A} \longrightarrow \Delta_{\mathcal{W}}.$$

We refer to α as *world channel* and denote the set of all world channels by $\Delta_{\mathcal{W}}^{\mathcal{W} \times \mathcal{A}}$.

We have defined three mechanisms that are involved in a (reactive) sensorimotor loop of an embodied agent. Clearly, the agent's embodiment poses constraints to this loop, which we attribute to the mechanisms β and α . The agent is equipped with these mechanisms, but they are both considered to be determined and not modifiable by the agent. On the other hand, the policy π can be modified by the agent in terms of learning processes. In order to describe the process of interaction of the agent with the world, we have to sequentially apply the individual mechanisms in the right order. Starting with an initial world state w^t at time t , first the sensor state s^t is generated in terms of the channel β . Then, based on the state of the sensor, an actuator state a^t is chosen according to the policy π . Finally, the world makes a transition, governed by α , from the state w to a new state w^{t+1} , which is influenced by the actuator state a^t of the agent. Altogether, this defines the combined mechanism

$$\mathbb{P}^{\pi}(w^t; ds^t, da^t, dw^{t+1}) := \beta(w^t; ds^t) \pi(s^t; da^t) \alpha(w^t, a^t; dw^{t+1}). \quad (1)$$

Note that we consider β and α fixed and therefore emphasise only the dependence on π . Now, with the new state w^{t+1} of the world, the three steps are iterated. This generates a process which is shown in Figure 2. Formally, the process is a probability distribution over trajectories that start with w^0 :

$$w^0, \quad s^0, a^0, w^1, \quad s^1, a^1, w^2, \quad s^2, a^2, w^3, \quad \dots, \quad s^{T-1}, a^{T-1}, w^T. \quad (2)$$

In order to describe this probability distribution, we have to iterate the mechanism (1) by multiplication:

$$\mathbb{P}^{\pi}(w^0; ds^0, da^0, dw^1, \dots, ds^{T-1}, da^{T-1}, dw^T) := \prod_{t=0}^{T-1} \mathbb{P}^{\pi}(w^t; ds^t, da^t, dw^{t+1}).$$

Now, what aspects of the sequence (2) represent the behavior of the agent? Let us consider, for instance, a walking behavior. It is given as a movement of the agent's body in physical space, which is completely determined by the world process. Remember that the body is part of the world. Clearly, the particular sequence of sensor and actuator states does not matter as long as they contribute to the generation of the same body movement. Therefore, we consider the world

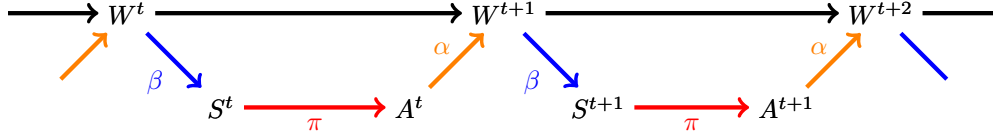


Figure 2: Causal structure of the reactive SML. The gray nuance groups the variables in one time step.

process w^t as the one in which behavior takes place and integrate out the other processes:

$$\begin{aligned} \mathbb{P}^\pi(w^0; dw^1, \dots, dw^T) &= \underbrace{\int_{\mathcal{S}} \int_{\mathcal{S}} \dots \int_{\mathcal{S}} \int_{\mathcal{S}}}_{T \text{ times}} \mathbb{P}^\pi(w^0; ds^0, da^0, dw^1, \dots, ds^{T-1}, da^{T-1}, dw^T) \\ &= \prod_{t=0}^{T-1} \mathbb{P}^\pi(w^t; dw^{t+1}). \end{aligned} \quad (3)$$

One can show that, with weak assumptions, the limit for $T \rightarrow \infty$ exists, so that we can write

$$\mathbb{P}^\pi(w^0; dw^1, dw^2, \dots),$$

which is a Markov kernel from an initial world state w^0 to the space of all infinite future sequences w^1, w^2, \dots . We denote the set of these Markov kernels by $\Delta_{\mathcal{W}\infty}^\pi$. This allows us to formalize the map that assigns to each policy the corresponding behavior:

$$\psi_\infty : \Delta_{\mathcal{S}}^\pi \longrightarrow \Delta_{\mathcal{W}\infty}^\pi, \quad \pi \longmapsto \mathbb{P}^\pi(w^0; dw^1, dw^2, \dots). \quad (4)$$

We refer to this map as the *policy-behavior map*. Two policies π_1 and π_2 will be considered equivalent, if they generate the same behavior, that is,

$$\psi_\infty(\pi_1) = \psi_\infty(\pi_2). \quad (5)$$

We argue that embodiment constraints render many equivalent policies. We can exploit this fact in order to design a concise control architecture. This will lead to a quantitative treatment of the notion of cheap design within the field of embodied intelligence. Let us treat this systems design problem in a more rigorous way. As we pointed out, the agent is equipped with the mechanisms β and α which constitute the embodiment of the agent. In a biological system these mechanisms will change due to developmental processes. However, we want to restrict our attention to the learning processes and disentangle them from developmental processes by assuming that the latter ones have already converged and therefore consider them as fixed. Learning refers to a process in which the policy is changing in time. Clearly, in order to model this change the agent has to be equipped with a family of possible policies, which we denote by \mathcal{M} , and refer to as *policy model*. For instance, we can consider neural networks as policy models that are parametrized by synaptic weights and threshold values for the individual neurons. Changing the weights and the thresholds will lead to a change of the policy (although there may be degeneracies, in general). In any case, going through all the possible parameter values will generate a set \mathcal{M} of policies with which the agent is equipped for its behavior.

We argue that, if the embodiment constraints lead to many equivalent policies, then it is possible to find a concise model \mathcal{M} that is capable of generating all behaviors. More precisely, we will require from a model

$$\psi_\infty(\mathcal{M}) = \psi_\infty(\Delta_{\mathcal{A}}^{\mathcal{S}}),$$

or, more precisely, a slight modification by taking limit points of \mathcal{M} into account. We refer to this property of \mathcal{M} as being an *embodied universal approximator*. In order to highlight the exploitation of embodiment constraints for cheap design, we compare this kind of universal approximation to the standard notion of universal approximation, which we refer to as *non-embodied universal approximation*.

3 Cheap Representation of Embodied Behaviors

Intuitively it is clear that the embodiment constraints cause restrictions in the set of behaviors that an agent can realize. For example, inertia restricts the pace at which an embodied system can change its direction of motion (imagine a train switching the traveling direction instantaneously). In turn, not all world-state transitions may be possible in a single time step, regardless of what the policy specifies as a desirable action to take. These restrictions create a bottleneck between the set of policies on the one side and the set of possible behaviors on the other. The consequence is that, generically, infinitely many policies parametrize the same behavior. If we understand the way in which different policies are mapped to the same, or to different, behaviors, then we can parametrize all the behaviors that can possibly emerge in the SML by a low-dimensional (or low-complexity) set of policies. We develop the necessary tools in this section. For clarity we will focus on the reactive SML with finite sensor and actuator state spaces but allowing the possibility of a continuous world state. In particular we will use $\beta(w; s)$ instead of $\beta(w; ds)$ and $\pi(s; a)$ instead of $\pi(s; da)$. Possible generalizations of these settings are discussed in Appendix C.

The condition (5) is clearly the same as

$$\mathbb{P}^{\pi_1}(w^0; dw^1, \dots, dw^T) = \mathbb{P}^{\pi_2}(w^0; dw^1, \dots, dw^T), \quad \text{for all } T = 1, 2, \dots,$$

and with equation (3), this is satisfied if and only if

$$\mathbb{P}^{\pi_1}(w; dw') = \mathbb{P}^{\pi_2}(w; dw').$$

Therefore, the mechanism $\mathbb{P}^\pi(w; dw')$ will play an important role in our analysis, and we consider the one-step formulation of the policy-behavior map:

$$\psi : \Delta_{\mathcal{A}}^{\mathcal{S}} \longrightarrow \Delta_{\mathcal{W}}^{\mathcal{W}}, \quad \pi \longmapsto \mathbb{P}^\pi(w; dw'), \quad (6)$$

where

$$\mathbb{P}^\pi(w; dw') = \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \beta(w; s) \pi(s; a) \alpha(a, w; dw'). \quad (7)$$

This is an affine map from the convex set $\Delta_{\mathcal{A}}^{\mathcal{S}}$ to the convex set $\Delta_{\mathcal{W}}^{\mathcal{W}}$. The image of this map represents the set of all possible behaviors that the SML can generate. We denote this set by Beh

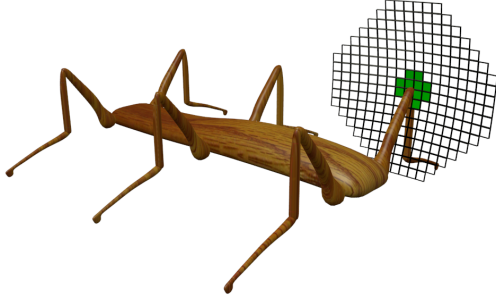


Figure 3: Locality of world-state transitions. This hexapod will be used in the experimental evaluation of our theory in Section 5.

