



Adaptive Tuning Curve Widths Improve Sample Efficient Learning

Florian Meier*, Raphaël Dang-Nhu and Angelika Steger

Department of Computer Science, ETH Zürich, Zurich, Switzerland

Natural brains perform miraculously well in learning new tasks from a small number of samples, whereas sample efficient learning is still a major open problem in the field of machine learning. Here, we raise the question, how the neural coding scheme affects sample efficiency, and make first progress on this question by proposing and analyzing a learning algorithm that uses a simple reinforce-type plasticity mechanism and does not require any gradients to learn low dimensional mappings. It harnesses three bio-plausible mechanisms, namely, population codes with bell shaped tuning curves, continuous attractor mechanisms and probabilistic synapses, to achieve sample efficient learning. We show both theoretically and by simulations that population codes with broadly tuned neurons lead to high sample efficiency, whereas codes with sharply tuned neurons account for high final precision. Moreover, a dynamic adaptation of the tuning width during learning gives rise to both, high sample efficiency and high final precision. We prove a sample efficiency guarantee for our algorithm that lies within a logarithmic factor from the information theoretical optimum. Our simulations show that for low dimensional mappings, our learning algorithm achieves comparable sample efficiency to multi-layer perceptrons trained by gradient descent, although it does not use any gradients. Furthermore, it achieves competitive sample efficiency in low dimensional reinforcement learning tasks. From a machine learning perspective, these findings may inspire novel approaches to improve sample efficiency. From a neuroscience perspective, these findings suggest sample efficiency as a yet unstudied functional role of adaptive tuning curve width.

Keywords: sample efficiency, neural tuning curves, population codes, gradient-free learning, reinforcement learning

OPEN ACCESS

Edited by:

Rava Azeredo da Silveira,
École Normale Supérieure, France

Reviewed by:

Artur Luczak,
University of Lethbridge, Canada

Cheng Qiu,
University of Pennsylvania,
United States

Jiang Mao,
University of Pennsylvania,
United States, in collaboration with
reviewer CQ

*Correspondence:

Florian Meier
meierflo@inf.ethz.ch

Received: 20 August 2019

Accepted: 29 January 2020

Published: 18 February 2020

Citation:

Meier F, Dang-Nhu R and Steger A
(2020) Adaptive Tuning Curve Widths
Improve Sample Efficient Learning.
Front. Comput. Neurosci. 14:12.
doi: 10.3389/fncom.2020.00012

1. INTRODUCTION

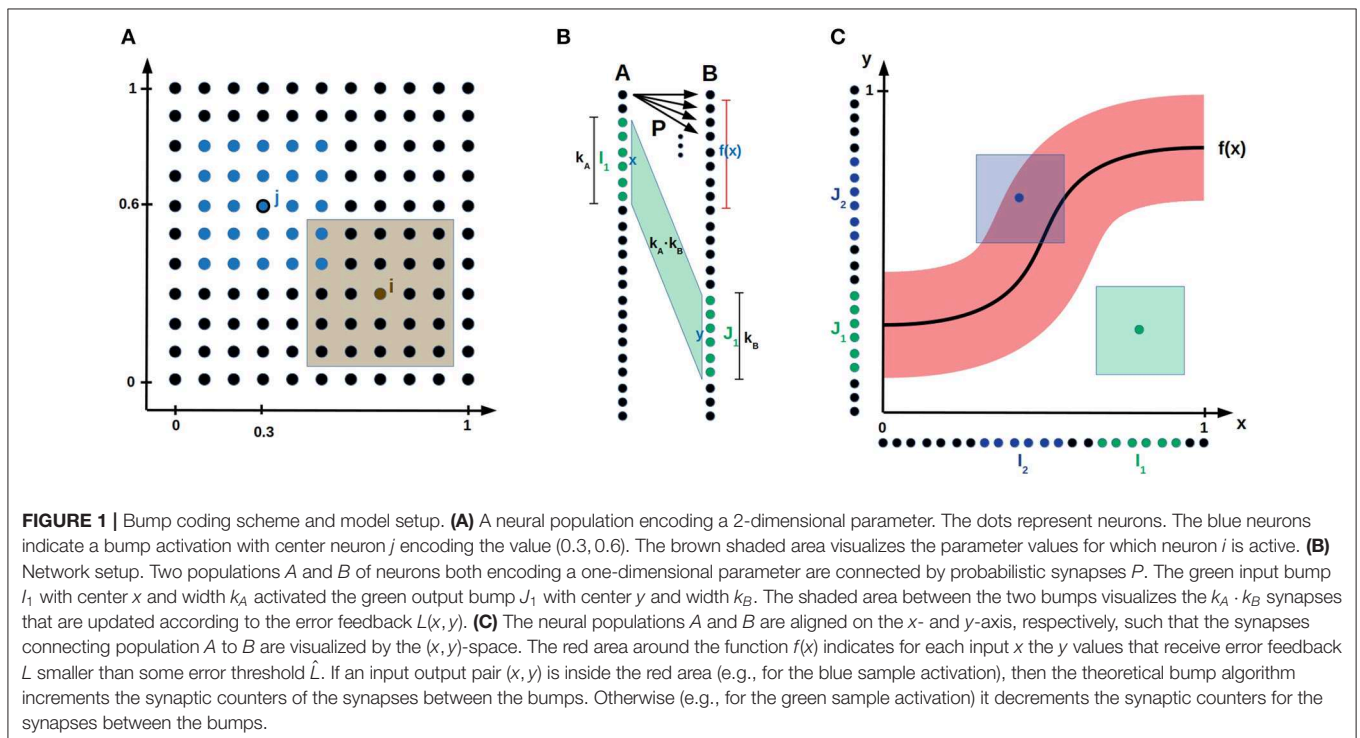
Humans operate in a rich and complex world and are extremely fast in learning new tasks and adapting to new environments. The level of generalization and speed of adaptation achieved by human brains remain unmatched by machine learning approaches, despite tremendous progress in the last years. How do real brains accomplish this outstanding skill of generalization and sample efficient learning, and what are the neural mechanisms that contribute to this ability of fast learning? Here, we investigate how neural coding supports sample efficient learning, by analyzing a learning algorithm that exploits three bio-plausible principles for sample efficient learning, namely, population codes of tuned neurons, continuous attractor mechanisms and probabilistic synapses.

From early on, neuroscience researchers characterized the first order response of single neurons by neural tuning curves (Adrian and Zotterman, 1926). The *neural tuning curve* is defined to be

the neurons mean firing rate as a function of some stimuli parameter. It typically peaks for a preferred parameter value and decays gradually as this parameter moves away from the preferred value, such as in orientation columns in the visual cortex (Hubel and Wiesel, 1959, 1962), spatially tuned cells in auditory cortex (Knudsen and Konishi, 1978), direction selective cells in motor cortex (Georgopoulos et al., 1988) and hippocampal place and head direction cells (O'Keefe, 1976; Ranck, 1985). In populations of tuned neurons, narrow (broad) tuning curves imply that a small (large) fraction of neurons is active for a given stimuli parameter. Here, we assume that such neural populations are geometrically ordered according to the neuron's preferred parameter value. Then, the neural activity resembles a localized bump activation like experimentally observed in the compass system of the drosophila fly (Seelig and Jayaraman, 2015; Kim et al., 2017) and theoretically studied in continuous attractor models (Wilson and Cowan, 1973; Amari, 1977; Ben-Yishai et al., 1995; Skaggs et al., 1995; Seeholzer et al., 2017). In this coding scheme, which we call *bump coding scheme*, the center of a bump activation corresponds to the parameter value encoded by the bump activation (see **Figure 1A**), and the width of the bump is determined by the number of active neurons in the bump. Further, we assume that neural populations are equipped with a *continuous attractor mechanism* that ensures that only one bump is active at a time. Continuous attractor mechanisms are an established model of cortical working memory (Seeholzer et al., 2017), and emerge from the wiring motive of local excitation and long range inhibition (Kim et al., 2017). As another bio-plausible ingredient, we use *probabilistic synapses*. We assume a simple synaptic model

consisting of a plastic synaptic probability p and a synaptic weight w , which we fix to 1 in order to concentrate on our main ideas. The synaptic probability corresponds to the pre-synaptic neuro-transmitter release probability and the weight w to the post-synaptic quantal amplitude (Llera-Montero et al., 2019). The neuro-transmitter release probability of synapses in the brain is highly variable and typically between 0.1 and 0.9 (Branco and Staras, 2009).

We explain now with a sample application, how these three bio-inspired principles are integrated into a reinforcement-type (Williams, 1992), gradient-free learning mechanism. Assume that a robot arm with two joints should learn to reach given target positions $x = (x_1, x_2)$ by applying the correct angles $\hat{y} = (\alpha, \beta)$ to the joints, such that they reach the given target position x . Consider two populations of neurons A and B connected by probabilistic synapses. Population A and B use the bump coding scheme to encode the target position x and the angles y that are applied to the two joints, respectively, see **Figure 1B**. The goal is to adapt the synaptic probabilities such that every target position x is mapped to the correct angles \hat{y} . A bump activation encoding x in population A is propagated via probabilistic synapses to population B , where an abstract continuous attractor mechanism ensures that a single bump remains active in B . Its center y encodes the applied angles (α, β) . Note that y usually varies from trial to trial since the synapses between population A and B are probabilistic. According to the final arm position, the network receives a scalar error feedback $L(x, y)$, that depends on the input x and the output y generated by the network. Then, the synaptic probabilities between the bumps in population A and B are updated depending to this error



feedback. The larger the bump width k is, the more synapses are between the two bumps, whose probabilities are all updated according to the same error feedback. In this way, the network exploits the continuity of the task for sample efficient learning.

Related to this work, populations of tuned neurons have been used to learn sensory-motor transformations (Bullock et al., 1993; Salinas and Abbott, 1995; Baraduc et al., 2001; Baraduc and Guigon, 2002; Sanger, 2003), extending and building on the investigation of radial basis networks conducted in the 1980s and 1990s (Klopfenstein and Sverdløve, 1983; Broomhead and Lowe, 1988; Sanger, 1991, 1997, 1998; Pouget et al., 1998). These studies use Hebbian plasticity mechanisms, that require simultaneous activation of inputs and target outputs, whereas we propose reinforce-type learning algorithms, that learn the correct outputs through exploration of the output space. Furthermore, the investigation of the relation between sample efficiency and tuning curve width is novel.

