

Network of Evolvable Neural Units: Evolving to Learn at a Synaptic Level

Paul Bertens^{*1}, Seong-Whan Lee^{1,2}

¹Department of Brain and Cognitive Engineering,

²Department of Artificial Intelligence,
Korea University, Seoul, South Korea

December 17, 2019

Abstract

Although Deep Neural Networks have seen great success in recent years through various changes in overall architectures and optimization strategies, their fundamental underlying design remains largely unchanged. Computational neuroscience on the other hand provides more biologically realistic models of neural processing mechanisms, but they are still high level abstractions of the actual experimentally observed behaviour. Here a model is proposed that bridges Neuroscience, Machine Learning and Evolutionary Algorithms to evolve individual soma and synaptic compartment models of neurons in a scalable manner. Instead of attempting to manually derive models for all the observed complexity and diversity in neural processing, we propose an Evolvable Neural Unit (ENU) that can approximate the function of each individual neuron and synapse. We demonstrate that this type of unit can be evolved to mimic Integrate-And-Fire neurons and synaptic Spike-Timing-Dependent Plasticity. Additionally, by constructing a new type of neural network where each synapse and neuron is such an evolvable neural unit, we show it is possible to evolve an agent capable of learning to solve a T-maze environment task. This network independently discovers spiking dynamics and reinforcement type learning rules, opening up a new path towards biologically inspired artificial intelligence.

1 Introduction

Much research has been done to understand how neural processing works [1, 2, 3], including synaptic learning [4] and overall network dynamics [5, 6]. However these processes are extremely complex and interact in a way that is currently still not well understood. Mathematical models are commonly used to approximate neurons, synapses and learning mechanisms [7, 4]. However, neurons can behave in many different ways, and no single model can accurately capture all experimentally observed behaviour. The neural complexity is mostly at the small scale, i.e. the cellular local interactions that can lead to an intelligent overall system. While information processing at the network level is generally well understood (and the basis for artificial neural networks [8, 9]), its how individual neurons are actually capable of giving rise to the overall networks behaviour and learning capability that is largely unknown.

^{*}Github: <https://github.com/paulbertens>, e-mail: p.m.w.a.bertens@gmail.com

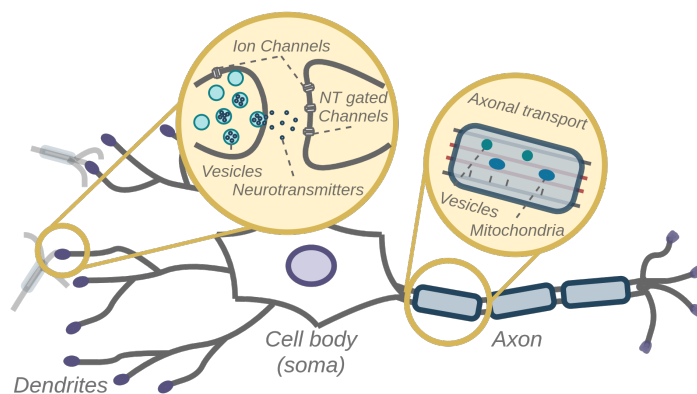


Figure 1: Abstract overview of a single neuron. Functionally the dendrites receive information, the cell body integrates and transforms that information and the axon transmits it to other neurons. Neurotransmitters allow for communication between neurons, and the number of neurotransmitters released affect the strength of the synaptic connection. Learning then occurs by updating this synaptic strength depending on the information received. Besides graded and action potentials (spikes), complex transport of mitochondria, vesicles used to contain neurotransmitters and other information that can change the behaviour of the neuron also occur along the axon and dendrite.

Biological Neural Networks When looking in more detail at the smaller scale of neurons and synapses, complexity quickly arises. Different ion-channels in the neuron are responsible for generating action potentials (spikes) or graded potentials (real valued), mainly sodium (Na^+), potassium (K^+) and calcium (Ca^{2+}) channels [10]. These ion channels are also responsible for triggering the release of chemical neurotransmitters, the main method of communication between neurons in biological neural networks. Many type of neurotransmitters exist that each have distinct roles, e.g. Dopamine [11] (related to learning), GABA_B [12] (inhibitory), Acetylcholine [13] (motor neurons) and Glutamate [14] (excitatory). Morphological differences across the brain are also common, where different type of neurons exist in different regions of the brain and across different cortical layers [15].

Axons and dendrites actively transport other information [16], for example vesicles [17] (which hold neurotransmitters) and mitochondria [18] (required for membrane excitability and neuroplasticity). Many protein types are also common [16], which are responsible for a diver set of functions, including increasing or decreasing the number of neuroreceptors in the synapses [19], and aiding in the release of neurotransmitters [17]. Furthermore, neurons can exhibit both regular or bursting spiking behaviour depending on their internal state and neuronal type [20, 21]. Deriving models capable of capturing all this behaviour is thus an active research area. These mechanisms might however simply be the result of biological constraints and evolutionary processes, and it is still uncertain which structures exactly are required to exist for intelligent behaviour to emerge.

Generally models in the computational neuroscience field focus on modeling individual neurons and synaptic learning behaviour. Simple Integrate and fire neurons [7] (IAF) assume input potential is summated over time, and once this reaches some threshold a spike occurs. More advanced models also try to model the sodium and potassium ion channels as done in the HodgkinHuxley (HH) model [22], which gives more realistic action potentials at the output.

Hebbian plasticity was one of the first proposed models describing synaptic learning [23], briefly it states that the synaptic weight increases if the pre-synaptic neuron fires before the post-synaptic neuron. This change in synaptic weight can then increase or decrease the firing rate of the post-synaptic neuron in the future when receiving a similar stimuli. A more extensive model that also takes the spike timing into account is Spike Timing Dependent Plasticity [4] (STDP), which updates the weights depending on the time between the pre and post synaptic neuron.