and refer to its dimension the *embodied behavior dimension* $d = \dim(\psi(\Delta_{\mathcal{A}}^{\mathcal{S}}))$. This is given by the number linearly independent vectors

$$\beta(w; s)(\alpha(w, a_0; dw') - \alpha(w, a; dw')) \quad (8)$$

in $\text{aff}(\psi(\Delta_{\mathcal{A}}^{\mathcal{S}}))$, where (s, a) are sensor-actuator states with $a \neq a_0$, for some fixed $a_0 \in \mathcal{A}$. These vectors are namely the images of a basis of $\text{aff}(\Delta_{\mathcal{A}}^{\mathcal{S}})$. See Appendix A for more details about this. This allows us to formulate a simple upper bound for the embodied behavior dimension d in terms of the image dimension of the maps β and α . Writing rank for the image dimension and regarding β as a linear map $\mathbb{R}^{\mathcal{S}} \rightarrow \mathbb{R}^{\mathcal{W}}$ (operating on the columns of π) and α as an affine map $\Delta_{\mathcal{A}} \rightarrow \Delta_{\mathcal{W}}^{\mathcal{W}}$ (operating on the rows of π), we have

$$d = \text{rank}(\psi) \leq \text{rank}(\beta) \cdot \text{rank}(\alpha). \quad (9)$$

For example, if \mathcal{W} is finite, $\text{rank}(\beta)$ is the rank of the matrix with entries $(\beta(w; s))_{w \in \mathcal{W}, s \in \mathcal{S}}$ and $\text{rank}(\alpha)$ is the rank of the matrix with entries $(\alpha(w, a_0; w') - \alpha(w, a; w'))_{a \in \mathcal{A}, (w, w') \in \mathcal{W} \times \mathcal{W}}$, for some fixed $a_0 \in \mathcal{A}$. Even though the upper bound (9) does not necessarily hold with equality (the rank of the combined map ψ , where β and α share the same w , can be significantly smaller than the product of the individual ranks), it gives us a clear picture of how the embodiment constraints, represented by β and α , can lead to an embodied behavior dimension d that is much smaller than $|\mathcal{S}|(|\mathcal{A}| - 1)$, the dimension of the set of all policies $\Delta_{\mathcal{A}}^{\mathcal{S}}$.

In an embodied system the sensors are usually insensitive to a large number of variations of the world-state w and, therefore, β is piece-wise constant with respect to w . Also, the sensors implement a certain degree of redundancy, meaning that, for each w , the probability distribution $\beta(w; \cdot) \in \Delta_{\mathcal{S}}$ has certain types of symmetries. This means that $\text{rank}(\beta)$ is much smaller than $|\mathcal{S}|$ (the maximum theoretically possible rank). In the case of α , usually several actions a produce the same world-state transition, such that, for any fixed world state w , $\alpha(w, \cdot; \cdot)$ is piece-wise constant with respect to a . Furthermore, for any given w , only very few states $w' \in \mathcal{W}$ are possible at the next time step, regardless of a , such that $\alpha(w, a; \cdot)$ assigns positive probability only to a very small subset of \mathcal{W} . This means that $\text{rank}(\alpha)$ is usually much smaller than $(|\mathcal{A}| - 1)$ (the maximum theoretically possible rank). An example for this kind of constraints on α is a robot's knee, which in a time step can only be moved to adjacent positions, as the one shown in Figure 3.

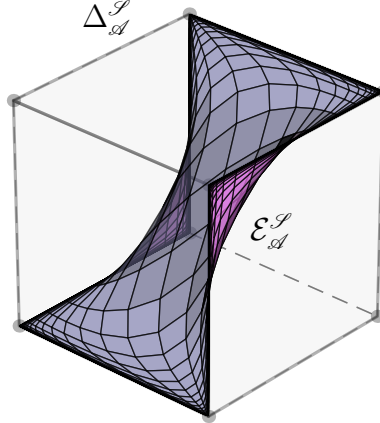


Figure 4: Illustration of the exponential family of policies described in Example 1. Here $|\mathcal{S}| = 3$ and $|\mathcal{A}| = 2$, such that the polytope $\Delta_{\mathcal{A}}^{\mathcal{S}}$ is the three-dimensional cube of 3×2 row stochastic matrices. The curved surface is the exponential family $\mathcal{E}_{\mathcal{A}}^{\mathcal{S}}$ from the example for some specific β and α with embodied behavior dimension $d = 2$. In the SML defined by β and α , the policies from $\mathcal{E}_{\mathcal{A}}^{\mathcal{S}}$ generate the same behaviors as the set of all policies $\Delta_{\mathcal{A}}^{\mathcal{S}}$.

So far we have discussed the embodied behavior dimension of an embodied system $d = \dim(\text{Beh})$ and why this can be much smaller than the dimension of the policy space $\dim(\Delta_{\mathcal{A}}^{\mathcal{S}}) = |\mathcal{S}|(|\mathcal{A}|-1)$. Since the policy-behavior map ψ is affine, for any generic behavior that can possibly emerge in the SML, there is a $(|\mathcal{S}|(|\mathcal{A}|-1) - d)$ -dimensional set (in fact a polytope) of equivalent policies generating that same behavior. By selecting representatives from each set of equivalent policies, we can define low-dimensional policy models which are just as expressive as the much higher dimensional set $\Delta_{\mathcal{A}}^{\mathcal{S}}$ of all possible policies, in terms of the representable behaviors. The following example shows that it is possible to define a smooth manifold of policies which translate in a one-to-one fashion to the set of all possible behaviors in the SML.

Example 1. Consider the matrix $E \in \mathbb{R}^{d \times (\mathcal{S} \times \mathcal{A})}$ that represents the policy-behavior map ψ with respect to some basis. Then the exponential family $\mathcal{E}_{\mathcal{A}}^{\mathcal{S}}$ of policies defined by

$$\pi_{\theta}(s; a) = \frac{\exp(\theta^{\top} E(s, a))}{\sum_{a' \in \mathcal{A}} \exp(\theta^{\top} E(s, a'))}, \quad \text{for all } a \in \mathcal{A} \text{ and } s \in \mathcal{S}, \quad \text{for all } \theta \in \mathbb{R}^d, \quad (10)$$

is an embodied universal approximator of dimension d . In fact, each behavior from the set Beh is realized by exactly one policy from the set $\overline{\mathcal{E}_{\mathcal{A}}^{\mathcal{S}}}$. See Figure 4 for an illustration and Appendix A for technical details.

The previous discussion shows that the set of behaviors that can possibly emerge in the SML usually has a much lower dimension than the set of all possible policies. Furthermore, it shows that it is possible to construct low-dimensional embodied universal approximators. Nonetheless, among all behaviors that are possible in the SML, we can expect that only a smaller subset $\mathcal{B} \subset \text{Beh}$ is actually relevant to the agent. For instance, among all locomotion gaits that an agent could possibly realize with its body in a given environment, we can expect that it will only utilize those which

are most successful (with respect to different tasks). The embodied behavior dimension can be directly generalized in order to capture such behavioral restrictions, in addition to the embodiment constraints α and β . We note that specifying a set \mathcal{B} of potentially interesting behaviors is in general not an easy task. For instance, one could be interested in behaviors that maximize the expected reward of an agent (besides from implicit definitions, say as the maximizers of a given objective function). As the main focus of this paper is not to find a good definition of interesting or natural behaviors, here we will consider a very simple but general classification of behaviors, in terms of their supports.

In relation to the information self-structuring mentioned in the introduction, for specific behaviors usually only a relatively small subset $\mathcal{S} \subseteq \mathcal{S}$ of sensor values emerges. In such situations, the policy only needs to be specified for the sensor values in the subset \mathcal{S} . Consider a set of behaviors \mathcal{B} that take place within a restricted set of world states $\mathcal{W} \subseteq \mathcal{W}$ and consider the set of sensor states that can be possibly measured from these world states, $\mathcal{S} := \{s \in \mathcal{S} : s \in \text{supp}(\beta(w; \cdot)) \text{ for some } w \in \mathcal{W}\}$. For the world states \mathcal{W} , the measurement by β always produces sensor values in \mathcal{S} , and the policy for states not in \mathcal{S} does not play any role. We denote by $\psi^{\mathcal{S}}$ the restriction of the policy-behavior map to $w \in \mathcal{W}$. This is given by the natural restriction of the kernels β and α to the domain $\mathcal{W} \subseteq \mathcal{W}$. In this case, the embodied behavior dimension is given by $d^{\mathcal{S}} := \dim(\psi^{\mathcal{S}}(\Delta_{\mathcal{A}}^{\mathcal{S}}))$. We will denote a model $\mathcal{M} \subseteq \Delta_{\mathcal{A}}^{\mathcal{S}}$ an embodied universal approximator on \mathcal{S} if $\psi^{\mathcal{S}}(\overline{\mathcal{M}}) = \psi^{\mathcal{S}}(\Delta_{\mathcal{A}}^{\mathcal{S}})$. This definition means that the model is powerful enough to control any behavior on \mathcal{W} just as well as the entire policy polytope. Given any set of behaviors \mathcal{B} , e.g., as the one described above, we are interested in the following problem.

Problem 1. For a given set of possible behaviors $\mathcal{B} \subseteq \text{Beh} = \psi(\Delta_{\mathcal{A}}^{\mathcal{S}})$ and a class of policy models \mathfrak{M} , what is the smallest model $\mathcal{M} \in \mathfrak{M}$ that can generate all these behaviors, such that $\mathcal{B} \subseteq \psi(\overline{\mathcal{M}})$?

Of particular interest are classes of policy models \mathfrak{M} defined in terms of neural networks. In Section 4 we will consider a class of policy models defined in terms of CRBMs. The following result gives us a simple and powerful combinatorial tool for addressing Problem 1.