The main contributions of this paper are summarized as follows. We introduce a reinforce-type learning algorithm that exploits the bump coding scheme, abstract continuous attractor mechanisms and probabilistic synapses for sample efficient learning of general low dimensional mappings. We show theoretically and by simulations that if the bump width is static during learning, then large bump width improves sample efficiency but harms the final precision, whereas small width impairs sample efficiency but improves final precision. Benefits of both are accomplished, if the bump width is dynamically decreased during the learning progress. Moreover, we show that the obtained sample efficiency is asymptotically optimal up to a $\log n$ factor in the limit of large population size n . For low dimensional mappings, the bump coding scheme achieves similar performance as a multi-layer perceptron trained by the backpropagation algorithm (Rumelhart et al., 1985), and it outperforms a multi-layer perceptron trained by the reinforce algorithm (Williams, 1992). It also achieves competitive performance on low dimensional reinforcement learning environments. Finally, we relate our findings to experimental observations of decreasing tuning curve width during learning and conclude that our findings propose sample efficiency as a functional role of the tuning curve width.

2. RESULTS

Assume that a network consisting of populations A and B should learn a mapping $f: [0, 1]^{d_A} \rightarrow [0, 1]^{d_B}$, e.g., mapping target position (\hat{x}_1, \hat{x}_2) to joint angles (α, β) as illustrated in the robotic arm task of the introduction. We consider general mappings f , that are only restricted to be Lipschitz continuous¹. This general framework, can be applied to many tasks including reinforcement learning as demonstrated in section 2.2. Populations A and B , are connected by probabilistic synapses and encode input x and output y , respectively, using a bump coding scheme, see **Figure 1** and section 3.1 for a formal

¹A function f is Lipschitz continuous with Lipschitz constant C , if $|f(x) - f(y)| \leq C|x - y|$ for all x, y . Intuitively, this is the case if the slope of f is everywhere smaller than C .

description. The goal is to learn the plastic synaptic probabilities, whereas synaptic weights are assumed to be fixed. The implicit goal of our learning algorithm is that a neuron x in population A keeps all synapses to the neurons y in B for which $|y - f(x)|$ is small and decreases the synaptic probabilities of all other synapses, see **Figure 1C**. This will ensure that a bump x in population A activates a bump with center y close to $f(x)$ in population B .

We begin by stating our theoretical results in section 2.1 before presenting the results obtained by simulations in section 2.2. For the theoretical results, we use a simplified version of the algorithm used in the simulations, because it allows a rigorous mathematical analysis and it illustrates the conceptual ideas of the algorithm. Both algorithm use the same basic principles and behave qualitatively the same.

2.1. Theoretical Results

We consider the following learning mechanism with fixed bump width k involving synaptic counters that are initialized with 0. We refer to a neuron with preferred parameter value x as neuron x and to a bump with center x as bump x . For every sample, a random input bump x is activated. The probabilistic synapses propagate the activity in A to population B , where an abstract continuous attractor mechanism activates the bump y in B that received highest synaptic input. The algorithm receives a scalar error feedback $L(x, y)$ that depends on input x and the output y , e.g., the euclidean distance between output y and target output $f(x)$. Then, the counters of the synapses between the two active bumps are decremented by 1 if the error feedback $L(x, y)$ is larger than some error threshold \hat{L} and incremented otherwise, see **Figure 1C**. After observing proportional $(\frac{n}{k})^{d_A+d_B} \log n$ many samples, we prune all synapses with non-positive counters, that is, we set their synaptic probabilities to 0. For a formal description of the algorithm, we refer to section 3.

We define the *error* of a learned network to be the expected error feedback $\mathbb{E}[L(x, y)]$ if the input is randomly chosen. If the mapping f is Lipschitz continuous and the network obtains the euclidean distance $L(x, y) = \|y - f(x)\|_2$ as error feedback, the following theorems hold.

Theorem 1 (Static bump width k). *The learning algorithm with static bump width k and euclidean error feedback learns a mapping $f: [0, 1]^{d_A} \rightarrow [0, 1]^{d_B}$ with error smaller than $\frac{3k}{n}$ after proportional to $(\frac{n}{k})^{d_A+d_B} \log n$ many random samples, where n is the population size.*

The approach with static bump width k ensures that each neuron x maintains the synapses to a small continuous interval of output neurons around value $f(x)$ and prunes away the other synapses. Thus, we can reapply the same learning mechanism, this time with smaller bump width k and smaller error threshold \hat{L} (for a formal description see section 3). Repeating this procedure will cause an input bump x to be mapped to a random bump y from a shrinking and shrinking interval around $f(x)$. This yields the following theorem.

Theorem 2 (Dynamic bump width k). *The learning algorithm with dynamic bump width and euclidean error feedback learns*

the mapping $f: [0, 1]^{d_A} \rightarrow [0, 1]^{d_B}$ with error smaller than ε after proportional to $\varepsilon^{-d_A} \log n$ many random samples, where n is the population size.

We conclude that in order to reach error $\varepsilon = \frac{3k}{n}$ the dynamic bump algorithm requires proportional to ε^{-d_B} times less samples than the static bump algorithm with bump width k . We remark that above results generalize to circular input and output spaces, that encode for example head direction or orientation angles, see the **Supplementary Material** for more details. Our proofs show that it is not necessary that above algorithms obtain the precise Euclidean distance as error feedback, but rather one bit of feedback suffices, if it indicates whether the Euclidean distance is larger than the error threshold \hat{L} or not. For such an algorithm, a lower bound on the sample efficiency can be obtained by an entropy argument.

Theorem 3 (Lower Bound). *For any algorithm that obtains only a single bit of feedback per sample, there are Lipschitz continuous mappings $f: [0, 1]^{d_A} \rightarrow [0, 1]^{d_B}$, such that the algorithm requires at least proportional to ε^{-d_A} many samples to learn f with error smaller than ε .*

Then, Theorem 2 and 3 imply that the learning mechanism with dynamic k accomplishes a sample efficiency that is asymptotically optimal up to a $\log n$ factor. For the proofs of these theorems, we refer to the **Supplementary Material**.

2.2. Empirical Results

For the simulations, we use a slightly more sophisticated learning algorithm that follows the same underlying principles, but differs from the algorithm that we analyze theoretically in five aspects. Firstly, it is designed to handle more general error feedback functions. For example, in the robotic arm task of the introduction, the error feedback is not given by the euclidean distance between the output angles, but by the distance between the reached position and the target position. In turn, the magnitude of the error feedback can change for different inputs, and a single error threshold \hat{L} will not allow fast learning for all inputs. To resolve this issue, we assume that every input neuron i keeps track of a running mean \hat{L}^i of the error feedbacks that were obtained when i was active. For the update of the outgoing synapses of neuron i , \hat{L}^i is compared to the error feedback L . Secondly, if $L \geq \hat{L}^i$ for a neuron i of the bump in population A , the synapses projecting from neuron i on the bump in population B are pruned away immediately. Thirdly, for the dynamic case, the bump width k is adapted continuously instead of repeatedly applying the static algorithm. Since the learning progress might vary for different input regimes, we allow k to depend on the input, and we set k for input x proportional to the number of outgoing synapses of neuron x . Note that this number is a reasonable measure of how well input x is already learned as the precision of the output depends on the magnitude of the interval of synapses connecting x to neurons around $f(x)$. Fourthly, long-time inactive synapses are pruned away, i.e., they are pruned if the post-synaptic neuron has not been active for a couple of times when the pre-synaptic neuron was active. Finally, synapses

of neuron i are consolidated (that is its probability is set to 1) if its number of synapses drops below a certain threshold value. We call this algorithm the *dynamic bump algorithm*. We refer to section 3 for a formal description of the algorithm and to **Figure 2** for an illustration of the evolution of the synaptic probabilities.

The *static bump algorithm* works analogously to the dynamic bump algorithm, except that k is held constant during the whole algorithm. **Figure 3** empirically confirms the trade-off between sample efficiency and final performance for static bump width k . Larger k leads to faster learning compared to smaller k , however reaches worse final error. The advantages of both large and small k can be exploited by adapting the bump width dynamically during the learning process, see **Figure 3**.

In order to put the sample efficiency of the bump coding scheme into context with other coding schemes, we compare it to the performance of a multi-layer perceptron (MLP), which encodes information with real valued units. **Figure 4** compares performance of the dynamic bump algorithm, with a MLP trained by the backpropagation algorithm (Rumelhart et al., 1985) and a reinforce algorithm as described in Williams (1992). Note that the backpropagation algorithm requires full access to the first order derivative of the error with respect to the parameters and thus is a first-order optimization technique, whereas the reinforce and dynamic bump algorithm only require a scalar error feedback and thus are zeroth-order optimization techniques. Nevertheless the dynamic bump algorithm achieves similar performance as the MLP trained by backpropagation and outperforms the MLP trained by the reinforce algorithm, **Figure 4**. For the backpropagation and reinforce algorithm, we used a hyper-parameter search to determine the best parameters. We note that this search yielded an untypically high learning rate and small batch size for the backpropagation algorithm. The learning rate is in the upper end of the recommended interval $[10^{-6}, 1]$ and much higher than the suggested default value of 0.01 (Bengio, 2012). This is necessary to achieve good performance after 1,000 samples, see **Figure 4**.

In **Figure 5**, we illustrate the performance of the bump coding scheme on reinforcement learning (RL) tasks. In RL environments an agent should learn to interact with an environment with the goal of maximizing some reward. At any time step, the agent observes the current state of the environment and outputs an action, which in turn affects the state of the environment. The agent obtains rewards for reaching certain states. It is unclear which actions lead to the reward, due to the well known credit-assignment problem. A classical RL method to mitigate this problem is the temporal difference learning method (Sutton et al., 1998), that relies on learning a *policy function* that maps states to actions and a *value function* that maps states to an estimate of the future expected reward. Then, the difference of the estimated expected future reward before and after each action can be computed. Combined with the obtained reward of that time step, one can estimate the reward that arose from that specific action, which allows to update the policy.