Artificial and Recurrent Neural Networks Deep Learning models have been very successful in many practical applications [9], and although spiking neural networks (SNNs) have been developed [24, 25], so far they have seen limited success due to their high computational demand and difficulty in training. Artificial Neural Networks, the foundation of deep learning models, are especially effective in performing function approximation, and are known to be universal function approximators [26]. Recurrent Neural Network (RNNs) architectures have also been investigated in order to process sequential data and allow memory to be stored between successive time steps [9], most commonly Long-Short Term Memory [27] (LSTMs) and Gated Recurrent Units [28] (GRUs). These add gating mechanisms to standard RNNs to allow for easier learning of long-range dependencies and to avoid vanishing gradients in back-propagation.

Deep learning based models however suffer from several limitations [1]. They require backpropagation of errors through the whole network, have difficulty performing one-shot learning (learning from a single example) and suffer from catastrophic forgetting of a previous task once trained on a new task. They are also extreme abstractions of biological neural networks, only considering a single weight value on the connections and having a single real valued number as their output.

Evolutionary Algorithms Evolutionary algorithms have been widely applied to a variety of domains [29, 30], and many related optimization algorithms exist that are based on the evolutionary principles of mutation, reproduction and survival of the fittest. Recently Evolution Strategies (ES) have been successfully applied to train deep neural networks, despite the large amount of parameters of most neural network architectures [31]. The advantage of ES is that it allows for non-differentiable objective functions, and when training RNNs does not require backpropagation through time. This means we can optimize and evolve the parameters of a model to solve any desired objective we want, and potentially learn over arbitrarily long sequences.

2 Objective

Past attempts on modeling biological neural networks have mostly been focused on manually deriving mathematical rules and abstractions based on experimental data, however in this paper a different approach is taken. The functionality of different components in neural networks are approximated using artificial Recurrent Neural Networks (RNNs), such that each synapse and neuron in the network is a recurrent neural network. To make this computationally feasible the RNN weight parameters across the global network are shared, while each RNN keeps their own internal state. This way the number of parameters to be learned are significantly reduced, and we can use evolutionary strategies to train such a network. Due to the universal

function approximation properties of RNNs, there is essentially no limitations on the neuronal behaviour we can obtain.

The RNN memory can store dynamic parameters that determine the behaviour of our compartments (e.g. synapse or neural cell body), while the shared weights are evolved. This is related to meta-learning approaches [32, 33, 34, 35, 36], where we are learning to learn. Since all RNN weights are shared across all synaptic or neural compartments, learning different behaviour can only occur by learning (evolving) to store and update dynamic parameters in the RNN memory state.

Such an RNN can thus be seen as modeling a function that expresses the neuronal behaviour. However, instead of having a fixed mathematical function derived from experimental observation (like e.g. hebbian learning, integrate-and-fire neurons or STDP rules), we evolve a function that could exhibit much more diverse behaviour. In our case it is not required that the evolved behaviour has to match experimental data or preexisting models, as cells in the real world are bound by more physical constraints than in our simulated setting. Ultimately we would like to evolve a type of mini-agent that when duplicated and connected together in an overall neural network exhibit intelligent learning behaviour. To this end, network agents containing these mini-agents are evolved using evolutionary algorithms, where their fitness is their ability to learn new tasks.

2.1 Contributions

To summarize, we are thus attempting to evolve a new type of Neural Network where each individual neuron and synapse in the network is by itself modeled by an evolvable Recurrent Neural Network. Our main contributions are as follows:

- Propose an Evolvable Neural Unit (ENU) that can be evolved through Evolution Strategies (ES) which is able to learn to store relevant information in its internal state memory and perform complex processing on the received input using that memory.
- Demonstrate such an ENU can be evolved to approximate the neural dynamics of simple Integrate and Fire neurons (IAF) and synaptic Spike-Timing-Dependent Plasticity (STDP).
- Combine multiple ENUs in a larger network, where the weight parameters of the ENUs to be evolved are shared across the neurons and synaptic compartments, but each internal state is unique and updated based on the local information they receive.
- Efficiently evolve and compute such a network of ENUs through large scale matrix multiplications on a single Graphical Processing Unit (GPU).
- Show it is possible to evolve a network of ENUs capable of reward based reinforcement learning purely through local dynamics, i.e. evolving to learn.

3 Proposed Model

3.1 Evolvable Neural Units (ENUs)

As the basis for an evolvable Recurrent Neural Network (RNN) to approximate neural and synaptic behaviour we build upon previous work on modeling long-term dependencies through Gated Recurrent Units [28] (GRUs). We extend upon this model and add an additional output

gate that applies a non-linear activation function and feeds this back to the input, which simultaneously also reduces the number of possible output channels (improving computational complexity). Additionally they are implemented in such a way that allows them to be combined and evolved efficiently in a larger overall network. We term these units Evolvable Neural Units (ENUs).