Lemma 1. *Any model $\mathcal{M} \subseteq \Delta_{\mathcal{A}}^{\mathcal{S}}$ with the following property is an embodied universal approximator on \mathcal{S} : for every policy $\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}$ whose \mathcal{S} -rows have a total of $|\mathcal{S}| + d^{\mathcal{S}}$ or less non-zero entries, there exists a policy $\pi^* \in \overline{\mathcal{M}}$ with $\pi(s; \cdot) = \pi^*(s; \cdot)$ for all $s \in \mathcal{S}$.*

This lemma states that for universal approximation of embodied behaviors it suffices to approximate the policies which, for a relevant set of sensor values, assign positive probability only to a limited number of actions. The number of actions is determined by the embodied behavior dimension.

It is worthwhile mentioning that Example 1 and Lemma 1 describe two complementary types of universal approximators of embodied behaviors. The first type, described in the example, is composed of maximum entropy policies, whereas the second type, described in the lemma, is composed of minimum entropy policies. If we consider the set of equivalent policies that map to a given behavior, Example 1 selects the one with the most random state-action assignments that are possible for generating that behavior. On the other hand, Lemma 1 selects the ones with the most deterministic state-action assignments that are possible for generating that behavior. Geometrically, the set of equivalent policies of a given behavior is the convex hull of the minimum entropy policies, with the maximum entropy policy lying in the center. The exponential family has nice geometric properties, but it is very specific to the kernels β and α , which define the sufficient statistics. The set

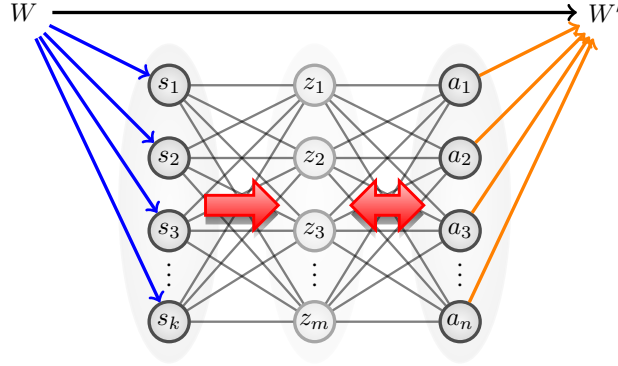


Figure 5: Illustration of a CRBM policy in the SML.

described in Lemma 1 can also be considered as a policy model. It offers several advantages that we will exploit later on. First, it has a very simple combinatorial description. Second, it only depends on the embodied behavior dimension d , irrespective of the specific kernels β and α (which are not directly accessible to the agent). Third, it selects policies with the minimum possible number of positive-probability actions, which seems natural from a concise controller.

4 A Case Study with Conditional Restricted Boltzmann Machines

4.1 Definitions

A Boltzmann machine (BM) is an undirected stochastic network with binary units, some of which may be hidden. It defines probabilities for the joint states of its visible units, given by the relative frequencies at which these states are observed, asymptotically, depending on the network parameters (interaction weights and biases). The probability of each joint state $x = (x_V, x_H)$ of the network's visible and hidden units can be described by the Gibbs-Boltzmann distribution $p(x) = \frac{1}{Z} \exp(-\mathcal{H}(x))$ with energy function $\mathcal{H}(x) = \sum_{i,j} x_i W_{ij} x_j + \sum_i b_i x_i$ and normalization partition function $Z(W, b) = \sum_{x'} \exp(-\mathcal{H}(x'))$. The probabilities of the visible states are given by marginalizing out the states of the hidden units, $p(x_V) = \sum_{x_H} p(x_V, x_H)$.

An RBM is a BM with the restriction that there are no interactions between the visible units nor between the hidden units, such that $W_{ij} \neq 0$ only when unit i is visible and j hidden. As any multivariate model of probability distributions, RBMs define models of conditional distributions, given by clamping the state of some of the visible units:

Definition 1. The conditional restricted Boltzmann machine model with k input, n output, and m hidden units, denoted $\text{RBM}_{n,m}^k$, is the set of all conditional distributions in $\Delta_{\mathcal{A}}^{\mathcal{S}}$, $\mathcal{S} = \{0, 1\}^k$, $\mathcal{A} = \{0, 1\}^n$, that can be written as

$$p(x|y) = \frac{1}{Z(W, b, Vy + c)} \sum_{z \in \{0,1\}^m} \exp(z^\top Wx + z^\top Vy + b^\top x + c^\top z), \quad \forall x \in \{0, 1\}^n, y \in \{0, 1\}^k,$$

where $Z(W, b, Vy + c)$ normalizes the probability distribution $p(\cdot|y) \in \Delta_{\mathcal{A}}$, for each $y \in \{0, 1\}^k$. Here, y , x , and z are state vectors of the input, output, and hidden units, respectively. Furthermore,

$V \in \mathbb{R}^{m \times k}$ is a matrix of interaction weights between hidden and input units, $W \in \mathbb{R}^{m \times n}$ is a matrix of interaction weights between hidden and output units, $c \in \mathbb{R}^m$ is a vector of biases for the hidden units, and $b \in \mathbb{R}^n$ is a vector of biases for the output units. Here \top denotes vector transposition.

The model $\text{RBM}_{n,m}^k$ has $mk + mn + m + n$ parameters (the interaction weights and biases). A bias term $a^\top y$ for the input units does not appear in the definition, as it would cancel out with the normalization function $Z(W, b, Vy + c)$. When there are no input units, i.e., $k = 0$, the conditional probability model $\text{RBM}_{n,m}^k$ reduces to the restricted Boltzmann machine probability model with n visible and m hidden units, which we denote by $\text{RBM}_{n,m}$.

4.2 Non-Embodied Universal Approximation

In this section we ask for the minimal number of hidden units m for which the model $\text{RBM}_{n,m}^k$ can approximate every conditional distribution from the set $\Delta_{\mathcal{A}}^{\mathcal{S}}$ with $\mathcal{S} = \{0, 1\}^k$ and $\mathcal{A} = \{0, 1\}^n$, denoted Δ_n^k , arbitrarily well. Later we will contrast this non-embodied universal approximation with the embodied case.

Note that each conditional distribution $p(x|y)$ can be identified with the set of joint distributions of the form $r(x, y) = q(y)p(x|y)$, with strictly positive marginals $q(y)$. By fixing $q(y)$ equal to the uniform distribution over \mathcal{S} , we obtain an identification of $\Delta_{\mathcal{A}}^{\mathcal{S}}$ with $\frac{1}{|\mathcal{S}|}\Delta_{\mathcal{A}}^{\mathcal{S}} \subseteq \Delta_{\mathcal{S} \times \mathcal{A}}$. In particular, we have that universal approximators of joint probability distributions define universal approximators of conditional distributions. This observation allows us to translate results on the representational power of RBMs to corresponding results for CRBMs. For example, we know that $\text{RBM}_{k+n,m}$ is a universal approximator of probability distributions on $\{0, 1\}^{k+n}$ whenever $m \geq \frac{1}{2}2^{k+n} - 1$ (see Montúfar and Ay 2011), and therefore:

Proposition 1. *The model $\text{RBM}_{n,m}^k$ can approximate every conditional distribution from Δ_n^k arbitrarily well whenever $m \geq \frac{1}{2}2^{k+n} - 1$.*

Now, since conditional models do not need to model the input-state distributions, in principle it is possible that $\text{RBM}_{n,m}^k$ is a universal approximator of conditional distributions even if $\text{RBM}_{n+k,m}$ is not a universal approximator of probability distributions. Therefore, we also consider a result by Montúfar et al. (2014), which improves Proposition 1 and does not follow from corresponding results for RBM probability models:

Theorem 1. *The model $\text{RBM}_{n,m}^k$ can approximate every conditional distribution from Δ_n^k arbitrarily well whenever $m \geq \frac{1}{2}2^k(2^n - 1) = \frac{1}{2}2^{k+n} - \frac{1}{2}2^k$.*

Theorem 1 represents a substantial improvement of Proposition 1 in that it reflects the structure of Δ_n^k as a 2^k -fold Cartesian product of the $(2^n - 1)$ -dimensional probability simplex Δ_n , in contrast to the proposition's bound, which rather reflects the structure of the $(2^{k+n} - 1)$ -dimensional joint probability simplex Δ_{k+n} . The full statement of the theorem is quite technical, and thus we refer the interested reader to (Montúfar et al. 2014). At this point let it suffice to say that some terms appearing in the bound on m decrease with increasing k , such that approximately the prefactor $\frac{1}{2}$ decreases to $\frac{1}{4}$ when k is large enough.

As expected, the asymptotic behavior of this result is exponential in the number of input and output units. We believe that the result is reasonably tight, although some improvements may still

be possible. A crude lower bound can be obtained by comparing the number of parameters with the dimension of the policy polytope (see details in Montúfar et al. 2014):

Proposition 2. *If the model $\text{RBM}_{n,m}^k$ can approximate every policy from Δ_n^k arbitrarily well, then necessarily $m \geq \frac{1}{(n+k+1)}(2^k(2^n - 1) - n)$.*

4.3 Embodied Universal Approximation

By Theorem 1, we can achieve embodied universal approximation by approximating only policies with a limited number of non-zero entries. Furthermore, as mentioned earlier, if we only care about a subset $\mathcal{S} \subseteq \mathcal{S}$ of sensor values, we can restrict the policy space $\Delta_{\mathcal{A}}^{\mathcal{S}}$ to $\Delta_{\mathcal{S}}^{\mathcal{S}}$. This means that the number of policy entries that we need to model is given by the number of interesting sensor values plus the corresponding embodied behavior dimension. On the other hand, we can use each hidden unit of a CRBM to model each relevant non-zero entry of the policy.