Our *RL bump algorithm* uses the temporal difference learning method. More precisely, it learns the policy with a static bump algorithm, while the value function is stored in tabular representation as done in the literature (Sutton et al., 1998).

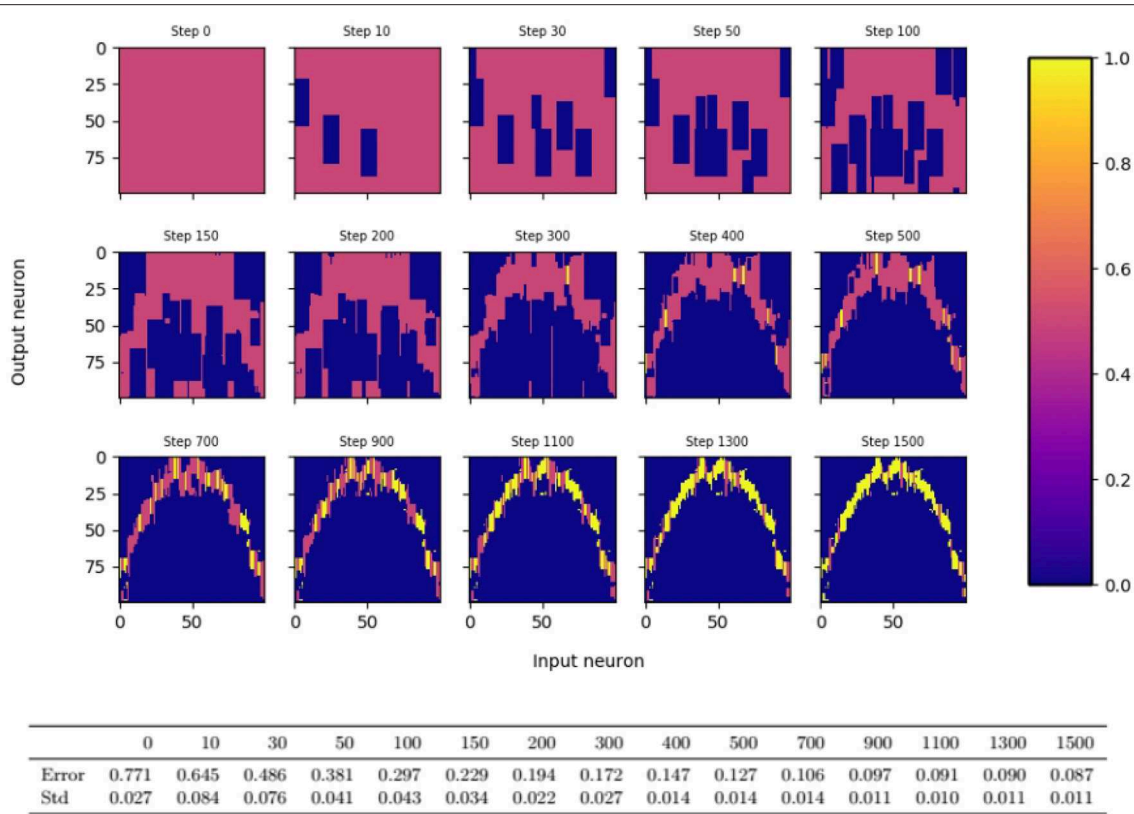


FIGURE 2 | Evolution of synaptic probabilities when learning a one dimensional mapping $f(x) = x^2 - 3x + 1$ with the dynamic bump algorithm. Each small color plot displays the synaptic probabilities, where input neuron is on the x-axis and output neuron on the y-axis. Blue areas visualize pruned synapses, yellow areas visualize consolidated synapses. Input and output population consist of 100 neurons each and the output bump is 3 times as large as the input bump.

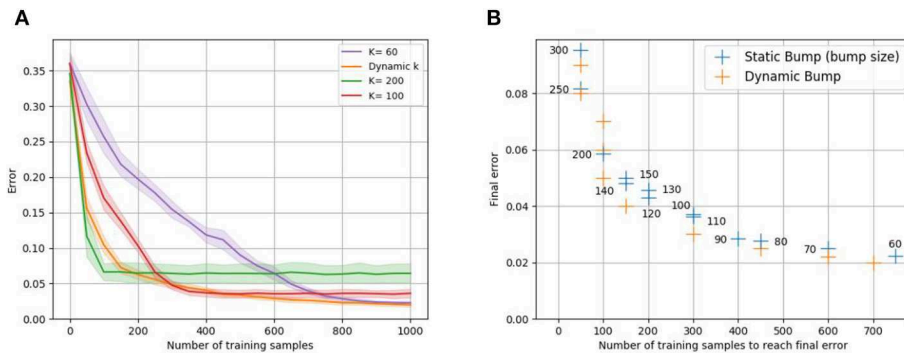


FIGURE 3 | Sample efficiency vs. final error trade-off for static bump width k . Recall that the error is defined to be the expected error feedback if the input is randomly chosen. Plot (A) shows the evolution of the error of the static bump algorithm for different k , when learning the one-dimensional identity mapping $f(x) = x$ with euclidean error feedback. The dynamic bump algorithm, is labeled as “dynamic k .” Plot (B) shows the minimal error against the number of training samples required to reach this error for different k (given as number next to the blue data points); to avoid taking into account the slow progress before final convergence, the number of samples required to achieve 1.5 times the final error is plotted. Analogously, we plot the number of samples required by the dynamic bump algorithm to achieve 1.5 times the shown error values. For both plots, populations A and B consist of 1,000 neurons each, and the mean of 10 trials is plotted.

At any time step, the temporal difference learning method provides an estimate of the reward arising from the action at this time step. This estimate is used as error feedback for the bump algorithm. Since this estimate might be off, we update the

synaptic probabilities more gradually, instead of pruning away synapses immediately as in above algorithms. The magnitude of the updates are chosen proportional to the reward estimates. We refer to section 3.4 for a detailed description.

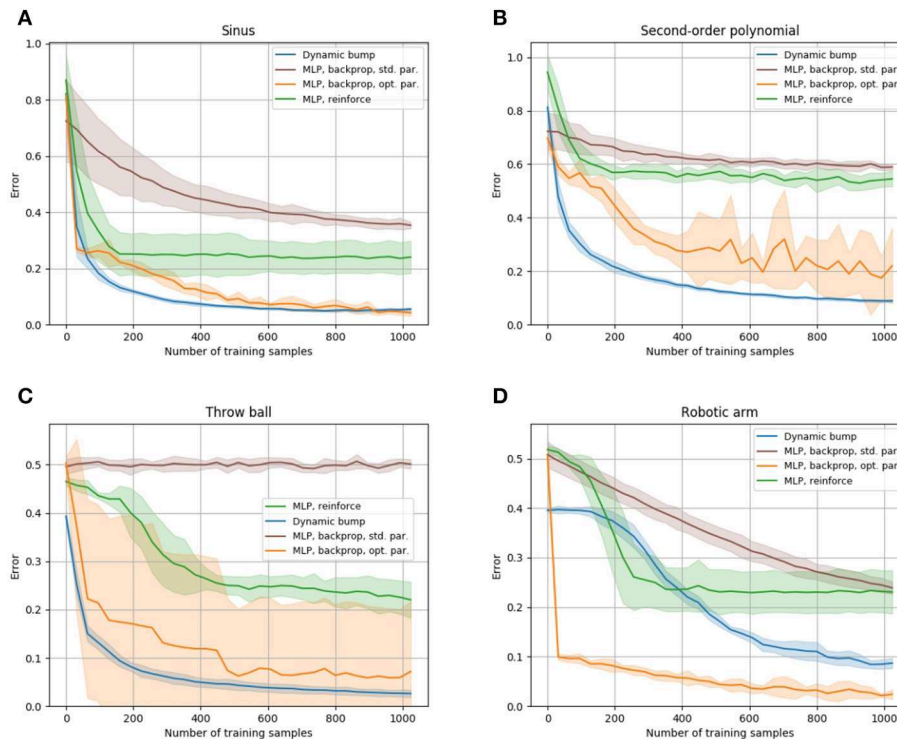


FIGURE 4 | Comparison of the sample efficiency to a multi-layer perceptron. The learning progress of the dynamic bump algorithm, of a MLP trained with the backpropagation algorithm and of a MLP trained with the reinforce algorithm is shown for the sinus, second-order polynomial, throw ball and robotic arm tasks in (A–D), respectively. The tasks are described in detail in section 3.5. We plot the learning curve that achieves best performance after 1,000 samples obtained by our hyper-parameter search. The hyper-parameters for each learning algorithm are given in the **Supplementary Material**. For each algorithm, we plot average and standard deviation of 10 runs.

We evaluate the RL bump algorithm on two RL tasks, the Mountain Car environment and the Inverted Pendulum, see section 3.5 for a description of the tasks. The resulting learning curves are displayed in **Figure 5**. We compare our algorithm with two different deep reinforcement learning algorithms: *deep deterministic policy gradient* (ddpg) (Lillicrap et al., 2015) and *proximal policy optimization* (ppo2) (Schulman et al., 2017). Both are state-of-the-art reinforcement learning algorithms that use neural networks as a representation of the policy. For both algorithms, we use the implementation and the standard parameters provided by the OpenAI Baselines (Dhariwal et al., 2017). Since these parameters do not aim at efficiently solving the exploration problem of the Mountain Car environment, we added an implementation of ddpq specifically tailored to achieve sample efficiency in the Mountain-Car environment: the parameters and implementation are taken from de Broissia and Sigaud (2016) and the resulting curve is labeled as *sample-efficient ddpq*.