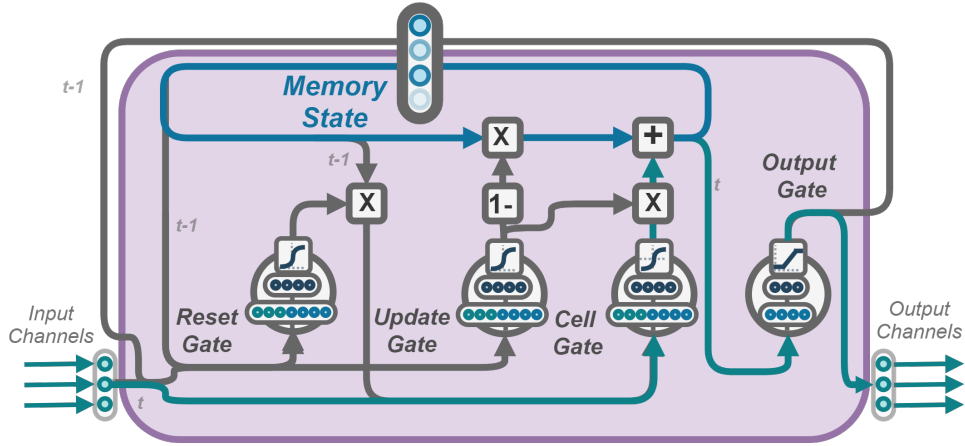


Figure 2: **An Evolvable Neural Unit.** Different evolvable gating mechanisms protect the internal memory of the neural unit (storing "dynamic parameters"). This memory can potentially encode and store the neuronal behaviour and determines how information is processed (analogous to e.g. the membrane potential, synaptic weights, protein states etc.). The reset gate can forget past information, the update gate determines how much we change the memory state (our dynamic parameters) and the cell gate determines the new value of the dynamic parameters. An additional output gate is also used to reduce the number of output channels and to allow for spike generation to potentially evolve. In this diagram example we have 4 dynamic parameters, 3 input channels and 3 output channels.

Using a gating structure allows us to have fine grained control over how different input influences the internal memory state (storing "dynamic parameters"), which in turn controls how that input is processed and how the dynamic parameters are updated. The input and output of the ENU is a vector, where each value in the vector can be considered a type of "channel" used to transmit information between different ENUs. The internal memory state is also a vector, that can store multiple values. It enables us to condition the received input on the dynamic parameters stored in the memory state, which control both the output and gating behaviour of the ENU (see figure 2). Each gate in the model is a standard single layer artificial neural network with k output units and an evolvable weight matrix w , where k is equal to the internal memory state size. These gates can process information from the current input channels and previous memory state. This allows us for example to evolve a function that adds some value to a dynamic parameter in the memory state, but only if there is a spike at a certain input channel. Evolving such units to perform complex functions then becomes significantly easier than in standard RNNs, which can suffer from vanishing values and undesired updates to their internal memory state [27]. Protecting updates to the memory state is especially important in our case, since they define the dynamic parameters that determine the behaviour of our ENU, and should be able to persist over their entire lifetime.

Intuitively this ENU type architecture allows us for example to evolve spike based behaviour through storing and summing received input into the memory state. Then once some threshold is reached the output gate can evolve to activate and output a value, a "spike". This spike is then fed back into the input allowing the reset and update gate to evolve to reset the internal memory state.

For synaptic compartments the internal state could evolve to memorize when a pre and post synaptic spike occurs , and also store and update some dynamic "weight" parameter that can change how the input is processed and passed to the post-synaptic neuron. Additionally, multiple input and output channels allow flexibility in evolving different type of information processing mechanisms, e.g. using a type of neurotransmitter, or even functions analogous to dendritic and axonal transport.

3.2 Network of ENUs

By combining multiple ENUs in a single network we can construct a ENU based Neural Network (ENU-NN), this network connects multiple neural and synaptic compartment ENU models together as seen in figure 3. All soma and synaptic neuronal compartments share the same ENU gate parameters, and we only evolve the weights of these shared gates. This means we have two "chromosomes" to evolve, one for the synapses and one for the neurons, shared across all synapses and neurons. However, they each have unique internal memory state variables which allows for complex signal processing and learning behaviour to occur, as the compartments can evolve to update their internal states depending on the local input they receive. Synaptic plasticity (the synaptic weights) in such a model would thus no longer be encoded in the weights of the network directly, but could instead evolve to be stored and dynamically updated in the internal states of the synapse ENUs as a function of the current and past input.

To evolve reinforcement type learning behaviour in such a network we can construct a sparse recurrent network with several input sensory neurons that detect e.g. different colors in front of the agent (which has an ENU network) in a given environment. We can also designate some ENU neurons as output motor neurons that determine the action the agent should take. Additionally, to allow reward feedback we can have a reward neuron that uses different ENU input channels to indicate environmental rewards obtained (see Figure 3).

The main components of the ENU network are as follows:

- **Neuron ENUs:** Analogous to the biological cell body and axon. Responsible for transforming the summated input from its connected synapses. Can evolve to use multiple channels to transmit information analogous to e.g. spikes or neurotransmitters. Also propogates it's output backwards to all connected input synapses, this way the synapse ENU compartment can potentially learn STDP type learning rules based on post-synaptic neural behaviour.
- **Synapse ENUs:** Analogous to the biological dendrite and synapse. Transforms information from the pre-synaptic to post-synaptic ENU neuron, similar to weights in artificial neural networks. However, it is allowed to use multiple channels to learn to transmit information analogous to e.g. spikes, graded potential or other types of dendritic transport.
- **Integration step:** Sums up the output per channel of each incoming synapse connection, which is then processed by the post-synaptic Neuron ENU.