Theorem 2. *The model $\text{RBM}_{n,m}^k$ is an embodied universal approximator on \mathcal{S} whenever $m \geq |\mathcal{S}| + d^{\mathcal{S}} - 1$.*

Proof. We use Lemma 1. The joint probability model $\text{RBM}_{k+n,m}$ can approximate any probability distribution with support of cardinality $m + 1$ arbitrarily well (see Montúfar and Ay 2011). Hence, with $m \geq |\mathcal{S}| + d^{\mathcal{S}} - 1$, $\text{RBM}_{n+k,m}$ can approximate any joint distribution with $|\mathcal{S}| + d^{\mathcal{S}}$ non-zero entries within the set $\{(s, a) : s \in \mathcal{S}, a \in \mathcal{A}\}$ arbitrarily well. These joint distributions represent a set \mathcal{M} of conditional distributions of a given s such that, for any policy π whose \mathcal{S} -rows have a total of $|\mathcal{S}| + d^{\mathcal{S}}$ or less non-zero entries, there is a conditional distribution $\pi^* \in \mathcal{M}$ with $\pi^*(\cdot|s) = \pi(\cdot|s)$ for all $s \in \mathcal{S}$. \square

This theorem gives a bound for the number of hidden units of CRBMs that suffices to obtain embodied universal approximation. The bound depends on the embodiment and behavioral constraints of the system, captured in the embodied behavior dimension $d^{\mathcal{S}}$. In general, this bound will be much smaller than the exponential bound from Theorem 1. We will test this result in the context of particular behavioral constraint on a hexapod in the next section.

5 Experiments with a Hexapod

In the previous sections we have derived a theoretical bound for the complexity of a CRBM based policy. In this section, we want to evaluate that bound experimentally. For this purpose, we chose a six-legged walking machine (hexapod) as our experimental platform (see Figure 6), because it is a well-studied morphology in the context of artificial intelligence, with one of its first appearances as Ghengis (Brooks 1989). The purpose of this section is *not* to develop an optimal walking strategy for this system. Contrary, this morphology was chosen, because the tripod gait (see Figure 6) is known to be one of the optimal locomotion behaviors, which can be implemented efficiently in various ways. This said, learning a control for this gait is not trivial, and hence, a good test bed to evaluate our complexity bound for CRBM based policies.

This section is organized in three parts. The first part presents the experimental set-up as far as it is required to understand the results. The second part describes how the CRBM complexity parameter m was estimated from the data. The third part presents the results of the experiment and compares them with the theoretical bound.

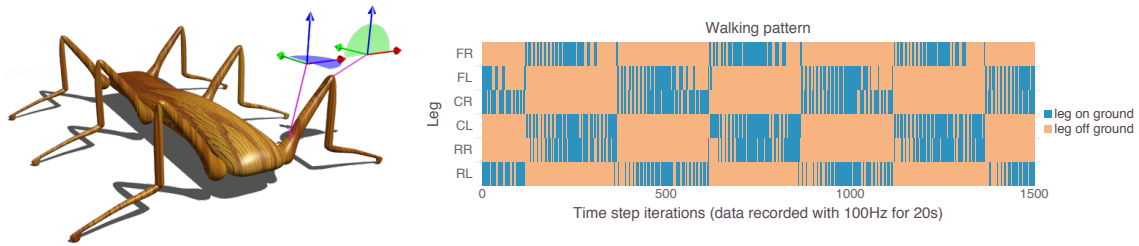


Figure 6: Hexapod set-up. Left-hand side: The simulated hexapod with a display of the joint configurations. For a detailed discussion, see Section 5.1. Right-hand side: Visualization of the target walking pattern. The plot shows which leg touched the ground at which point in time. Blue areas refer to a contact with a the ground, while orange areas refer to points in time during which the correspond leg did not touch the ground. The different legs are plotted over the y -axis, while each point on the x -axis refers to a single point in time.

5.1 Simulation

The hexapod was simulated with YARS (Zahedi et al. 2008), which is a mobile robot simulator based on the bullet physics engine (Coumans 2012). Each segment of the hexapod is defined by its physical properties (dimension, weight, etc.) and each actuator is defined by its force, velocity and its angular range. In the case of the hexapod shown in Figure 6 (left-hand side), the main body’s dimension (bounding box) is 4.4m length, 0.7m width, 0.5m height, and the weight is 2kg. Each leg consists of three segments (femur, tarsus, tibia), of which the two lower segments (tarus, tibia) are connected by a fixed joint. The leg segments were freely modeled with respect to the dimensions of an insect leg. The actuator which connects the femur and tarsus (knee actuator) only allows rotations around the local y -axis of the femur segment (see Figure 6). The maximal deviation for the femur-tarsus actuator is limited to $\omega_{\text{fe-ta}} \in [-15^\circ, 25^\circ]$. For actuators which connect the main body with the femur (ma-fe), the maximal deviation is limited to $\omega_{\text{ma-fe}} \in [\pm 35^\circ]$. The rotation axis of the ma-fe actuator is limited to the local z -axis of the main body. In bullet, an actuator is defined by its impulse (set to 1Ns for both actuators) and its maximal velocity (set to $0.75\text{rad/s} \approx 42^\circ/\text{s}$ for both actuators). It must be noted here, that the sensors and actuators are all mapped onto the interval $[-1, 1]$, which means that a sensor value of $\tilde{S}_i^t = 1$ refers to the maximal current deviation of the corresponding joint. In the same sense, an actuator value $\tilde{A}_j^t = 1$ refers to the motor command to deviate the corresponding joint to its maximal position. For simplicity, we refer to the two types of joints as shoulder (main body-femur) and knee (femur-tarus).

The policy update frequency was set to 10Hz, i.e., the controller received ten sensor values and generated ten actuator values per second. The target behavior of the hexapod (a tripod walking gait, see Figure 6) was generated by an open-loop controller which applied phase shifted sinus oscillations to the actuators. For each leg, the sinus oscillations were discretized into 50 pairs of actuator values, which means that one locomotion step requires 50 time steps (5 seconds) to complete.

For the training and analysis, the sensor and actuator value data was discretized into 16 equidistant bins for each sensor and actuator. This corresponds to four binary input units for each sensor and four binary output units for each actuator. Combined into two random variables $S = (S_1, S_2, \dots, S_{12}), A = (A_1, A_2, \dots, A_{12})$, this leads to a total of 16^{12} possible values ($|\mathcal{S}| =$

$|\mathcal{A}| = 16^{12}$) corresponding to a total of 48 binary input and 48 binary output units. In the following sections, we only refer to this pre-processed data, which means that calculations and the training of the CRBMs described in the remainder of this section refer to the two random variables S, A .

The next section will discuss the estimation of the sufficient controller complexity that is able to reproduce the desired tripod walking gait.

5.2 Estimation of the Sufficient Complexity

Before the estimation procedure and results are presented, we restate the inequality given in Theorem 2, which is given by

$$m \geq |\mathcal{S}| + d^{\mathcal{S}} - 1. \quad (11)$$

This means that a CRBM should not require more hidden units (m) than the sum of the support set cardinality $|\mathcal{S}|$ and embodied behavior dimension d (minus 1). The following paragraphs explain how these two values were calculated from the recorded data.

The first step in estimating the sufficient controller complexity m of the CRBM policy model is the estimation of the support’s cardinality $|\mathcal{S}|$. It was mentioned above that there are 16^{12} possible sensor values. The necessary complexity of a CRBM policy for a specific behavior depends on the actually used number of sensor values, which is also known as the sensor support set. By estimating the cardinality of the support set, we know how many relevant sensor values the CRBM needs to handle to reproduce the behavior of interest.

The estimation of the support set cardinality depends on the quality of the sample. Therefore, we sampled 10^5 sensor values to ensure a sufficient convergence of the relative frequencies of the sensor values. The upper left plot in Figure 7 shows the histogram for all recorded sensor values. The orange vertical line shows where we have pruned the data so that 80% of the recorded data was kept. The lower left plot in Figure 7 shows the remaining data. With this procedure (recording, estimating relative frequencies, pruning the data to 80%), we estimated the cardinality of the sensor support set at $|\mathcal{S}| = 63$. The pruning threshold of 80% might appear arbitrary here. To clarify, estimating the support from data is an interesting research topic by itself, which, however, goes beyond the scope of this work. Our underlying assumption for the pruning is that the sampling is noisy. We estimated the noise empirically by analyzing the histogram and decided for a threshold that seemed reasonable to us. We want to point out that the threshold was estimated before the results of the experiments (see next section) were available.

The next step is to estimate the embodied behavior dimension, which is done here based on the affine rank of the empirically estimated internal world model $\gamma(s, a; s')$. For the sake of readability, we defer the justification for the replacement of the embodiment-behavior dimension by the affine rank of the internal world model to the appendix (see Appendix B).