We observe that our algorithm performs comparable or better than the baselines on both environments, see **Figure 5**. In the Mountain Car experiment, we observe that the OpenAI baselines implementation are unable to make substantial progress on the observed time scale. The sample efficient ddpq implementation from de Broissia and Sigaud (2016) is able to reach higher rewards, but it is very quickly outperformed by our algorithm

in terms of final performance. In the Pendulum experiment, we observe a first phase during which the learning curve of our algorithm is very similar to ddpq, whereas ppo2 does not make any progress. In a second phase, the ppo2 learning curve catches up our algorithm, while ddpq is outperformed. We close this section with a word of caution. **Figure 5** seems to indicate that our algorithm outperforms deep policy gradient methods for reinforcement learning tasks. However, note that both, Mountain Car and Inverted Pendulum, have a low dimensional input space (2 dimensions for Mountain Car, 4 for Inverted Pendulum). Currently, our algorithm does not scale up to a higher number of dimensions in terms of computational cost, whereas deep policy gradient algorithms have been engineered to deal with high-dimensional spaces.

3. METHODS

In sections 3.1–3.4, we describe the bump coding scheme and our algorithms formally. Section 3.5, contains a description of all the tasks used for evaluation of the algorithms.

3.1. Bump Coding Scheme

We assume that a population of binary neurons is arranged in a grid. Intuitively, a d -dimensional parameter x is encoded by a bump of active neurons that lay within the d -dimensional cube

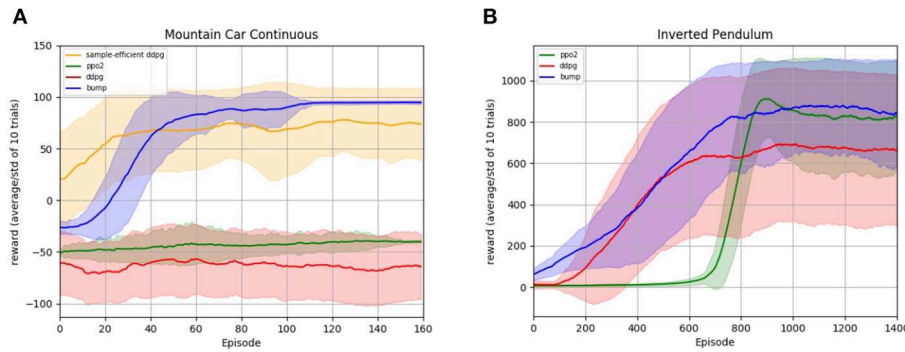


FIGURE 5 | Reinforcement learning experiments. The plots show the learning progress of different learning algorithms on the Mountain Car **(A)** and Inverted Pendulum **(B)** task. The tasks are described in more detail in section 3.5. The RL bump algorithm is compared with the *deep deterministic policy gradient* algorithm (ddpg) and the *proximal policy optimization* algorithm (ppo2), cf. text for implementation details. For every algorithm, the average reward per episode and its standard deviation for 10 different random seeds are smoothed for better readability. The algorithm hyper-parameters are optimized to maximize the mean reward on the last 30 episodes for the Mountain Car, and the last 300 episodes for the Inverted Pendulum, and are given in the **Supplementary Material**.

with side length k and center x , see **Figure 1**. We note that our results qualitatively do not depend on the precise shape of the bump, that is, whether it is the d -dimensional cube or ball, however our cube shaped bumps facilitate efficient simulations. Formally, a population of n^d neurons encodes values in the d -dimensional interval $[0, 1]^d$. Define $A_i = [n]$, where $[n]$ denotes the set of integers $\{1, \dots, n\}$, and define the set of neurons to be $A = A_1 \times \dots \times A_d$. Then, a neuron $i \in A$ represents the value $(i_1/n, \dots, i_d/n)$ in $[0, 1]^d$. Moreover, for $a \in [n]$ define $Int_k(a)$ to be the set of integers in the interval $[a - k/2, a + k/2]$; to be precise, we actually define $Int_k(a)$ as the set of integers in the interval $[\max\{0, a - k/2\}, \min\{n, a + k/2\}]$ in order to take care of cases close to the boundary of the interval. We define the *bump of width k with center neuron $i \in A$* to be $Int_k(i) := Int_k(i_1) \times \dots \times Int_k(i_d)$. In the bump coding scheme of width k the value a is encoded by the $Int_k(a)$ -activation, where an I -activation is defined to be the state, where the neurons in I are active and the ones not in I are inactive.

3.2. Network Architecture, Activation Distribution, and Feedback Error Measure

In order to learn a mapping $f: [0, 1]^{d_A} \rightarrow [0, 1]^{d_B}$, we consider a network consisting of populations A and B equipped with probabilistic synapses and an abstract continuous attractor mechanism in population B . Intuitively, given a bump activation I in population A , the attractor mechanism activates the bump J in population B that received most synaptic input from the bump activation in A . This way, the probabilistic synapses enable exploration of the output space. More formally, the network consists of population A and B consisting of $[n]^{d_A}$ and $[n]^{d_B}$ neurons, respectively. These encode values according to the bump coding scheme defined above. They are fully connected by probabilistic synapses with weights w_{ij} fixed to value 1 and plastic synaptic probabilities p_{ij} . Given an input x , the bump $I = Int_{k_A}(x)$ with width k_A and center x , the matrix of synaptic probabilities P and the bump width k_B in population B , we define the following *activation distribution* $Act(I, k_B, P)$ that returns (J, y) , where J is

the sampled bump in population B with center y . Assuming that the bump I is active in population A , we explain now how a bump activation J in population B is sampled. Denote by $\mathcal{J} = \{Int_{k_B}(b) | b \in B\}$ the set of all possible bump activations in B with width k_B and by X_{ij} Bernoulli random variables that are 1 with probability p_{ij} and 0 otherwise, i.e., they indicate whether synapse (i, j) is active or not. Then, $s(I, J) = \sum_{(i, j) \in I \times J} X_{ij}$ is the total synaptic input that neurons in J receive (recall that we assumed that all synaptic weights are constant and equal to 1), and the *abstract continuous attractor mechanism* activates the output bump $J \in \mathcal{J}$ with maximal $s(I, J)$, where ties are broken uniformly at random. Formally, we write $(J, y) \sim Act(I, k, P)$, where $J = \arg \max_{J \in \mathcal{J}} \{s(I, J)\}$ and y is the center of bump J . In a machine-learning context this can be efficiently implemented by adding a convolutional layer with suitable weights on top of layer B . We remark that for the theoretical analysis, we change the activation distribution slightly to be able to deal with the dependencies between distributions $s(I, J)$ and $s(I, J')$, see **Supplementary Material**.

Given a bump I in population A with center x and a sampled activation $(J, y) \sim Act(I, k_B, P)$ the network receives some *error feedback* $L(x, y)$, where L is a function that depends on the output y and the target output $f(x)$, e.g., the euclidean error feedback returns the euclidean distance between output value y and the target value $f(x)$. Note that the learning task is more difficult if the precise definition of the function L is unknown to the algorithm. The *error* of a network with learned synaptic probabilities \hat{P} is defined to be the expected error feedback $\mathbb{E}[L(x, y)]$, where x is chosen uniformly at random in A and y is sampled according to $Act(Int_{k_A}(x), k_B, \hat{P})$.

3.3. The Static and Dynamic Bump Algorithm

We first explain the algorithms used for the theoretical analysis and then the ones used for the simulations. The goal of these algorithms is to prune away all synapses for every input neuron x , except the ones connecting x to some small continuous area

around x 's target neuron $f(x)$. The basic mechanism to do so is to prune away synapses between the input and output bump whenever the error feedback is larger than some error threshold, as then the target neuron $f(x)$ is not contained in the output bump.

3.3.1. Algorithms for Theoretical Analysis

For the following algorithms, every synapse (i, j) has a synaptic counter c_{ij} that indicates whether a synapse should be pruned away. Further, the input and output bump widths are set such that the fraction k_B/k_A is equal to the Lipschitz constant C of the mapping f that is to be learned.

The *static bump algorithm (theory)* fixes the input width k_A , the output bump width k_B and the error threshold \hat{L} proportional to the desired final error ℓ and observes $M = c \left(\frac{n}{k}\right)^{d_A+d_B} \log n$ many samples with random input x , where $c = (d_A + d_B + 1)(2C)^{d_A}(\sqrt{d_A} + \sqrt{d_B})^{d_A+d_B}$. For every sample activation $(J, y) \sim \text{Act}(\text{Int}_{k_A}(x), P, k_B)$, the counters c_{ij} are incremented by 1 for all synapses between the two bumps $I = \text{Int}_{k_A}(x)$ and J , if $L(x, y) \leq \hat{L}$ and otherwise decremented by 1. Finally, all synapses with non-positive counters are pruned away. Intuitively, the choice of sample size M ensures with high probability² that any input neuron x remains to be connected to output neurons y with $\|f(x) - y\|_2 \leq \hat{L}$.

The *dynamic bump algorithm (theory)* proceeds in phases and repeatedly applies the static version, see Algorithm 1. The bump widths k and error threshold \hat{L} are initially chosen large and divided by 2 in every phase. Intuitively, this causes any input neuron x to be connected to a shrinking and shrinking area around target output neuron $f(x)$.

3.3.2. Algorithms for Simulations

The following algorithms can deal with error feedback functions that differ in magnitude for different input regimes and adapt the bump width in a more continuous manner. We give a short overview of the mechanisms used in the algorithms. Firstly, every neuron keeps track of its own error threshold \hat{L}^i . \hat{L}^i is the running average with decay factor α of the error feedbacks obtained when neuron i was active. Further, a mechanism to prune long-time inactive synapses is implemented with synaptic counters d_{ij} .

The *dynamic bump algorithm* sets for every sample the input bump width k_A and output bump width k_B both equal to a constant fraction of the number of synapses of neuron x . The *static bump algorithm* sets the bump widths to some fixed constant but otherwise proceeds analogously as follows. Any sample activation consists of input and output bumps I and J with centers x and y , respectively, and error feedback $L(x, y)$. For each sample, the neurons $i \in I$ update their error threshold according to $\hat{L}^i = \alpha L(x, y) + (1 - \alpha)\hat{L}^i$. Synapses (i, j) between bumps I and J are pruned away if $L(x, y) \geq \hat{L}^i$. Further, synapses (i, j) with $i \in I$ and $j \notin J$ increase their synaptic counter d_{ij} by 1 and the ones with $i \in I$ and $j \in J$ reset $d_{ij} = 0$. Then, long-time inactive synapses (i.e., $d_{ij} \geq \theta_{prune}$) are pruned away. Finally, synapses (i, j) are consolidated, that is, p_{ij} is set to 1, if the number

of synapses of input neuron i drops below a threshold value. The procedure stops as soon the mean of the \hat{L}^i drops below the desired precision ℓ .