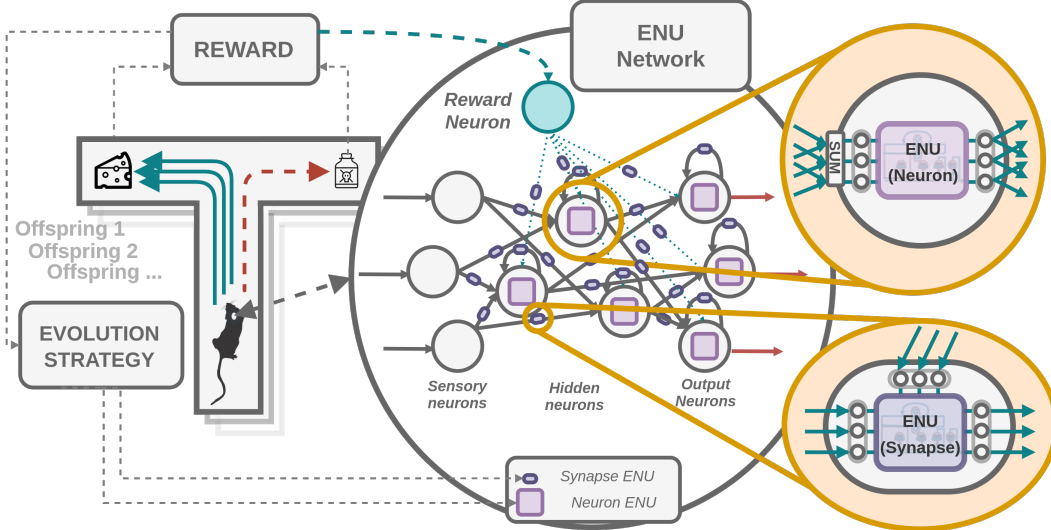


Figure 3: **Network of ENUs.** Example of evolving reinforcement based learning to solve a T-maze environment task. Reward from the environment determines the fitness of the agent. This reward is also passed to the reward neuron in the ENU Network, which connects to the other ENU neurons, allowing them to evolve to update their incoming synapse ENUs according to the reward obtained. Several sensory input neurons are used to detect colors in front of the agent, and the output neurons determine the action the agent has to take (forward, left or right). Each ENU is capable of updating their unique internal dynamic parameters depending on the current and past input, which changes how information is processed over multiple channels and compartments within the overall ENU network.

4 Method

4.1 Evolution Strategies

Since our desired goal is to solve Reinforcement Learning (RL) based environments and the task is not directly differentiable, we cannot use standard Back-propagation through time (BPTT) to optimize our network. We would also like to be able to learn over long time-spans, which makes BPTT impractical due to vanishing gradient issues. This makes Evolution Strategies (ES) ideally suited [37, 38, 39, 40], and we use a similar approach taken in previous work [31]. The gradient is approximated through random Gaussian sampling across the parameter space, and the weights are updated following this approximate gradient direction (see Algorithm 1). We can also use standard Stochastic Gradient Descent (SGD) methods like momentum [41] on the resulting approximate gradients obtained. Due to weight sharing we only have to evolve the gate parameters of two ENUs, one for the synapse, and one for the neuron, giving us a relatively small parameter space.

Fitness Ranking We use fitness ranking in order to reduce the effect of outliers [31]. We can sort and assign rank values to each offspring according to their fitness, which determines their relative weight in calculating our approximate gradient. We then get a transformed fitness function $F_r(\theta_i) = \frac{\text{rank}(F(\theta_i))^5}{\sum_{i=0}^n \text{rank}(F(\theta_i))^5}$, where $\text{rank}()$ assigns a linear ranking from 1.0 to 0.0. This results in around the top 20% best performing agents to account for 80% of the gradient estimate.

Algorithm 1 Evolution Strategies

Input: base parameters θ_k , learning rate α , standard deviation σ
for $k = 1$ **to** N **do**
 Sample offspring mutations $\epsilon_i, \dots, \epsilon_n$ from $N(0, I)$
 Compute fitness $F_i = F(\theta_k + \sigma\epsilon_i)$ for $i = 1, \dots, n$
 Calculate approximate gradient $G_k = \sum_{i=1}^n F_i\epsilon_i$
 Update base parameters $\theta_{k+1} = \theta_k + \alpha G_k$
end for

Batching Mini-batches were also used to better estimate the fitness of the same offspring across multiple environments, so the fitness of each offspring was determined by the mean fitness across m mini-batch environments.

4.2 Gating in an Evolvable Neural Unit

The exact equations for the ENU gates and updating of the memory state are as follows:

$$\begin{aligned} z_t &= \sigma(W_z \cdot [h_{t-1}, o_{t-1}, x_t]) \\ r_t &= \sigma(W_r \cdot [h_{t-1}, o_{t-1}, x_t]) \\ \tilde{h}_t &= \tanh(W_c \cdot [r_t \odot h_{t-1}, o_{t-1}, x_t]) \\ h_t &= (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t \\ o_t &= \text{clip}(W_o \cdot h_t, 0, 1) \end{aligned} \tag{1}$$

Where x_t is the input, z_t the update gate, r_t the reset gate, \tilde{h}_t the cell gate, h_t the new memory state and o_t the output gate. σ is the sigmoid function, \cdot is matrix multiplication and \odot is element-wise multiplication. W_* are the weight matrices of each gate, and the parameters to be evolved. We can also apply clipping (restricting output values to be between 0 and 1), since we use evolution strategies to optimize our parameters and do not require a strict differentiable activation function. This clipping results in a thresholded output and prevents exploding values.

4.3 Implementation

Figure 4 shows a detailed computation diagram, which relies heavily on large scale matrix multiplication. We can reshape each neuron and synaptic compartment output to batches and perform standard matrix multiplication to compute the full ENU-NN network in parallel (since our ENU parameters are shared). The connection matrix determines how neurons are connected, and in this case can be either random or have a fixed sparse topology. It determines how we broadcast the output of a neuron to each synapse connected to it (Broadcasting "copies" the output of neurons to different synapses as multiple synapses can be connected to the same neuron). However, the denser the connection matrix, the more synapses we have and thus the higher our computational cost.

ES can be performed efficiently on a Graphical Processing Unit (GPU) through multi-dimensional matrix multiplications, allowing us to evaluate thousands of mutated networks in parallel. The weight matrix in this case is 3 dimensional, and the 3th dimension stores each offspring's mutated weights.