Given the internal world model γ , the affine rank is calculated in the following way:

$$d^{\mathcal{S}} = \sum_{s \in \mathcal{S}} \text{rank}((\gamma(s, a_0; s') - \gamma(s, a; s'))_{s' \in \mathcal{S}, a \in \mathcal{A}}). \quad (12)$$

To estimate the internal world model $\gamma(s, a; s')$, we pruned the data in accordance with the estimated support set cardinality. This means that we removed all pairs of S, A for which S is not the in pruned support set \mathcal{S} . For the remaining data, we counted the occurrences of s^{t+1}, s^t, a^t and filled the matrix $\gamma(s^t, a^t; s^{t+1})$. The matrix is initialized with zero and each row is normalized by the row

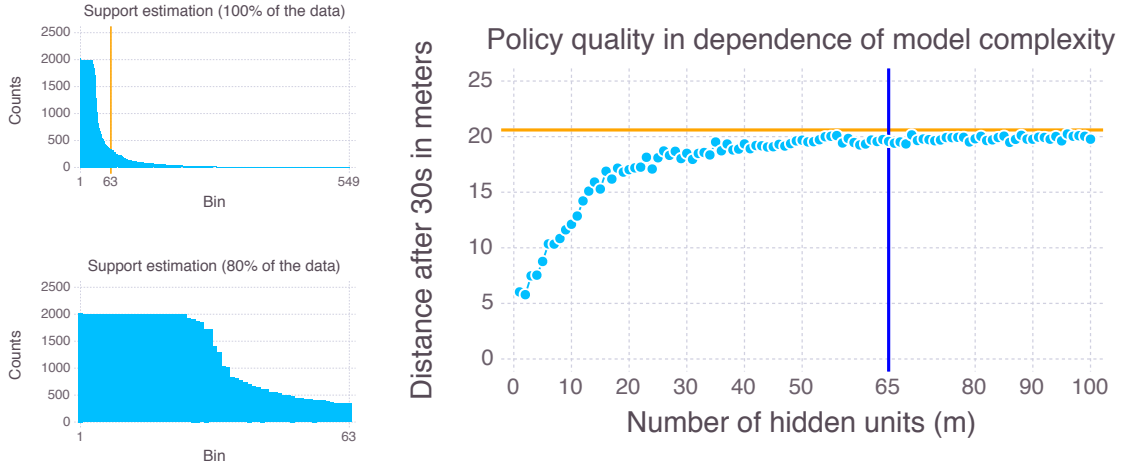


Figure 7: Experimental results. Left-hand side: Estimation of the support set cardinality (before and after pruning). Right-hand side: Performance of the best CRBM for different complexity parameters m in comparison to the performance of the target behavior (horizontal orange line). The vertical blue line indicates the m estimated from the data (see Section 5.2).

sum. The resulting matrix does not model a conditional probability distribution, because many rows have row sum zero. As we are only interested in the affine rank of the matrix, this is of no matter to us. The resulting value is $d^S = 3$.

Resulting estimation of the model complexity: It follows that the CRBM is able to represent the target behavior whenever the number of hidden units satisfies

$$m \geq |\mathcal{S}| + d^S - 1 = 63 + 3 - 1 = 65. \tag{13}$$

To evaluate the tightness of this bound, we conducted a series of experiments, which are explained in the following section.

5.3 Experiments to Evaluate the Tightness of the Complexity Estimation

Before the experiments can be described, there is an important note to make. This work is concerned with the minimally required complexity that is sufficient to control an embodied agent, such that it is able to produce a set of desirable behaviors (which also includes all behaviors as well as one specific behavior). Here, we are *not* concerned with the question how these CRBMs should be trained optimally. This is why we used a standard training algorithm for RBMs (Hinton 2002; 2012) and conducted a large scan over different complexity parameters m . For each $m = 1, 2, 3, \dots, 100$ we trained 100 CRBMs with the following learning parameters: epochs = 20000, batch size = 50, learning rate $\alpha = 1.0$, momentum = 0.1, Gaussian distributed noise on sensor data = 0.01, weight cost = 0.001, bins = 16, CRBM update iterations = 10, on a data set of 10^4 pairs of sensor and actuator values. Each trained CRBM was evaluated ten times, by applying it to the hexapod and recording the covered distance for 30 seconds. The performance of the CRBMs is measured against the target tripod walking gait, which achieves 20.6 meters during the same time. As we are not concerned

with the training of the CRBMs in this work but instead with the best performance of a CRBM for a given m , we choose to take the maximally covered distance over all 1000 evaluations (100 CRBMs, 10 evaluation runs for each CRBM) to estimate the quality of a CRBM with a given complexity m . The plot on the right-hand side of Figure 7 shows the resulting performance for all scanned values of m . The results show that our estimation is fairly tight, which means the performance of the CRBMs converges to the optimal behavior close to the estimated value of $m = 65$.

6 Conclusions

We presented an approach for studying and implementing cheap design in the context of embodied artificial intelligence. In this context, we referred to cheap design as the reduction of the controller complexity that is possible through an exploitation of the agent’s body and environment. We developed a theory to determine the minimal controller complexity that is sufficient to generate a given set of desired behaviors. Being more precise, we studied the way in which embodiment constraints induce equivalent policies in the sense that they generate the same observable behaviors. This led to the definition of the effective dimension of an embodied system, the *embodied behavior dimension*. In this way, we were able to define low-dimensional policy models that can generate all possible behaviors. Such policy models are related to the classical notion of universal approximation.

We used CRBMs as a platform of study, for which we presented non-trivial universal approximation results in both the non-embodied and the embodied settings. While the non-embodied universal approximation requires an enormous number of hidden units (exponentially many in the number of input and output units), embodied universal approximation can be achieved using essentially only as many hidden units as the effective dimension of the system. Notably, our construction depends only on the embodied behavior dimension and, therefore, is independent of the specific embodiment constraints.

Experiments conducted on a walking machine demonstrate the tightness of the estimated number of hidden units for a CRBM controller. This shows the practical utility of our theoretical analysis. To the best of our knowledge, the presented formalism and results are amongst the first quantitative contributions to cheap design in embodied artificial intelligence.

Appendix

A Technical Proofs

Details of Equation (9). In order to obtain the dimension of the image $\text{Beh} = \psi(\Delta_{\mathcal{A}}^{\mathcal{S}})$, we first look at the affine space $\text{aff}(\Delta_{\mathcal{A}}^{\mathcal{S}})$. A basis of this space can be obtained by building differences of vertices of $\Delta_{\mathcal{A}}^{\mathcal{S}}$. The vertices are the deterministic policies $\pi^f(s; a) := \delta_{f(s)}(a)$, for all $s \in \mathcal{S}$ and $a \in \mathcal{A}$, each of which is characterized by a function $f: \mathcal{S} \rightarrow \mathcal{A}$. We fix a deterministic policy π^f , with $f(s) = a_0$ for all s , for some $a_0 \in \mathcal{A}$, and consider the differences $e_{(s,a)} := \pi^f - \pi^{f_{(s,a)}}$ for all possible pairs (s, a) with $a \neq a_0$, where $f_{(s,a)}$ is the function that differs from f only at s , where it takes value $f_{(s,a)}(s) = a$. This set of $e_{(s,a)}$ ’s is a basis of $\text{aff}(\Delta_{\mathcal{A}}^{\mathcal{S}})$. There are $|\mathcal{S}|(|\mathcal{A}| - 1)$ of these vectors, which corresponds with the dimension of $\Delta_{\mathcal{A}}^{\mathcal{S}}$.

Now, the image $\text{Beh} = \psi(\Delta_{\mathcal{A}}^{\mathcal{S}})$ has the dimension of $\text{aff}(\psi(\Delta_{\mathcal{A}}^{\mathcal{S}})) = \psi(\text{aff} \Delta_{\mathcal{A}}^{\mathcal{S}})$, which is the vector space spanned by $\psi(e_{(s,a)}) = \beta(w; s)(\alpha(w, a_0; dw') - \alpha(w, a; dw'))$ for all possible

pairs (s, a) with $a \neq a_0$. Therefore, the dimension is given by the number of linearly independent $\psi(e_{(s,a)})$'s (the rank of the matrix with rows $e_{(s,a)}$). \square

Details of Example 1. We consider an exponential family \mathcal{E} of probability distributions on the set $\mathcal{A}^{\mathcal{S}}$ of functions $f: \mathcal{S} \rightarrow \mathcal{A}$. Let this exponential family be specified by the sufficient statistic $F := \psi \circ \eta$, where $\eta: \mathcal{A}^{\mathcal{S}} \rightarrow \Delta_{\mathcal{A}^{\mathcal{S}}}$; $f \mapsto \pi^f$, $\pi^f(s; a) := \delta_{f(s)}(a)$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}$, and ψ is the policy-behavior map, represented by the matrix $E \in \mathbb{R}^{d \times (\mathcal{S} \times \mathcal{A})}$. Note that, given a basis of $\text{aff}(\psi(\Delta_{\mathcal{A}^{\mathcal{S}}}))$, composed of d vectors in $\text{aff}(\psi(\Delta_{\mathcal{A}^{\mathcal{S}}}))$, we can represent each $e_{(s,a)}$ and $\psi(\pi)$ with respect to this basis by a vector of length $d \leq |\mathcal{S}|(|\mathcal{A}| - 1)$. The exponential family \mathcal{E} consists of all probability distributions of the form

$$p_{\theta}(f) = \frac{\exp(\theta^{\top} F(f))}{\sum_{f'} \exp(\theta^{\top} F(f'))}, \quad \text{for all } f \in \mathcal{A}^{\mathcal{S}}, \text{ for all } \theta \in \mathbb{R}^d.$$