The hyper-parameters of the algorithm need to be tuned to yield good performance. We optimized them using a coarse grid search for each task. We discuss the influence of the hyper-parameters on performance and give the hyper-parameters used for **Figure 4**.

Algorithm 1: Dynamic bump algorithm for theoretical analysis. It learns a Lipschitz mapping $f: [0, 1]^{d_A} \rightarrow [0, 1]^{d_B}$. Hyper/parameters: Lipschitz constant or upper bound on Lipschitz constant C , desired precision ℓ , $0 < c_k \leq 2/(3(\sqrt{d_A} + \sqrt{d_B}))$, $c = 2^{d_A}(d_A + d_B + 1) \cdot \text{Vol}(S_{d^B}(2/c_k + \sqrt{d_A}/2 + \sqrt{d_B}/2))$, where $\text{Vol}(S_d(r))$ denotes the euclidean volume of the d -dimensional ball with radius r .

```

1  $\hat{L}_0 = 1$ 
2 while  $\hat{L}_s \geq \ell/2$  do
3    $\hat{L}_s = \hat{L}_{s-1}/2$ ;    $k_s^B = c_k \hat{L}_s$ ;    $k_s^A = k_s^B/C$ ;    $c_{ij} = 0$ 
   for all  $(i, j) \in A \times B$ ;
4   for  $M = c \left(\frac{n}{k_s^A}\right)^{d_A} \log n$  samples do
5      $x \in_{u.a.r.} A$ ;    $I = \text{Int}_{k_s^A}(x)$ ;    $(J, y) \sim \text{Act}(I, P, k_s^B)$ ;
6     if  $L(x, y) \geq \hat{L}$  then  $c_{ij} = c_{ij} - 1$  for all  $(i, j) \in I \times J$ ;
     else  $c_{ij} = c_{ij} + 1$  for all  $(i, j) \in I \times J$ ;
7   for  $(i, j) \in I \times J$  do
8     if  $c_{ij} \leq 0$  then  $p_{ij} = 0$ 
```

Algorithm 2: Dynamic bump algorithm used for the simulations. It learns a Lipschitz mapping $f: [0, 1]^{d_A} \rightarrow [0, 1]^{d_B}$. Hyper/parameters: Lipschitz constant or upper bound on Lipschitz constant C , desired precision ℓ , running average factor α , c_A , c_B , θ_{prune} , θ_{syn} , \hat{L}_{init} .

```

1 for  $i \in A$  do  $\hat{L}^i = \hat{L}_{init}$ 
2 while  $\frac{\sum \hat{L}^i}{n} \geq \ell$  do
3    $x \in_{u.a.r.} A$ ;    $k_B = \text{number-synapses}(x)/c_B$ ;
    $k_A = \text{number-synapses}(x)/c_A$ ;
4    $I = \text{Int}_{k_A}(x)$ ;    $(J, y) \sim \text{Act}(I, P, k_B)$ ;
5   for  $i \in I$  do
6      $\hat{L}^i = \alpha L(x, y) + (1 - \alpha)\hat{L}^i$ ;
7      $d_{ij} = 0$  for all  $j \in J$ ;    $d_{ij} = d_{ij} + 1$  for all  $j \notin J$ ;
8     if  $L(x, y) \geq \hat{L}^i$  then  $p_{ij} = 0$  for all  $j \in J$  for  $j \in B$  do
9       if  $d_{ij} \geq \theta_{prune}$  then  $p_{ij} = 0$ 
10    if  $\text{number-synapses}(i) \leq \theta_{syn}$  then
      Consolidate synapses  $(i, j)$  of neuron  $i$ 
```

3.4. Reinforcement Learning Bump Algorithm

In this section, we describe our *RL bump algorithm*, which combines the classical temporal difference learning method (Sutton et al., 1998) with the bump coding scheme. The building

²With high probability means with probability tending to 1 as n tends to ∞ .

blocks of temporal difference learning are learning a *policy* that maps states to actions and learning a *value function* v that estimates the expected future reward $v(S_t)$ for states S_t . Here, it is assumed that the agent observes the whole state S_t of the environment. To learn these mappings, one computes the *return* G_t , that intuitively is the difference between the estimated future reward $v(S_t)$ and the real future reward. The k -step temporal difference method estimates the real future reward using the bootstrap estimate $v(S_{t+k})$ and computes G_t as

$$G_t \leftarrow R_{t+1} + \dots + R_{t+k} + v(S_{t+k}) - v(S_t),$$

where R_u denotes the reward received at time u . Then, the return G_t is used to update the value function, which is stored in a tabular representation (Sutton et al., 1998).

$$v(S_t) \leftarrow v(S_t) + \alpha G_t,$$

where α is a hyper parameter regulating the magnitude of the update. Further, the return G_t is used to update the policy, that in our case is learned with a bump coding algorithm with static bump width that updates the synaptic probabilities gradually. If at time step t , the activated bumps were I and J , we update the probabilities p_{ij} for all $(i, j) \in I \times J$ according to

$$p_{ij} \leftarrow \begin{cases} p_{i,j} + (0.9 - p_{i,j}) \min(\beta G_t, 1) & \text{if } G_t \geq 0 \\ p_{i,j} + (p_{i,j} - 0.1) \max(\beta G_t, -1) & \text{if } G_t < 0 \end{cases},$$

where β is a hyper parameter regulating the magnitude of the update. Intuitively, this update increases p_{ij} proportional to G_t and proportional to the distance of p_{ij} to 0.9 if $G_t \geq 0$, and decreases p_{ij} proportional to G_t and proportional to the distance of p_{ij} to 0.1 if $G_t \leq 0$. Note that due to the clipping of βG_t to $[-1, 1]$, the invariant is maintained that all probabilities belong to the interval $[0.1, 0.9]$. The structure of these updates is very similar to policy gradient methods (Sutton et al., 1998), using the reinforce algorithm. However, the reinforce algorithm does not directly apply to our model as it is internally non differentiable, see Williams (1992).

3.5. Description of Tasks

3.5.1. Low Dimensional Mappings With Immediate Feedback

As 1-dimensional mappings f , we consider the identity function $f(x) = x$, a sinus function $f(x) = \sin(x)$ and a second order polynomial $f(x) = x^2 - 3x + 1$. For all functions, we consider the absolute distance $L(x, y) = |f(x) - y|$ as error feedback, where x and y are input and output of the network.

In the *throw ball task*, the network has to learn to throw a ball to a certain distance. The target distance is given as a 1-dimensional input to the network. The networks gives a 2-dimensional output, that consists of the vertical throwing angle and the initial speed of the ball. The error feedback is the absolute difference between the distance where the ball touches the ground and the target distance. Note that the optimal output is underdetermined, as for any angle in $(0, \pi/2)$ there exists a speed such that any target distance can be hit.

In the *robotic arm task*, the network learns how to control a simple robotic arm with two degrees of freedom. The arm is composed of two rigid moving parts and is connected to a fixed anchorage point. As we restrict to movements in the plane, the agent only has to control angles, one at each joint. The target position is given as input in Cartesian coordinates and the network outputs two angles that are applied to the two joints. The feedback is given by the euclidean distance between the actual position of the arm and the target position.

3.5.2. Reinforcement Learning Tasks

The RL environments used for assessing the RL bump algorithm are from the OpenAI Gym (Brockman et al., 2016) toolkit: the continuous version of the classical Mountain Car control problem (MountainCarContinuous-v0), and the Inverted Pendulum environment from the MuJoCo suite (InvertedPendulum-v2). In the *mountain car task* the goal is to reach the top of a hill, that can only be reached by obtaining momentum when driving down the neighboring hill. The network receives as input the position and speed of the car and outputs the acceleration that is applied to the car. There is a large positive reward if the car reaches the top of the hill and a small negative reward for the fuel use in every time step. In the *inverted pendulum task* the goal is to balance an inverted pendulum on a cart. The network receives 4-dimensional input describing position and velocity of the cart and pendulum, and it outputs the acceleration applied to the cart. As long as the pendulum does not fall to the ground, there is a positive reward in every time step.

4. DISCUSSION

In this section, we first relate our work to related work from the field of machine learning that studies sample efficient learning, then we discuss the bio-plausibility of our proposed coding scheme and learning algorithms, and finally we discuss the insights gained about the bio-plausible mechanisms used in this study.

4.1. Sample Efficient Learning in Machine Learning

How to improve sample efficiency of learning algorithms is a major topic in the field of machine learning in general and the field of reinforcement learning in particular (Botvinick et al., 2019). Data samples are often limited, and training of artificial networks is computationally costly. A common approach to improve sample efficiency is to handcraft artificial networks to the task at hand. The most famous example are convolutional neural networks (LeCun and Bengio, 1995; Krizhevsky et al., 2012), where the translational invariance property of images is hand-wired into the convolutional network architecture. Another successful approach is to store all observed samples or the neural states that encode these samples. Then, inputs are classified according to the most similar samples in storage. This idea is present in non-parametric approaches such as the nearest neighbor methods (Bishop, 2006), as well as in deep neural networks augmented with external memory systems (Graves

et al., 2014) and attention mechanisms (Bahdanau et al., 2014; Vaswani et al., 2017). In reinforcement learning this idea is known as episodic reinforcement learning (Lengyel and Dayan, 2008; Blundell et al., 2016; Gershman and Daw, 2017; Pritzel et al., 2017). Further, an approach to improve sample efficiency is meta learning (Schaul and Schmidhuber, 2010), which is also often referred to as “learning to learn.” In the meta learning setting, an outer learning system adjusts parameters or learning mechanisms of an inner learning system in order to improve the performance and efficiency of the later (Schmidhuber et al., 1996; Baxter, 1998; Thrun and Pratt, 1998; Hochreiter et al., 2001; Schweighofer and Doya, 2003). The outer learning system usually performs updates in a slow timescale, whereas the inner learning system can adapt fast to new environments, e.g., evolutionary algorithms can optimize learning architectures or loss functions to improve their sample efficiency (Stanley and Miikkulainen, 2002; Jaderberg et al., 2018).