Normal matrix multiplication is of the form $(N, K) \times (K, M) \rightarrow (N, M)$, where $(N$ is the batch size, K the number of inputs units and M the number of output units), in the 3D case we get $(P, N, K) \times (P, K, M) \rightarrow (P, N, M)$, where P is the number of offspring. PyTorch [42] was used for computing and evolving our ENU-NN, while the rest was implemented mainly in Numpy [43], including custom vectorized experimental environments.

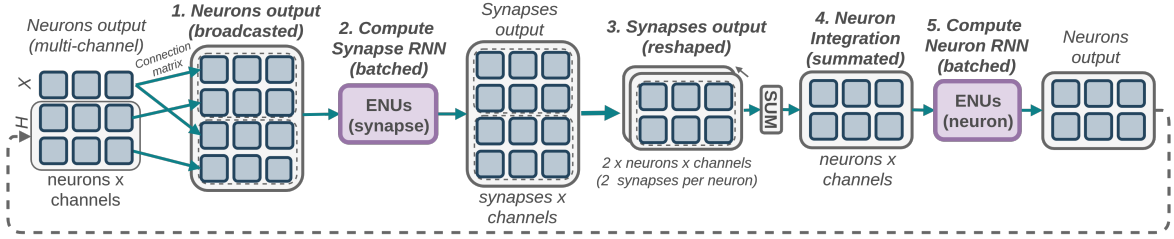


Figure 4: **Computation flow diagram.** The sensory input neurons X are concatenated with all the ENU neurons H to get our input batch. A connection matrix is then applied that broadcasts (copies) the neurons output to each connected synapse (1). On this resulting matrix we can then apply standard matrix multiplication and compute our synapse ENU output in parallel (2). We can reshape this and sum along the first axis, as we have the same number of synapses for each neuron (3). This gives us the integrated synaptic input to each neuron (4). Finally, we apply the neuron ENUs on this summed batch and obtain the output for each neuron in the ENU network (5).

4.4 Experiments

First we evolve an ENU to mimic the behaviour of a simple Integrate and Fire (IAF) and Spike Timing Dependent Plasticity (STDP) model. We then combine multiple neural and synaptic ENUs into a single network, and evolve reinforcement type learning behaviour purely through local dynamics.

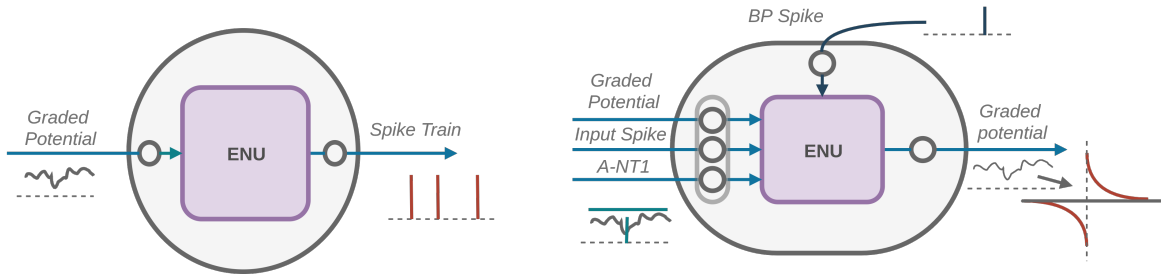


Figure 5: **IAF and STDP experimental setup.** For evolving the IAF ENU a single random graded potential is given as input, with as goal to mimic the IAF spiking model (left). In case of evolving the STDP rule (right) multiple input channels are used. The graded input potential, the input spike, the neuromodulation signal (A-NT1) and the backpropagating spike. The target is then to output the modified graded input potential matching the STDP rule.

Evolving Integrate and Fire Neurons (IAF) The ENU receives random uniform noise as graded input potential, and has to mimic the IAF model receiving the same input. Once the sum of the input potentials reaches a certain threshold, we reset the internal state of the IAF and output a spike (See also figure 5). The simplest IAF variant with linear summation is give below:

$$IAF(x, t, h) = \begin{cases} 0; h_t = h_{t-1} + x, & \text{if } h < th \\ 1; h=0, & \text{if } h \geq th \end{cases} \quad (2)$$

Where x is the input potential, t the current time, h the membrane potential and th the threshold.

Learning spike like behaviour can generally be difficult since the task is not directly differentiable. Optimizing the mean squared error between the IAF model and ENU output directly is infeasible since even a small shift in spikes would cause a decremental effect on the error. However, since we are using Evolution Strategies we can have more flexible loss functions and optimize the timing and intensity of each spike. Therefore, we optimize the Inter-Spike Interval (ISI) instead, which allows for incremental fitness improvements and gradually gets each spike to better match the desired timing of the IAF model. We also add an additional term that requires each spike to be as close as 1 as possible (matching the IAF spike).

Evolving Spike-timing dependent plasticity (STDP) For evolving a STDP type learning rule we require multiple input channels for the model. The basic equation for the STDP rule is as follows:

$$W(t) = \begin{cases} Ae^{(-\frac{t}{\tau})}, & \text{if } t > 0 \\ -Ae^{(\frac{t}{\tau})}, & \text{if } t < 0 \end{cases} \quad (3)$$

where A and τ are constants that determine the shape of the STDP function, t is the relative timing difference between the pre and post synaptic spike and W is the resulting synaptic weight value.

First we input random uniform noise as graded potential, we then also generate a random input spike from the pre-synaptic neuron at some time t_i and a random backpropagating spike from the post-synaptic neuron at t_b . We also use a neuromodulated variant of the standard STDP rule, which is known to play a role in synaptic learning [44, 45], and only allows STDP type updates to occur if a given input neurotransmitter is present (See figure 5).