The *moment map* μ maps probability distributions to the corresponding expectation value of the sufficient statistic,

$$\mu: \Delta_{\mathcal{A}^{\mathcal{S}}} \rightarrow \mathbb{R}^d; p \mapsto \sum_f F(f)p(f) = \sum_f \psi(\eta(f))p(f) = \psi\left(\sum_f \eta(f)p(f)\right) = \psi(\pi^p),$$

where $\pi^p := \eta(p) = \sum_f \pi^f p(f) \in \Delta_{\mathcal{A}^{\mathcal{S}}}$. A key property of the moment map is that it maps the closure $\bar{\mathcal{E}}$ of \mathcal{E} bijectively to the set $\mu(\Delta_{\mathcal{A}^{\mathcal{S}}})$ of all possible expectation values. We have

$$\psi(\eta(\bar{\mathcal{E}})) = \mu(\bar{\mathcal{E}}) = \mu(\Delta_{\mathcal{A}^{\mathcal{S}}}) = \psi(\Delta_{\mathcal{A}^{\mathcal{S}}}).$$

Now we only need to show that the set $\eta(\mathcal{E}) = \{\pi^p = \sum_f \pi^f p(f) : p \in \mathcal{E}\}$ is contained in $\mathcal{E}_{\mathcal{A}^{\mathcal{S}}}$. That this is true can be seen from

$$\begin{aligned} \pi^{p_{\theta}}(s; a) &= \sum_f \pi^f(s; a) p_{\theta}(f) = \sum_{f: f(s)=a} p_{\theta}(f) = p_{\theta}(\{f: f(s)=a\}) = \frac{\exp(\theta^{\top} E(s, a))}{\sum_{a'} \exp(\theta^{\top} E(s, a'))} \\ &= \pi_{\theta}(s; a), \quad \text{for all } a \in \mathcal{A} \text{ and } s \in \mathcal{S}, \quad \text{for all } \theta \in \mathbb{R}^d, \end{aligned}$$

where we used

$$\begin{aligned} p_{\theta}(f) &= \frac{\exp(\theta^{\top} E(\pi^f))}{\sum_{f'} \exp(\theta^{\top} E(\pi^{f'}))} = \frac{\exp(\theta^{\top} \sum_s E(s, f(s)))}{\sum_{f'} \exp(\theta^{\top} \sum_{s'} E(s', f'(s')))} = \prod_s \frac{\exp(\theta^{\top} E(s, f(s)))}{\sum_a \exp(\theta^{\top} E(s, a))} \\ &= \prod_s p_{\theta}(\{f': f'(s) = f(s)\}). \end{aligned}$$

In fact, since μ is a bijection between $\bar{\mathcal{E}}$ and $\mu(\Delta_{\mathcal{A}^{\mathcal{S}}}) = \psi(\Delta_{\mathcal{A}^{\mathcal{S}}})$, we have that ψ is a bijection between $\bar{\mathcal{E}}_{\mathcal{A}^{\mathcal{S}}} = \eta(\bar{\mathcal{E}})$ and $\psi(\Delta_{\mathcal{A}^{\mathcal{S}}})$. \square

Proof of Lemma 1. Assume first that $\mathcal{S} = \mathcal{S}$. Geometrically, the policy-behavior map ψ projects the policy polytope linearly into a polytope of dimension d . A d -dimensional projection of a polytope is equal to the projection of its d -dimensional faces. This implies that the result of applying any given policy from $\Delta_{\mathcal{A}^{\mathcal{S}}}$ can be achieved equally well by applying a policy from a d -dimensional face of $\Delta_{\mathcal{A}^{\mathcal{S}}}$. See Figure 8 for an illustration of what we mean.

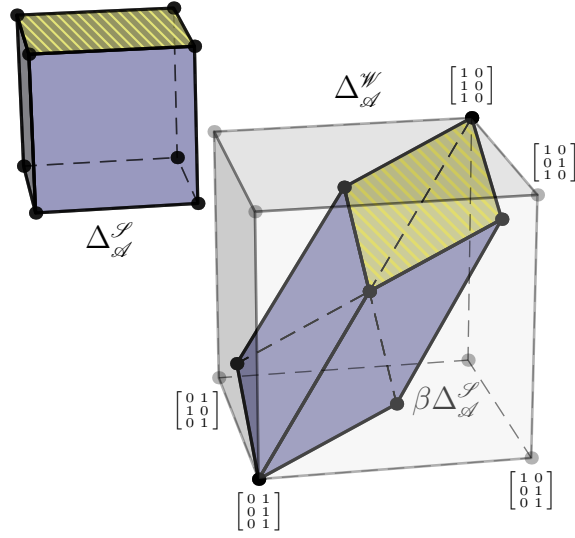


Figure 8: Illustration of the projection of the set of policies $\Delta_{\mathcal{A}}^{\mathcal{S}}$ into a lower dimensional polytope. In this example $|\mathcal{S}| = 3$ and $|\mathcal{A}| = 2$, such that the policy polytope $\Delta_{\mathcal{A}}^{\mathcal{S}}$ is a cube. The set of all kernels $\Delta_{\mathcal{A}}^{\mathcal{W}}$, $|\mathcal{W}| = 3$, is also a cube. The set of all realizable kernels within that set, $p(w; a) = \sum_s \beta(w; s)\pi(s; a)$, for all choices of the policy π , is the projection of $\Delta_{\mathcal{A}}^{\mathcal{S}}$ by β . In this example $\text{rank}(\beta) = 2$ and the projection of the policy polytope is two-dimensional polygon (the blue hexagon). This projection by β represents one part of the projection by the policy-behavior map. The two-dimensional faces of the policy polytope (one of them highlighted in dashed yellow) have the same image as the entire policy polytope.

Now, the d -dimensional faces of the policy polytope $\Delta_{\mathcal{A}}^{\mathcal{S}}$ consist of those policies with at most $|\mathcal{S}| + d$ non-zero entries. The arguments for this are as follows. The policy polytope is a product of simplices $\Delta_{\mathcal{A}}^{\mathcal{S}} = \times_{s \in \mathcal{S}} \Delta_{\mathcal{A}}^s$, where the s -th factor corresponds to the set of all possible probability distributions $\pi(s; \cdot)$. The faces of $\Delta_{\mathcal{A}}^{\mathcal{S}}$ are products of faces of its factors and have the form $\times_{s \in \mathcal{S}} \Delta_{\mathcal{A}_s}^s$, where $\mathcal{A}_s \subseteq \mathcal{A}$ for all $s \in \mathcal{S}$. Each face of $\Delta_{\mathcal{A}}^{\mathcal{S}}$ corresponds to a choice of positions $\mathcal{A}_s \subseteq \mathcal{A}$ of the non-zero entries of $\pi(s; \cdot)$ for all $s \in \mathcal{S}$. The d -dimensional faces are those for which $\sum_{s \in \mathcal{S}} (|\mathcal{A}_s| - 1) = d$, meaning that they consist of policies which have at most $\sum_{s \in \mathcal{S}} |\mathcal{A}_s| = |\mathcal{S}| + d$ non-zero entries.

Consider now $\mathcal{S} \subseteq \mathcal{S}$. The projection of the policy polytope by the \mathcal{S} embodiment matrix can be regarded as a composition which first projects $\Delta_{\mathcal{A}}^{\mathcal{S}}$ to $\Delta_{\mathcal{A}}^{\mathcal{S}}$ and then projects $\Delta_{\mathcal{A}}^{\mathcal{S}}$ by $\psi^{\mathcal{S}}$. Now we can use the same arguments as above, with the difference that now only need to represent the $d^{\mathcal{S}}$ -dimensional faces of the polytope $\Delta_{\mathcal{A}}^{\mathcal{S}}$. \square

B Estimation of the Embodied Behavior Dimension based on the Internal World Model

In many situations, the embodied behavior dimension is not available from a perspective that is intrinsic to the agent, as the agent does not have direct access to the sensor kernel β nor to the world

kernel α . From that perspective, only an internal version of the world model is accessible, which we refer to as *internal world model*. It is defined as a kernel $\gamma \in \Delta_{\mathcal{S}}^{\mathcal{S} \times \mathcal{A}}$, assigning to each sensor state s with positive probability and each actuator state a the next sensor state s' , that is,

$$\gamma(s, a; s') = \int_{\mathcal{W}} \left\{ \int_{\mathcal{W}} \beta(w'; s') \alpha(w, a; dw') \right\} \mathbb{P}(s; dw), \quad (14)$$

where

$$\mathbb{P}(s; dw) := \frac{\beta(w; s)}{\int_{\mathcal{W}} \beta(w''; s) \mathbb{P}(dw'')} \mathbb{P}(dw).$$

Note that the internal world model is not completely determined by β and α . It also depends on the distribution $\mathbb{P}(dw)$ of the world states w . If we choose this distribution to be a fixed reference distribution of world states, then the world model will be determined by α and β only. However, in order to describe the actual distribution of world states, we have to take into account the contribution of the agent's policy π . This implies that, if the policy is subject to changes in terms of a learning process, then, in general, the world model will also be time dependent.