Our approach is orthogonal to all these approaches. In essence, our work shows that the coding scheme of a network affects its sample efficiency and that adapting the coding scheme during learning can improve its sample efficiency.

4.2. Bio-plausibility of Our Coding Scheme and Learning Mechanisms

The proposed coding scheme and learning algorithms are of abstract nature and we do not intend to argue that they might be implemented in biological systems precisely in this form. However, we do claim that neural implementations of the basic concepts used by our model are plausible and the brain might use similar mechanisms for computation.

Our primary assumption that information is encoded and processed by populations of tuned neurons is supported by the abundance of such neurons across brain areas (Hubel and Wiesel, 1962; O’Keefe, 1976; Knudsen and Konishi, 1978; Ranck, 1985; Georgopoulos et al., 1988). The bump coding scheme described in section 3 requires that geometrically close-by neurons have close-by preferred stimuli parameters. Such geometrically ordered networks are indeed present in real neural network, such as the drosophila fly compass system (Seelig and Jayaraman, 2015; Kim et al., 2017). The underlying network-wiring that gives rise to bump-like activation patterns is generally believed to follow the circuit-motive of local excitation and long-range inhibition, as suggested by experimental evidence (Kim et al., 2017) and theoretical findings (Wilson and Cowan, 1973; Amari, 1977; Ben-Yishai et al., 1995). Note however, that the geometrical ordering of the neurons is not necessary for the results presented in this work. Indeed, the geometrical arrangement can be arbitrary if network-wiring supports activity patterns consisting of neurons with similar preferred stimuli parameters. Such network wiring consists of excitatory connections between neurons that are active for similar stimuli and inhibition that limits the total activity. It can be found across brain areas and animal species (Weliky et al., 1995; Mysore et al., 2010; Ko et al., 2011), and is often assumed by theoretical studies (Ben-Yishai et al., 1995; Skaggs et al., 1995; Knierim and Zhang, 2012). It is conceivable that such experimentally observed wiring motives implement a

version of the abstract continuous attractor mechanism used in this paper.

The dynamic bump algorithm requires a dynamic adaptation of the bump width during the learning process. Experimental and theoretical studies give evidence that the tuning curve width is controlled by inhibition (Suga et al., 1997; Knierim and Zhang, 2012; Lee et al., 2012). Thus, controlling the strength of inhibition in the system yields a straight forward explanation of how the bump width could be adjusted during the learning progress.

Moreover, the assumption of constant weights and binary neurons are mere abstractions for mathematical simplicity. Due to the on-off nature of binary neurons, we approximated the bell shaped tuning curves by rectangular tuning curves. It seems plausible that the results would translate qualitatively to networks of rate neurons with bell shaped tuning curves. Stable bump like activity patterns also can be produced by spiking networks (Seeholzer et al., 2017). We leave extensions of our algorithms that are more bio-plausible for future investigations. Moreover, we note that all results from this work also hold if the populations are sparsely instead of fully connected. As long as the bump width is broad enough, sufficiently large population codes give rise to stable learning mechanisms for sparsely connected populations (Gauy et al., 2017).

Furthermore, our plasticity rules are plausible in the sense that they solely depend on pre- and post-synaptic activity, a global reward feedback and memory traces of these quantities. All the neuronal and synaptic counters used in our algorithms require only local storage of activity and reward feedback traces.

4.3. Functional Role of Tuning Curve Width

A large body of literature in theoretical and experimental neuroscience investigated tuning curve shape under the aspect of optimal coding (Seung and Sompolinsky, 1993; Brunel and Nadal, 1998; Panzeri et al., 1999; Eurich et al., 2000; Bethge et al., 2002, 2003; Todorov, 2002; Sanger, 2003; Harper and McAlpine, 2004; Johnson and Ray, 2004; Seriès et al., 2004; Lánský and Greenwood, 2005; Brown and Bäcker, 2006; Montemurro and Panzeri, 2006; Toyozumi et al., 2006; McDonnell and Stocks, 2008; Geisler et al., 2009; Nikitin et al., 2009; Yarrow and Seriès, 2015). Also the width of tuning curves was analyzed from an information theoretical viewpoint. Hinton et al. (1986) and Zhang and Sejnowski (1999) established a dependence between optimal tuning width and the dimensionality of the encoded parameter, Pouget et al. (1999) and Butts and Goldman (2006) showed that optimality of tuning width heavily depends on the level of noise and covariance of the noise in the system, and Yaeli and Meir (2010) found that optimal tuning width depends on the prior uncertainty and on the length of the decoding time window. Such studies can explain the sharpening of tuning curves that is observed in a variety of experimental set-ups (Spitzer et al., 1988; Wagner, 1990; Ringach et al., 1997; Menz and Freeman, 2003; Wang et al., 2005; Samonds et al., 2009).

In this work, we introduce sample efficiency as a novel notion of optimality. If the neural code is optimized for sample efficient learning, then the model analyzed in this paper predicts that the tuning curves sharpen during the process of learning. In fact, this phenomena is known to occur in the inferior temporal

cortex, where tuning curves of shape selective neurons sharpen during acquaintance to new objects (Booth and Rolls, 1998; Freedman et al., 2005), as well as in many sensory areas during development (Brugge et al., 1981; Tavazoie and Reid, 2000; Mrsic-Flogel et al., 2003). The precise relation between tuning curve width and sample efficiency likely depends on the applied plasticity mechanisms. Nonetheless, for any plasticity mechanism that requires pre- and post-synaptic activity, the tuning curve width yields an upper bound on the number of synaptic weight updates, because it limits the number of active neurons per sample. Therefore, the basic principle that larger tuning curve width leads to more synaptic updates per sample and thus faster learning, may apply to many plasticity mechanisms.

4.4. Functional Role of Probabilistic Synapses

The functional role of probabilistic synapses is highly debated (Llera-Montero et al., 2019). The proposed functional roles include regularization and improved generalization in deep neural networks (Wan et al., 2013; Blundell et al., 2015) and energy saving constraints (Harris et al., 2012). Further, probabilistic synapses can give rise to a good exploration exploitation trade-off in reinforcement learning (Seung, 2003; Blundell et al., 2015; Kappel et al., 2018), and synaptic sampling can be seen as sampling from some posterior distribution (Aitchison and Latham, 2015; Kappel et al., 2015, 2018). Our model is in line with the last two proposals. In our model, probabilistic synapses combined with an continuous attractor mechanism encode the uncertainty of the learned input-output mapping and implement the variability and exploration that is required for reward-based learning.

4.5. Conclusion

In this work, we asked how sample efficient learning is affected by the neural coding scheme. We showed that population codes with tuned neurons support sample efficient learning for low dimensional tasks with immediate reward feedback and low dimensional reinforcement learning tasks. For these tasks, our gradient-free learning algorithm is competitive to multi-layer perceptrons trained by backpropagation. These findings might inspire an integration of tuning curve coding schemes into

machine learning approaches, especially, if data-samples are limited and no access to gradient information is given. For our learning mechanisms, we found that tuning curve width severely influences the sample efficiency. We showed that for static tuning widths, there is a trade-off between sample efficiency and final precision. Broad tuning curves give rise to sample efficient learning, whereas narrow tuning curves account for high final precision. Moreover, we showed that dynamic adaptation of the tuning width results in both high sample efficiency and high final accuracy. These results propose sample efficient learning as a functional role of the tuning curve width.

DATA AVAILABILITY STATEMENT

The code used for the simulations of the learning algorithms is provided under the following link (https://github.com/rdang-nhu/Sample_Efficient_Tuning_Curves).

AUTHOR CONTRIBUTIONS

FM and AS development of model set-up and research question, design, development and analysis of learning algorithms, and writing of the manuscript. RD-N implementation and simulation of the learning algorithms, development of the RL-bump algorithm, and proofreading of the manuscript.

FUNDING

Research supported by grant no. CRSII5173721 of the Swiss National Science Foundation.

ACKNOWLEDGMENTS

We thank Ulysse Schaller for preliminary work on the model set-up and research question during his master thesis.

SUPPLEMENTARY MATERIAL

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/fncom.2020.00012/full#supplementary-material>

REFERENCES

- Adrian, E. D., and Zotterman, Y. (1926). The impulses produced by sensory nerve endings: part 3. Impulses set up by touch and pressure. *J. Physiol.* 61, 465–483. doi: 10.1113/jphysiol.1926.sp002308
- Aitchison, L., and Latham, P. E. (2015). Synaptic sampling: a connection between psp variability and uncertainty explains neurophysiological observations. *arXiv [Preprint]*. arXiv: 1505.04544.
- Amari, S. (1977). Dynamics of pattern formation in lateral-inhibition type neural fields. *Biol. Cybern.* 27, 77–87. doi: 10.1007/bf00337259
- Bahdanau, D., Cho, K., and Bengio, Y. (2014). Neural machine translation by jointly learning to align and translate. *arXiv [Preprint]*. arXiv: 1409.0473.
- Baraduc, P., and Guigon, E. (2002). Population computation of vectorial transformations. *Neural Comput.* 14, 845–871. doi: 10.1162/089976602317318983
- Baraduc, P., Guigon, E., and Burnod, Y. (2001). Recoding arm position to learn visuomotor transformations. *Cereb. Cortex* 11, 906–917. doi: 10.1093/cercor/11.10.906
- Baxter, J. (1998). “Theoretical models of learning to learn,” in *Learning to Learn*, eds S. Thrun and L. Pratt (Boston, MA: Springer), 71–94.
- Bengio, Y. (2012). “Practical recommendations for gradient-based training of deep architectures,” in *Neural Networks: Tricks of the Trade*, eds G. Montavon, G. B. Orr, and K.-B. Müller (Berlin; Heidelberg: Springer), 437–478.
- Ben-Yishai, R., Bar-Or, R. L., and Sompolinsky, H. (1995). Theory of orientation tuning in visual cortex. *Proc. Natl. Acad. Sci. U.S.A.* 92, 3844–3848. doi: 10.1073/pnas.92.9.3844
- Bethge, M., Rotermund, D., and Pawelzik, K. (2002). Optimal short-term population coding: when fisher information fails. *Neural Comput.* 14, 2317–2351. doi: 10.1162/08997660260293247
- Bethge, M., Rotermund, D., and Pawelzik, K. (2003). Optimal neural rate coding leads to bimodal firing rate distributions. *Network* 14, 303–319. doi: 10.1088/0954-898X_14_2_307