In total we thus have 4 input channels, and a single output channel that is the transformed graded input potential multiplied by some weight value w . The weight value is updated according to the standard STDP rule which is dependent on the timing between the pre and post synaptic spike. The fitness is then the mean squared error between the desired STDP model output and ENU output.

Evolving Reinforcement Learning in a network of ENUs In order to evolve reinforcement type learning behaviour, we design an experimentally commonly used maze task [46] (see figure 3). The goal of the agent (i.e. the mouse), is to explore the maze and find food, once the agent eats the food he will receive a positive reward and be reset to the initial starting location. On the other hand, if the food was actually poisonous he will receive a negative reward (and also be reset). It is then up to the agent to remember where the food was and revisit the previous location. After the agent has eaten the non-poisonous food several times

there is a random chance that the food and poison will be switched. The agent thus has to evolve to learn to detect from it’s sensory input whether the food is poisonous or not, and can only know this by eating the food at least once to receive the associated negative or positive reward.

The fitness of the agent is determined by the reward obtained in the environment, the better it is at remembering where food is and at avoiding eating poison, the higher the fitness. We also add an additional term that decays an agent’s ‘energy’ every time step to encourage exploration. Eating food refreshes the energy again and gives the agent around 40 time steps to get new food. Once the agent runs out of energy it dies and will no longer be able to move or gather more rewards.

To this end we provide several inputs to the agent, one that detects the wall of the maze, one that detects green and one that detects red. These inputs are passed to the first channel of the connected synapse ENU. We also have a neuron that provides the resulting reward of eating the food or poison. Positive rewards go to the second channel, while negative rewards to the third channel. For the output the agent can either go forward, left or right or do nothing (if no output neuron activates).

The ENU network consists of 6 ENU neurons of which 3 are output neurons. Each neuron has 8 ENU synapses that sparsely connect to the other neurons (2 to the sensory neurons, 2 to the hidden neurons, 2 to the output neurons 1 to the reward neuron and 1 to itself). The neuron that has the highest output activity over 4 time steps determines the action of the agent (this allows for sufficient time of sensory input to propagate through the network). The output is taken from the first channel of the output neuron. We also add noise to the output gate of the ENU, since all ENUs share the same parameters and we want each neuron and synapse to behave slightly different upon initialization.

It is important to note that each generation the agents are reset and have no recollection of the past, each time they have to relearn which sensory neuron has what meaning and which output neuron performs what motor command. Also the network input and output neurons get randomly shuffled each generation, this avoids agents potentially exploiting the topology to learn fixed behaviour. We are thus evolving an ENU network to perform reinforcement type learning behaviour (evolving to learn), instead of directly learning fixed behaviour in the synaptic weights.

Experimental Details The experimental parameters were chosen through initial experimentation such that we achieved a good balance between computation time and performance. For optimization we used 1024 offspring and Gaussian mutations with standard deviation of 0.01, a learning rate of 1 and momentum of 0.9. For both the synaptic and neuronal compartments an ENU with a memory size of 32 units was used, with 16 output units (i.e. 16 output channels). The experiments were run on a single Titan V GPU with an i7 12-core CPU. In case of the T-maze experiment the mean fitness was taken across 8 random environments for each offspring. For evolving the IAF and STDP model the mean over 32 environments was taken (allowing us to plot a STDP type curve from multiple observations). 1 episode in the environment is 1 generation and each episode the ENU internal memory states are reset. In case of the IAF and STDP environment each episode lasts for 100 time steps, while for the T-maze experiment the episode lasted for 400 time steps.

5 Results

Single ENU: Evolving Integrate and Fire Neurons Results of evolving Integrate and Fire neurons for 3000 generations are shown in figure 6. It can be seen that the ENU is able to accurately mimic the underlying IAF model, properly integrating the received graded input, and outputting a spike at the right time. This shows an ENU is capable of evolving to approximate spiking behaviour.

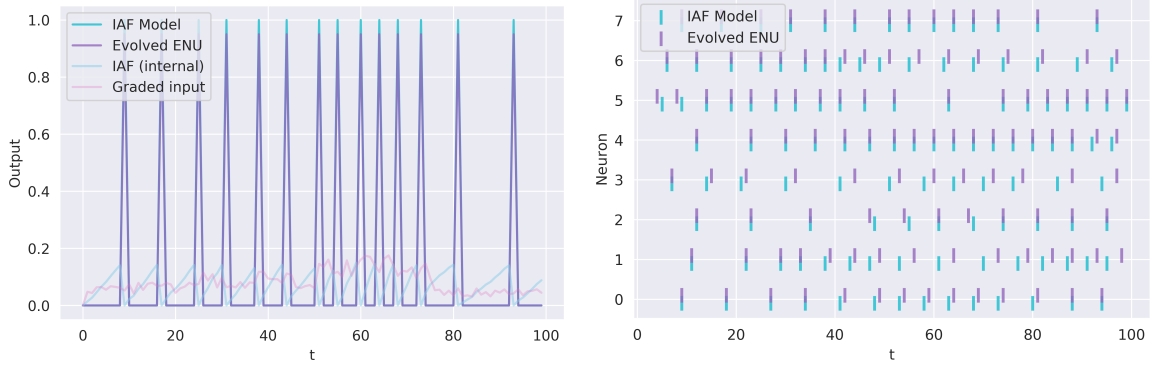


Figure 6: **Result of evolving Integrate and Fire neurons**, for a single input example (left) and of multiple inputs (right). The evolved ENU model closely matches the actual IAF model, properly integrating information and spiking at the right time once a threshold is reached. It is also able to correctly process random graded input potentials that result in slower or faster spiking patterns.