On the other hand, γ is the only information about world dynamics that is intrinsically available to the agent. The extent to which γ is not a good replacement for α depends on how much the agent can “see” from the world with its sensors. If the agent has direct access to the world state, that is $W^t = S^t$, then α and γ coincide. However, this is not very realistic. Generically, only partial observation of the world is possible. Now the question arises whether it is possible to determine the embodied behavior dimension d in terms of γ even in cases where the agent has only partial access to the world state. This is indeed possible under specific conditions which are satisfied in our experimental setup. We first present these conditions in general terms, before we then relate them to our experiment at the end of this section. Let the world state w consist of two parts s and r , where s is directly accessible to the agent and r is the remaining part of the world, which is hidden to the agent. The situation is illustrated in Figure 9.

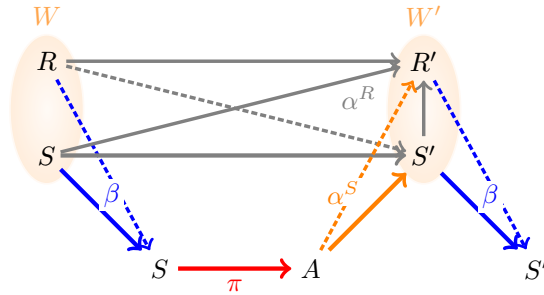


Figure 9: Special causal structure of the sensorimotor loop. The dashed arrows are the ones that we omit within our assumptions.

In this interpretation of the world state, the sensor kernel β is simply the identity map $s \mapsto s$. Furthermore, this interpretation sets structural constraints on the world transition kernel α , which assigns a probability distribution of the next world state $w' = (s', r')$ given the current world state $w = (s, r)$ and an actuator value a . As r is assumed to be hidden to the agent, s' should not depend on r . This leads to the following natural factorization of α :

$$\alpha(s, r, a; s', dr') = \alpha^S(s, a; s') \cdot \alpha^R(r, s, s'; dr'). \quad (15)$$

With this assumption, we obtain as internal world model

$$\gamma(s, a; s') = \alpha^S(s, a; s'), \quad \text{whenever } \mathbb{P}(s) > 0, \quad (16)$$

and the following transition probabilities from a world state $w = (s, r)$ with positive probability to a world states $w' = (s', r')$:

$$\begin{aligned} \mathbb{P}^\pi(s, r; s', dr') &= \sum_{s'', a} \beta(s, r; s'') \pi(s''; a) \alpha(s, r, a; s', dr') \\ &= \sum_a \pi(s; a) \alpha^S(s, a; s') \alpha^R(r, s, s'; dr') \\ &= \alpha^R(r, s, s'; dr') \underbrace{\sum_a \pi(s; a) \alpha^S(s, a; s')}_{=: \mathbb{Q}^\pi(s; s')}. \end{aligned} \quad (17)$$

This shows that $\mathbb{P}^\pi = \mathbb{P}^{\pi^*}$ if and only if $\mathbb{Q}^\pi = \mathbb{Q}^{\pi^*}$, and therefore the embodied behavior dimension is given by the dimension of the image of $\pi \mapsto \mathbb{Q}^\pi$. This is given by the affine rank of the kernel α^S , which coincides with the affine rank of γ :

$$d = \sum_s \text{rank}((\gamma(s, a_0; s') - \gamma(s, a; s'))_{s' \in \mathcal{S}, a \in \mathcal{A}}), \quad (18)$$

where a_0 is any fixed value in \mathcal{A} , for all s .

This applies to our hexapod experiment discussed in Section 5 for the following reason. In the special case of the tripod gait of a hexapod on an even and otherwise featureless plane, the next joint angles S^{t+1} are only determined by the current joint angles S^t and the current action A^t . The rest of the world, here denoted by R , contains information such the contact points of the legs with the ground. This information is carried from one time step to the next, as it determines how the hexapod walks along the plane. Nevertheless, the contact points of the legs do not influence the joint angles. Hence, in our experiment, S' is conditionally independent of R given S and A . Furthermore, R' is conditionally independent of A given R , S and S' as the contact points of the legs with the ground are only determined by the relative joint angles, and not by the current action. Therefore, we can estimate the embodied behavior dimension by the rank of the internal world model.

C Generalizations

C.1 SML with Internal State

In the main part of the paper we considered reactive SMLs. In a more general setting, the agent may be equipped with some sort of memory or internal representation of the world. In this case, besides from the world state, the sensor state, and the actuator state, the SML also includes an internal state variable. As in the reactive SML, the dynamics of these variables are governed by Markov transition kernels, but the causality structure is slightly different. Let \mathcal{W} , \mathcal{S} , \mathcal{C} , and \mathcal{A} denote the sets of possible states of the world, the sensors, the internal state, and the actuators. Then the

Markov kernels are

$$\begin{aligned}
 \beta: \mathcal{W} &\rightarrow \Delta_{\mathcal{S}}, \\
 \varphi: \mathcal{C} \times \mathcal{S} &\rightarrow \Delta_{\mathcal{C}}, \\
 \pi: \mathcal{C} &\rightarrow \Delta_{\mathcal{A}}, \\
 \alpha: \mathcal{W} \times \mathcal{A} &\rightarrow \Delta_{\mathcal{W}}.
 \end{aligned} \tag{19}$$

See Figure 10 for an illustration of this causality structure.

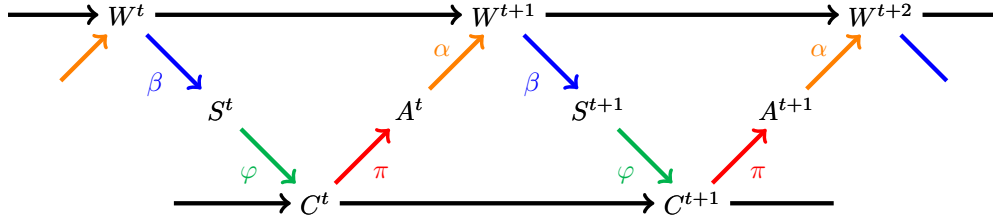


Figure 10: Causal structure of a SML with internal state. Here W^t, S^t, C^t, A^t are the states of the world, sensors, internal variable, and actuators at the discrete time t .

As in the reactive case discussed in the main text, here we also want to consider the (extrinsic) behavior of the agent, which is described in terms of the stochastic process of world states. In this case, however, we condition these processes not only on an initial world state but also on an initial internal state. The difficulty arising here is that, in general, in the presence of an internal state the stochastic process over the world states is not Markovian. The world state transition at each time step does not only depend on the previous world state but it also depends on a longer history, encoded in the internal state.

For example, when navigating a territory, a robot endowed with an internal state could operate in the following way. If at a given time step the robot detects an obstacle ahead, $s = \text{“obstacle”}$, then, in conjunction with a current internal state $c = \text{“safe”}$, the new internal state could become $c' = \text{“attentive”}$, in which case the policy would choose $a = \text{“maintain direction”}$. However, if the current internal state was $c = \text{“attentive”}$, the new internal state could become $c' = \text{“alert”}$, in which case the policy would choose between the actions $a' = \text{“turn left”}$ and $a' = \text{“turn right”}$ with probability $\frac{1}{2}$. This example shows that the internal state may contain information about the history of world and sensor states, which is not available from the current world and sensor states alone.

Nonetheless, for any fixed choice of the kernels $\beta, \varphi, \pi, \alpha$ and a starting value (w^0, c^0) at time $t = 0$, the SML defines a (discrete-time homogeneous) Markov chain with state space $\mathcal{W} \times \mathcal{S} \times \mathcal{C} \times \mathcal{A}$. The transition probabilities of this chain are given by

$$\mathbb{P}^\pi(w^0, s^0, c^0, a^0; dw^1, s^1, c^1, a^1) = \alpha(w^0, a^0; dw^1) \beta(w^1; ds^1) \varphi(c^0, s^1; dc^1) \pi(c^1; da^1).$$

Furthermore, the process with state space $\mathcal{W} \times \mathcal{C}$ is also Markovian. The transition probabilities of this chain are given by

$$\psi(w, c; dw', dc') = \int_{\mathcal{S}'} \int_{\mathcal{A}} \pi(c; da) \alpha(w, a; dw') \beta(w'; ds') \varphi(c, s'; dc').$$

The process on \mathcal{W} (extrinsic behavior) is the marginal of the process on $\mathcal{W} \times \mathcal{C}$ (extrinsic-intrinsic behavior). We can study some properties of the extrinsic behavior in terms of the properties of the extrinsic-intrinsic behavior. The latter is easier to analyze, since it is Markovian. In particular, we can study it in the same way we studied the extrinsic behavior in the reactive SML.

More explicitly we have the following. Writing $\xi(w, c; dw') = \int_{\mathcal{A}} \pi(c; da) \alpha(w, a; dw')$ and $\phi(w', c; dc') = \int_{\mathcal{S}'} \beta(w'; ds') \varphi(c, s'; dc')$, the transition probabilities for the process on $\mathcal{W} \times \mathcal{C}$ are given by

$$\psi(w, c; dw', dc') = \xi(w, c; dw') \phi(w', c; dc'). \quad (20)$$

For each (w, c) , the probability distribution $\xi(w, c; \cdot) \in \Delta_{\mathcal{W}}$ is the projection of $\pi(c; \cdot) \in \Delta_{\mathcal{A}}$ by the linear map defined by $\alpha(w, \cdot; \cdot)$. If the intersection of the null-spaces of $\alpha(w, \cdot; \cdot)$ for all w has a positive dimension, then there is a positive dimensional set of policies π that are mapped to the same ξ and hence to the same behavior. In order to obtain that behavior, it is sufficient to represent one of the policies that map to ξ , in contrast to the potentially much larger set of all policies that map to the same ξ . A similar observation applies to φ . This shows that already when considering the process on $\mathcal{W} \times \mathcal{C}$ (the combined extrinsic-intrinsic behavior of the agent), many policies may be identified. Embodiment constraints restrict the possible behaviors. When considering only the process over \mathcal{W} (the extrinsic behavior of the agent), many more policies may be identified with the same behavior. The detailed study of projections from combined behaviors to extrinsic behaviors is left for future work.