- Bishop, C. (2006). *Pattern Recognition and Machine Learning*. New York, NY: Springer.
- Blundell, C., Cornebise, J., Kavukcuoglu, K., and Wierstra, D. (2015). Weight uncertainty in neural networks. *arXiv [Preprint]*. arXiv: 1505.05424.
- Blundell, C., Uria, B., Pritzel, A., Li, Y., Ruderman, A., Leibo, J. Z., et al. (2016). Model-free episodic control. *arXiv [Preprint]*. arXiv: 1606.04460.
- Booth, M. C., and Rolls, E. T. (1998). View-invariant representations of familiar objects by neurons in the inferior temporal visual cortex. *Cereb. Cortex (New York, NY: 1991)* 8, 510–523. doi: 10.1093/cercor/8.6.510
- Botvinick, M., Ritter, S., Wang, J. X., Kurth-Nelson, Z., Blundell, C., and Hassabis, D. (2019). Reinforcement learning, fast and slow. *Trends Cogn. Sci.* 23, 408–422. doi: 10.1016/j.tics.2019.02.006
- Branco, T., and Staras, K. (2009). The probability of neurotransmitter release: variability and feedback control at single synapses. *Nat. Rev. Neurosci.* 10:373. doi: 10.1038/nrn2634
- Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., et al. (2016). OpenAI gym. arXiv:1606.01540.
- Broomhead, D. S., and Lowe, D. (1988). *Radial Basis Functions, Multi-variable Functional Interpolation and Adaptive Networks*. Technical report, Royal Signals and Radar Establishment Malvern.
- Brown, W. M., and Bäcker, A. (2006). Optimal neuronal tuning for finite stimulus spaces. *Neural Comput.* 18, 1511–1526. doi: 10.1162/neco.2006.18.7.1511
- Brugge, J. F., Kitzes, L. M., and Javel, E. (1981). Postnatal development of frequency and intensity sensitivity of neurons in the anteroventral cochlear nucleus of kittens. *Hear. Res.* 5, 217–229. doi: 10.1016/0378-5955(81)90047-2
- Brunel, N., and Nadal, J.-P. (1998). Mutual information, fisher information, and population coding. *Neural Comput.* 10, 1731–1757. doi: 10.1162/089976698300017115
- Bullock, D., Grossberg, S., and Guenther, F. H. (1993). A self-organizing neural model of motor equivalent reaching and tool use by a multijoint arm. *J. Cogn. Neurosci.* 5, 408–435. doi: 10.1162/jocn.1993.5.4.408
- Butts, D. A., and Goldman, M. S. (2006). Tuning curves, neuronal variability, and sensory coding. *PLoS Biol.* 4:e92. doi: 10.1371/journal.pbio.0040092
- de Broissia, A. D. F., and Sigaud, O. (2016). Actor-critic versus direct policy search: a comparison based on sample complexity. *arXiv [Preprint]*. arXiv: 1606.09152.
- Dhariwal, P., Hesse, C., Klimov, O., Nichol, A., Plappert, M., Radford, A., et al. (2017). *Openai Baselines*. Available online at: <https://github.com/openai/baselines>
- Eurich, C. W., Wilke, S. D., and Schwegler, H. (2000). “Neural representation of multi-dimensional stimuli,” in *Advances in Neural Information Processing Systems* (Cambridge, MA; London: MIT Press), 115–121.
- Freedman, D. J., Riesenhuber, M., Poggio, T., and Miller, E. K. (2005). Experience-dependent sharpening of visual shape selectivity in inferior temporal cortex. *Cereb. Cortex* 16, 1631–1644. doi: 10.1093/cercor/bhj100
- Gauy, M. M., Meier, F., and Steger, A. (2017). Multiassociative memory: recurrent synapses increase storage capacity. *Neural Comput.* 29, 1375–1405. doi: 10.1162/NECO_a_00954
- Geisler, W. S., Najemnik, J., and Ing, A. D. (2009). Optimal stimulus encoders for natural tasks. *J. Vis.* 9:17. doi: 10.1167/9.13.17
- Georgopoulos, A. P., Kettner, R. E., and Schwartz, A. B. (1988). Primate motor cortex and free arm movements to visual targets in three-dimensional space. II. Coding of the direction of movement by a neuronal population. *J. Neurosci.* 8, 2928–2937. doi: 10.1523/JNEUROSCI.08-08-02928.1988
- Gershman, S. J., and Daw, N. D. (2017). Reinforcement learning and episodic memory in humans and animals: an integrative framework. *Annu. Rev. Psychol.* 68, 101–128. doi: 10.1146/annurev-psych-122414-033625
- Graves, A., Wayne, G., and Danihelka, I. (2014). Neural Turing machines. *arXiv [Preprint]*. arXiv: 1410.5401.
- Harper, N. S., and McAlpine, D. (2004). Optimal neural population coding of an auditory spatial cue. *Nature* 430:682. doi: 10.1038/nature02768
- Harris, J. J., Jolivet, R., and Attwell, D. (2012). Synaptic energy use and supply. *Neuron* 75, 762–777. doi: 10.1016/j.neuron.2012.08.019
- Hinton, G. E. (1986). “Learning distributed representations of concepts,” in *Proceedings of the Eighth Annual Conference of the Cognitive Science Society* (Amherst, MA: Lawrence Erlbaum Associates), 1–12.
- Hochreiter, S., Younger, A. S., and Conwell, P. R. (2001). “Learning to learn using gradient descent,” in *International Conference on Artificial Neural Networks* (Berlin; Heidelberg: Springer), 87–94.
- Hubel, D. H., and Wiesel, T. N. (1959). Receptive fields of single neurones in the cat's striate cortex. *J. Physiol.* 148, 574–591. doi: 10.1113/jphysiol.1959.sp006308
- Hubel, D. H., and Wiesel, T. N. (1962). Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. *J. Physiol.* 160, 106–154. doi: 10.1113/jphysiol.1962.sp006837
- Jaderberg, M., Czarnecki, W. M., Dunning, I., Marris, L., Lever, G., Castaneda, A. G., et al. (2018). Human-level performance in first-person multiplayer games with population-based deep reinforcement learning. *arXiv [Preprint]*. arXiv: 1807.01281.
- Johnson, D. H., and Ray, W. (2004). Optimal stimulus coding by neural populations using rate codes. *J. Comput. Neurosci.* 16, 129–138. doi: 10.1023/B:JCNS.0000014106.09948.83
- Kappel, D., Habenschuss, S., Legenstein, R., and Maass, W. (2015). Network plasticity as bayesian inference. *PLoS Comput. Biol.* 11:e1004485. doi: 10.1371/journal.pcbi.1004485
- Kappel, D., Legenstein, R., Habenschuss, S., Hsieh, M., and Maass, W. (2018). A dynamic connectome supports the emergence of stable computational function of neural circuits through reward-based learning. *eNeuro* 5:ENEURO.0301-17.2018. doi: 10.1523/ENEURO.0301-17.2018
- Kim, S. S., Rouault, H., Druckmann, S., and Jayaraman, V. (2017). Ring attractor dynamics in the drosophila central brain. *Science* 356, 849–853. doi: 10.1126/science.aal4835
- Klopfenstein, R. W., and Sverdlöve, R. (1983). “Approximation by uniformly spaced gaussian functions,” in *Approximation Theory IV*, eds C.K. Chui, L. L. Schumaker, and J. D. Ward (New York, NY: Academic Press), 575–580.
- Knierim, J. J., and Zhang, K. (2012). Attractor dynamics of spatially correlated neural activity in the limbic system. *Annu. Rev. Neurosci.* 35, 267–285. doi: 10.1146/annurev-neuro-062111-150351
- Knudsen, E. I., and Konishi, M. (1978). Center-surround organization of auditory receptive fields in the owl. *Science* 202, 778–780. doi: 10.1126/science.715444
- Ko, H., Hofer, S. B., Pichler, B., Buchanan, K. A., Sjöström, P. J., and Mrsic-Flogel, T. D. (2011). Functional specificity of local synaptic connections in neocortical networks. *Nature* 473:87. doi: 10.1038/nature09880
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). “Imagenet classification with deep convolutional neural networks,” in *Advances in Neural Information Processing Systems* 25, eds F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger (Red Hook, NY: Curran Associates, Inc.), 1097–1105.
- Lánský, P., and Greenwood, P. E. (2005). Optimal signal estimation in neuronal models. *Neural Comput.* 17, 2240–2257. doi: 10.1162/0899766054615653
- Lecun, Y., and Bengio, Y. (1995). “Convolutional networks for images, speech, and time-series” in *The handbook of BRAIN theory and Neural Networks* ed M. A. Arbib (Cambridge, MA: MIT Press), 10.
- Lee, S.-H., Kwan, A. C., Zhang, S., Phoumthipphavong, V., Flannery, J. G., Masmanidis, S. C., et al. (2012). Activation of specific interneurons improves v1 feature selectivity and visual perception. *Nature* 488:379. doi: 10.1038/nature11312
- Lengyel, M., and Dayan, P. (2008). “Hippocampal contributions to control: the third way,” in *Advances in Neural Information Processing Systems*, 889–896.
- Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., et al. (2015). Continuous control with deep reinforcement learning. *arXiv [Preprint]*. arXiv: 1509.02971.
- Llera-Montero, M., Sacramento, J., and Costa, R. P. (2019). Computational roles of plastic probabilistic synapses. *Curr. Opin. Neurobiol.* 54, 90–97. doi: 10.1016/j.conb.2018.09.002
- McDonnell, M. D., and Stocks, N. G. (2008). Maximally informative stimuli and tuning curves for sigmoidal rate-coding neurons and populations. *Phys. Rev. Lett.* 101:058103. doi: 10.1103/PhysRevLett.101.058103
- Menz, M. D., and Freeman, R. D. (2003). Stereoscopic depth processing in the visual cortex: a coarse-to-fine mechanism. *Nat. Neurosci.* 6:59. doi: 10.1038/nn986
- Montemurro, M. A., and Panzeri, S. (2006). Optimal tuning widths in population coding of periodic variables. *Neural Comput.* 18, 1555–1576. doi: 10.1162/neco.2006.18.7.1555
- Mrsic-Flogel, T. D., Schnupp, J. W., and King, A. J. (2003). Acoustic factors govern developmental sharpening of spatial tuning in the auditory cortex. *Nat. Neurosci.* 6:981. doi: 10.1038/nn1108