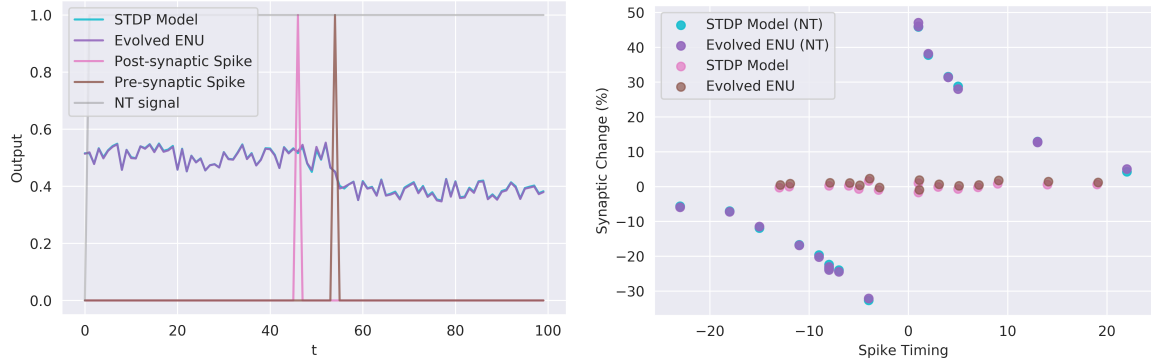


Figure 7: **Results of evolving neuromodulated Spike-Timing-Dependent Plasticity**, for a single input example (left) and multiple input observations (right). It evolved to update the synaptic weight only when the neurotransmitter signal is present in the input, it also evolved to update the weights relative to the timing of the pre and post synaptic spike, matching the original STDP function.

Single ENU: Evolving Neuromodulated STDP Figure 7 illustrates the results after evolving an STDP type learning rule for 10000 generations. If the pre-synaptic spike occurs after the post-synaptic spike the graded input potential reduces in output intensity as in the standard STDP rule (and vice versa). This demonstrates an ENU is capable of evolving

complex synaptic type learning through memorizing pre and post synaptic spikes in its internal memory and by storing and updating some dynamic parameter that changes how the incoming input is processed (analogous to a synaptic weight). It also had to evolve the multiplication operation of this dynamic parameter with the input, as initially the ENU has no such operation.

Network of ENUs: Evolving Reinforcement Learning Results for evolving Reinforcement Learning behaviour in a network of ENUs after 30000 generations in a T-maze environment are shown in Figure 8. A single generation is one episode in the environment (one simulation run) and last for 400 time steps. Each generation all the dynamic parameters of each synapse and neuron reset, meaning it always has to relearn which sensory neuron leads to a negative or positive rewards and which output neuron performs what action, as if it is reborn (a blank slate).

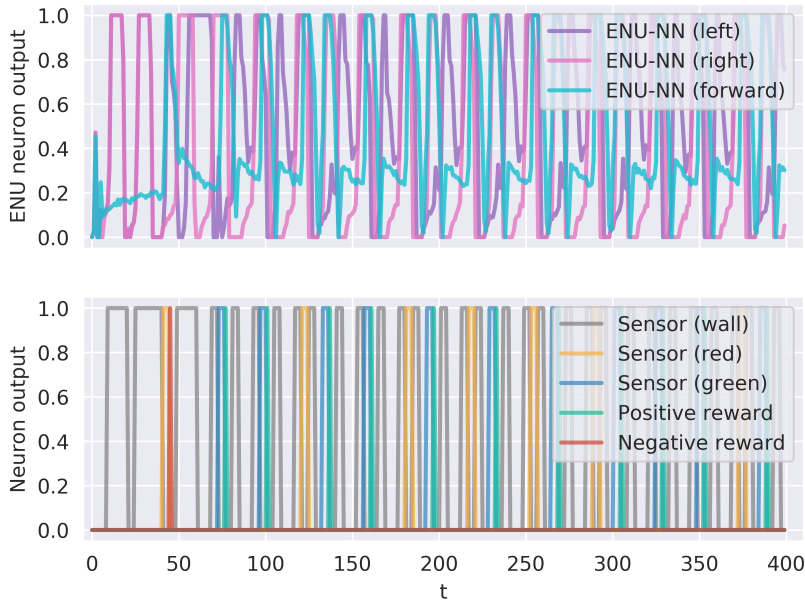


Figure 8: **Results of evolving Reinforcement learning in a network of ENUs**, showing both the output of the agent (top) and input (bottom). Spike like patterns on the output neuron evolved completely independently, even though we did not directly optimize for such behaviour. Furthermore, it can be seen that the agent performs one-shot learning when eating the poison, only eating it once and subsequently always avoiding it.

Interestingly we obtain spike like patterns on the output neurons even though we never strictly enforced it. We only optimize to maximize the reward in the environment. This could partially be explained by the sensory input neurons outputting spikes as well. However, the resulting output is not identical and it still had to evolve to process those input spikes and integrate them properly across the network and at the output.

An example of the steps taken by the agent are also given in figure 9. We can see that the agent at first detects red, eats the poison, and subsequently receives a negative reward. After that the agent learns to go the other way to get the food instead (detecting green). Once the food is switched it detects the poison but does not eat it (and so does not receive a negative reward), since it has now learned to turn around and obtain food on the other side instead.