C.2 Continuous Sensor and Actuator State Spaces

We have considered systems where \mathcal{S} and \mathcal{A} are finite sets. In some case it can be more natural to consider continuous sensor and actuator spaces. The continuous case brings some subtleties with it. In particular, the set of policies with continuous state spaces is infinite dimensional. In this case one has to depart from linear algebra and use functional analysis. Furthermore, in the setting of continuous sensor and actuator spaces usually it is not possible to achieve universal approximation by one fixed model. Rather, one says that a class of models has the universal approximation property, meaning that for each given error tolerance, there is a model in that class, that can approximate to within that error tolerance. Nonetheless, one can measure the approximation performance in terms of the (finite) number of parameters or hidden variables that a model needs in order to satisfy a given error tolerance. Continuous policy models can be defined in terms of stochastic feedforward neural networks with continuous variables or also in terms of CRBMs with Gaussian output units. Here, the complexity of a model can be measured in terms of the number of hidden variables.

Acknowledgment

We would like to acknowledge support for this project from the DFG Priority Program Autonomous Learning (DFG-SPP 1527). G. M. and K. G.-Z. would like to thank the Santa Fe Institute for hosting them during the work on this article.

References

N. Ay and K. Zahedi. On the causal structure of the sensorimotor loop. In M. Prokopenko, editor, *Guided Self-Organization: Inception*. Springer, 2014.

- H. Bauer. *Probability Theory*. De Gruyter studies in mathematics. Bod Third Party Titles, 1996. ISBN 9783110139358. URL <http://books.google.com/books?id=w76IHsPHybcC>.
- Y. Bengio. Learning deep architectures for ai. *Foundations and Trends in Machine Learning*, 2(1): 1–127, 2009. ISSN 1935-8237. doi: 10.1561/22000000006. URL <http://dx.doi.org/10.1561/22000000006>.
- W. Bialek and N. Tishby. Predictive information. SEE Bialek, Nemenman, Tishby, 2001, 1999.
- V. Braitenberg. *Vehicles*. MIT Press, Cambridge MA, 1984.
- R. A. Brooks. A robot that walks; emergent behaviors from a carefully evolved network. *Neural Comput.*, 1(2):253–262, June 1989. ISSN 0899-7667. doi: 10.1162/neco.1989.1.2.253. URL <http://dx.doi.org/10.1162/neco.1989.1.2.253>.
- R. A. Brooks. Intelligence without reason. In J. Myopoulos and R. Reiter, editors, *Proceedings of the 12th International Joint Conference on Artificial Intelligence (IJCAI-91)*, pages 569–595, Sydney, Australia, 1991a. Morgan Kaufmann publishers Inc.: San Mateo, CA, USA.
- R. A. Brooks. Intelligence without representation. *Artificial Intelligence*, 47(1-3):139–159, 1991b.
- D. D. Clark and L. Sokoloff. Circulation and energy metabolism of the brain. In G. J. Siegel, B. W. Agranoff, R. W. Albers, S. K. Fisher, and M. D. Uhler, editors, *Basic Neurochemistry: Molecular, Cellular and Medical Aspects*, chapter 31. Lippincott-Raven, Philadelphia, 6th edition, 1999.
- E. Coumans. Bullet physic sdk manual. www.bulletphysics.org, 2012.
- Y. Freund and D. Haussler. *Unsupervised Learning of Distributions of Binary Vectors Using Two Layer Networks*. Technical report. Computer Research Laboratory, University of California, Santa Cruz, 1994.
- G. Hinton. A practical guide to training restricted Boltzmann machines. In G. Montavon, G. Orr, and K.-R. Müller, editors, *Neural Networks: Tricks of the Trade*, volume 7700 of *Lecture Notes in Computer Science*, pages 599–619. Springer Berlin Heidelberg, 2012. ISBN 978-3-642-35288-1. doi: 10.1007/978-3-642-35289-8_32. URL http://dx.doi.org/10.1007/978-3-642-35289-8_32.
- G. E. Hinton. Training products of experts by minimizing contrastive divergence. *Neural Computation*, 14(8):1771–1800, 2002.
- A. S. Klyubin, D. Polani, and C. L. Nehaniv. Tracking information flow through the environment: Simple cases of stigmerg. In J. Pollack, editor, *Artificial Life IX: Proceedings of the Ninth International Conference on the Simulation and Synthesis of Living Systems*. MIT Press, 2004.
- H. Larochelle and Y. Bengio. Classification using discriminative restricted Boltzmann machines. In W. W. Cohen, A. McCallum, and S. T. Roweis, editors, *Proceedings of the 25th International Conference on Machine Learning (ICML 2008)*, volume 307, pages 536–543, 2008.
- N. Le Roux and Y. Bengio. Representational power of restricted Boltzmann machines and deep belief networks. *Neural Computation*, 20(6):1631–1649, June 2008.

- P. M. Long and R. A. Servedio. Restricted Boltzmann machines are hard to approximately evaluate or simulate. In J. Fürnkranz and T. Joachims, editors, *Proceedings of the 27th International Conference on Machine Learning (ICML 2010)*, pages 703–710. Omnipress, 2010.
- M. Lungarella and O. Sporns. Information self-structuring: Key principle for learning and development. In IEEE, editor, *Proceedings. The 4th International Conference on Development and Learning, 2005.*, pages 25–30, San Diego, CA, 2005. IEEE Press.
- J. Martens, A. Chattopadhyaya, T. Pitassi, and R. Zemel. On the expressive power of restricted Boltzmann machines. In C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Weinberger, editors, *NIPS 26*, pages 2877–2885, 2013.
- T. McGeer. Passive dynamic walking. *International Journal of Robotic Research*, 9(2):62–82, 1990.
- G. Montúfar and N. Ay. Refinements of universal approximation results for deep belief networks and restricted Boltzmann machines. *Neural Computation*, 23(5):1306–1319, 2011.
- G. Montúfar, J. Rauh, and N. Ay. Expressive power and approximation errors of restricted Boltzmann machines. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K. Weinberger, editors, *NIPS 24*, pages 415–423, 2011.
- G. Montúfar, N. Ay, and K. Ghazi-Zahedi. Expressive power of conditional restricted Boltzmann machines. *arXiv preprint arXiv:1402.3346*, 2014.
- R. Pfeifer and J. C. Bongard. *How the Body Shapes the Way We Think: A New View of Intelligence*. The MIT Press (Bradford Books), Cambridge, MA, 2006.
- R. Pfeifer, M. Lungarella, and F. Iida. Self-organization, embodiment, and biologically inspired robotics. *Science*, 318(5853):1088–1093, 2007.
- D. Polani, C. Nehaniv, T. Martinetz, and J. T. Kim. Relevant Information in Optimized Persistence vs. Progeny Strategies. In L. M. Rocha, M. Bedau, D. Floreano, R. Goldstone, A. Vespignani, and L. Yaeger, editors, *Proc. Artificial Life X*, pages 337–343, Cambridge, MA, 2006. MIT Press.
- R. Salakhutdinov, A. Mnih, and G. E. Hinton. Restricted Boltzmann machines for collaborative filtering. In *Proceedings of the 24th International Conference on Machine Learning (ICML 2007)*, pages 791–798, 2007.
- B. Sallans and G. E. Hinton. Reinforcement learning with factored states and actions. *J. Mach. Learn. Res.*, 5:1063–1088, 2004.
- J. Schmidhuber. Driven by compression progress: A simple principle explains essential aspects of subjective beauty, novelty, surprise, interestingness, attention, curiosity, creativity, art, science, music, jokes. *Anticipatory Behavior in Adaptive Learning Systems*, pages 48–76, 2009.
- P. Smolensky. Information processing in dynamical systems: foundations of harmony theory. In *Symposium on Parallel and Distributed Processing*, 1986.
- L. Sokoloff, R. Mangold, R. Wechsler, C. Kennedy, and S. Kety. Effect of mental arithmetic on cerebral circulation and metabolism. *J. Clin. Invest.*, 34(7):1101–1108, 1955.

- D. Sol, N. Garcia, A. Iwaniuk, K. Davis, A. Meade, W. A. Boyle, and T. Székely. Evolutionary divergence in brain size between migratory and resident birds. *PLoS ONE*, 5(3):e9617, 03 2010.
- I. Sutskever and G. E. Hinton. Learning multilevel distributed representations for high-dimensional sequences. *Proceeding of the 11th International Conference on Artificial Intelligence and Statistics*, 2007.
- G. W. Taylor, G. E. Hinton, and S. Roweis. Modeling human motion using binary latent variables. In *NIPS 19*, pages 1345–1352. MIT Press, 2007.
- K. Zahedi, A. von Twickel, and F. Pasemann. Yars: A physical 3d simulator for evolving controllers for real robots. In S. Carpin, I. Noda, E. Pagello, M. Reggiani, and O. von Stryk, editors, *SIMPAR 2008*, LNAI 5325, pages 71–82. Springer, 2008.