- Mysore, S. P., Asadollahi, A., and Knudsen, E. I. (2010). Global inhibition and stimulus competition in the owl optic tectum. *J. Neurosci.* 30, 1727–1738. doi: 10.1523/JNEUROSCI.3740-09.2010
- Nikitin, A. P., Stocks, N. G., Morse, R. P., and McDonnell, M. D. (2009). Neural population coding is optimized by discrete tuning curves. *Phys. Rev. Lett.* 103:138101. doi: 10.1103/PhysRevLett.103.138101
- O'Keefe, J. (1976). Place units in the hippocampus of the freely moving rat. *Exp. Neurol.* 51, 78–109. doi: 10.1016/0014-4886(76)90055-8
- Panzeri, S., Treves, A., Schultz, S., and Rolls, E. T. (1999). On decoding the responses of a population of neurons from short time windows. *Neural Comput.* 11, 1553–1577. doi: 10.1162/089976699300016142
- Pouget, A., Deneve, S., Ducom, J.-C., and Latham, P. E. (1999). Narrow versus wide tuning curves: what's best for a population code? *Neural Comput.* 11, 85–90. doi: 10.1162/089976699300016818
- Pouget, A., Zhang, K., Deneve, S., and Latham, P. E. (1998). Statistically efficient estimation using population coding. *Neural Comput.* 10, 373–401. doi: 10.1162/089976698300017809
- Pritzel, A., Uria, B., Srinivasan, S., Puigdomènech, A., Vinyals, O., Hassabis, D., et al. (2017). "Neural episodic control," in *Proceedings of the 34th International Conference on Machine Learning, Vol. 70* (Sydney: JMLR. Org), 2827–2836.
- Ranck, J. B. J. (1985). "Head direction cells in the deep cell layer of dorsolateral pre-subiculum in freely moving rats," in *Electrical Activity of the Archicortex*, eds G. Buzsáki and C. H. Vanderwolf (Budapest, Akademiai Kiado).
- Ringach, D. L., Hawken, M. J., and Shapley, R. (1997). Dynamics of orientation tuning in macaque primary visual cortex. *Nature* 387:281. doi: 10.1038/387281a0
- Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1985). *Learning Internal Representations by Error Propagation*. Technical report, California Univ San Diego La Jolla Inst for Cognitive Science.
- Salinas, E., and Abbott, L. F. (1995). Transfer of coded information from sensory to motor networks. *J. Neurosci.* 15, 6461–6474. doi: 10.1523/JNEUROSCI.15-10-06461.1995
- Samonds, J. M., Potetz, B. R., and Lee, T. S. (2009). Cooperative and competitive interactions facilitate stereo computations in macaque primary visual cortex. *J. Neurosci.* 29, 15780–15795. doi: 10.1523/JNEUROSCI.2305-09.2009
- Sanger, T. D. (1991). Optimal hidden units for two-layer nonlinear feedforward neural networks. *Int. J. Patt. Recogn. Artif. Intell.* 5, 545–561. doi: 10.1142/S0218001491000314
- Sanger, T. D. (1997). A probability interpretation of neural population coding for movement, *Adv. Psychol.* 119, 75–116. doi: 10.1016/S0166-4115(97)80005-2
- Sanger, T. D. (1998). Probability density methods for smooth function approximation and learning in populations of tuned spiking neurons. *Neural Comput.* 10, 1567–1586. doi: 10.1162/089976698300017313
- Sanger, T. D. (2003). Neural population codes. *Curr. Opin. Neurobiol.* 13, 238–249. doi: 10.1016/S0959-4388(03)00034-5
- Schaul, T., and Schmidhuber, J. (2010). Metalearning. *Scholarpedia* 5:4650. doi: 10.4249/scholarpedia.4650
- Schmidhuber, J., Zhao, J., and Wiering, M. (1996). Simple principles of metalearning. *Technical Report IDSIA* 69, 1–23.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. (2017). Proximal policy optimization algorithms. *arXiv [Preprint]*. arXiv: 1707.06347.
- Schweighofer, N., and Doya, K. (2003). Meta-learning in reinforcement learning. *Neural Netw.* 16, 5–9. doi: 10.1016/S0893-6080(02)00228-9
- Seeholzer, A., Deger, M., and Gerstner, W. (2017). Efficient low-dimensional approximation of continuous attractor networks. *arXiv [Preprint]*. arXiv: 1711.08032.
- Seelig, J. D., and Jayaraman, V. (2015). Neural dynamics for landmark orientation and angular path integration. *Nature* 521:186. doi: 10.1038/nature14446
- Seriès, P., Latham, P. E., and Pouget, A. (2004). Tuning curve sharpening for orientation selectivity: coding efficiency and the impact of correlations. *Nat. Neurosci.* 7:1129. doi: 10.1038/nn1321
- Seung, H. S. (2003). Learning in spiking neural networks by reinforcement of stochastic synaptic transmission. *Neuron* 40, 1063–1073. doi: 10.1016/S0896-6273(03)00761-x
- Seung, H. S., and Sompolinsky, H. (1993). Simple models for reading neuronal population codes. *Proc. Natl. Acad. Sci. U.S.A.* 90, 10749–10753. doi: 10.1073/pnas.90.22.10749
- Skaggs, W. E., Knierim, J. J., Kudrimoti, H. S., and McNaughton, B. L. (1995). "A model of the neural basis of the rat's sense of direction," in *Advances in Neural Information Processing Systems* (Red Hook, NY: Curran Associates, Inc.), 173–180.
- Spitzer, H., Desimone, R., and Moran, J. (1988). Increased attention enhances both behavioral and neuronal performance. *Science* 240, 338–340. doi: 10.1126/science.3353728
- Stanley, K. O., and Miikkulainen, R. (2002). Evolving neural networks through augmenting topologies. *Evol. Comput.* 10, 99–127. doi: 10.1162/106365602320169811
- Suga, N., Zhang, Y., and Yan, J. (1997). Sharpening of frequency tuning by inhibition in the thalamic auditory nucleus of the mustached bat. *J. Neurophysiol.* 77, 2098–2114. doi: 10.1152/jn.1997.77.4.2098
- Sutton, R. S., and Barto, A. G. (1998). *Introduction to Reinforcement Learning*. Cambridge: MIT Press.
- Tavazoie, S. F., and Reid, R. C. (2000). Diverse receptive fields in the lateral geniculate nucleus during thalamocortical development. *Nat. Neurosci.* 3:608. doi: 10.1038/75786
- Thrun, S., and Pratt, L. (eds.). (1998). "Learning to learn: introduction and overview," in *Learning to Learn* (Boston, MA: Springer).
- Todorov, E. (2002). Cosine tuning minimizes motor errors. *Neural Comput.* 14, 1233–1260. doi: 10.1162/089976602753712918
- Toyoizumi, T., Aihara, K., and Amari, S. (2006). Fisher information for spike-based population decoding. *Phys. Rev. Lett.* 97:098102. doi: 10.1103/PhysRevLett.97.098102
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., et al. (2017). "Attention is all you need," in *Advances in Neural Information Processing Systems 30*, eds I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Red Hook, NY: Curran Associates, Inc.), 5998–6008.
- Wagner, H. (1990). Receptive fields of neurons in the owl's auditory brainstem change dynamically. *Eur. J. Neurosci.* 2, 949–959. doi: 10.1111/j.1460-9568.1990.tb00007.x
- Wan, L., Zeiler, M., Zhang, S., Le Cun, Y., and Fergus, R. (2013). "Regularization of neural networks using dropconnect," in *International Conference on Machine Learning (JMLR.org)*, 1058–1066.
- Wang, X., Lu, T., Snider, R. K., and Liang, L. (2005). Sustained firing in auditory cortex evoked by preferred stimuli. *Nature* 435:341. doi: 10.1038/nature03565
- Weliky, M., Kandler, K., Fitzpatrick, D., and Katz, L. C. (1995). Patterns of excitation and inhibition evoked by horizontal connections in visual cortex share a common relationship to orientation columns. *Neuron* 15, 541–552. doi: 10.1016/0896-6273(95)90143-4
- Williams, R. J. (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Mach. Learn.* 8, 229–256. doi: 10.1007/BF00992696
- Wilson, H. R., and Cowan, J. D. (1973). A mathematical theory of the functional dynamics of cortical and thalamic nervous tissue. *Kybernetik* 13, 55–80. doi: 10.1007/bf00288786
- Yaeli, S., and Meir, R. (2010). Error-based analysis of optimal tuning functions explains phenomena observed in sensory neurons. *Front. Comput. Neurosci.* 4:130. doi: 10.3389/fncom.2010.00130
- Yarrow, S., and Seriès, P. (2015). The influence of population size, noise strength and behavioral task on best-encoded stimulus for neurons with unimodal or monotonic tuning curves. *Front. Comput. Neurosci.* 9:18. doi: 10.3389/fncom.2015.00018
- Zhang, K., and Sejnowski, T. J. (1999). Neuronal tuning: to sharpen or broaden? *Neural Comput.* 11, 75–84. doi: 10.1162/089976699300016809

Conflict of Interest: The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Copyright © 2020 Meier, Dang-Nhu and Steger. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.