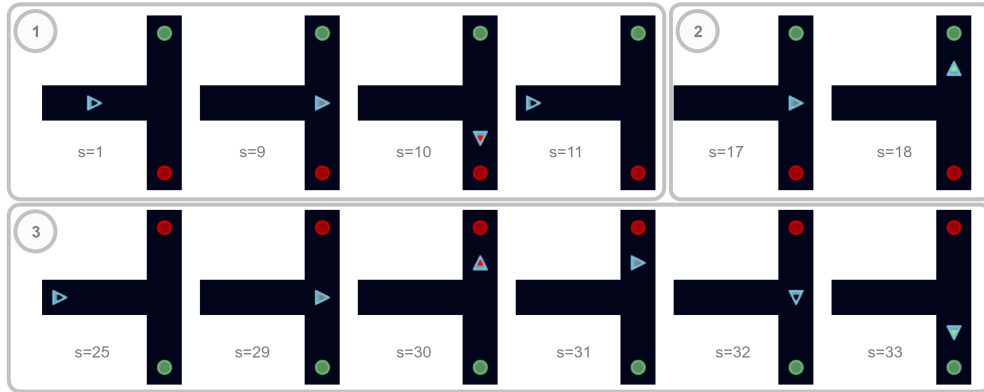


Figure 9: **Example of the ENU Network evolved behaviour.** It evolved to learn to remember which food gives a negative reward (1), avoid such reward in the future by taking a different route (2), and adapt to turn around once the food and poison are switched (3).

The performance (and fitness) is measured by how good the agent is at getting food and avoiding poison. The ENU networks of the agents with random mutations better able to learn and store the right information in the internal memory of each individual ENU will thus have a higher fitness, able to learn quicker from negative or positive rewards. Evolution Strategies then updates the shared parameters of the ENU gates in that direction. This leads it to improve its ability to process information in the ENU network and learn from reward feedback. It has to evolve in a way such that all ENU compartments with the shared gate parameters can cooperate and update their unique internal state to achieve an higher fitness. It thus evolves a type of universal function the same across all neurons and synapses that updates their internal dynamic parameters depending on the input.

The agent is able to learn through the evolved ENU network that the sensory neuron that detects the color red and the activation of the output neuron that results in eating it lead to a negative reward. This means that it evolved a mechanism for updating the ENU synapse between those neurons to have an inhibitory effect, such that next time it detects red, the previous action that led to eating the poison is not chosen (it is inhibited), and instead another action is taken.

Comparison to other models Standard Deep Neural Networks (DNNs) use only a single value for the synaptic weight parameter (which is static) and neuron output. The same holds for spiking neural network models with STDP type learning rules [24]. The activation function and update rules of these models are also fixed. Similarly, previous meta-learning approaches that learn to learn synaptic update rules still use the same fundamental underlying design of DNN computations [35, 36]. In our case, each dynamic parameter is a vector capable of storing multiple values that can influence the "operation" performed by the synaptic and neuronal ENU compartments. We evolve the ability of each ENU to update those dynamic parameters based on the reward obtained and local input received, unlike DNNs that use backpropagation over the entire network given some learning signal. In order for learning to occur in our ENU network, the synapse ENU has to evolve the performed operation on the input over multiple channels, which is not necessarily a multiplication as in standard DNNs.

Each generation our "dynamic" parameters also get reset, while in DNNs the parameters of the weights persist. DNNs are thus not learning-to-learn but learn fixed behaviour instead, only applicable to the current environment. Additionally, to learn from reward signals using deep neural networks, reinforcement based learning methods like Deep Q-learning [47] or Actor-Critic models [48] are required. In these methods however the agent requires millions of observations to update the actions to take in that specific environment.

In our case the final evolved model is able to perform one-shot learning from a single experience to modify its behaviour. This is because a learning rule evolves that is able to update the synaptic compartments given the reward, which consequently enables it to alter the way information flows through the synapses and neurons (over multiple channels). It also implicitly evolved random exploration like behaviour to seek out rewards, since this is required to achieve a higher overall fitness. Having multiple channels allows for more flexibility in learning and information processing behaviour, and could also explain why biological networks might have evolved to use so many types of different neurotransmitters [49, 50].

6 Discussion

We showed that we were able to successfully train the proposed Evolvable Neural Units (ENUs) to mimic Integrate and Fire Neurons and Spike-Timing Dependent Plasticity. This demonstrated that in principle an ENU is flexible enough to evolve neural like dynamics and complex processing mechanisms. We then evolve an interconnected network where each synapse and neuron in the network is such an evolvable neural unit. This network of ENUs can be evolved to solve a T-maze environment task, which results in an agent capable of reinforcement type learning behaviour. This environment requires the agent to evolve the ability to learn to remember the color and location of food and poison depending on the reward obtained, change its behaviour to avoid such poison, and dynamically adapt to changes in the environment. These neuronal and synaptic ENUs in the network thus have to learn to work together cooperatively through local dynamics to maximize their overall fitness and survival, dynamically updating the way they integrate and process information.

Interesting future directions include simulating and evolving entire cortical columns by evolving smaller networks that are robust to changes in network size, allowing us to scale up the number of neurons and synaptic compartments after the evolutionary process. Furthermore, we could use separate ENUs to evolve e.g. dendrites, axons or the behaviour of different cell types seen across cortical layers. It is however still uncertain what aspects of biological neural networks are actually necessary or just byproducts of evolutionary processes. It also remains an open question what type of environment would be required in order for higher level intelligence to emerge, and research into open-ended evolution and artificial life might therefore be critical to achieve such intelligence.

The Evolvable Neural Units presented in this paper offer a new direction for potentially more powerful and biologically realistic neural networks, and could ultimately not only lead to a greater understanding of neural processing, but also to an artificially intelligent system that can learn and act across a wide variety of domains.

References

- [1] Demis Hassabis, Dharshan Kumaran, Christopher Summerfield, and Matthew Botvinick. Neuroscience-inspired artificial intelligence. *Neuron*, 95(2):245–258, 2017.
- [2] Nelson Spruston. Pyramidal neurons: dendritic structure and synaptic integration. *Nature Reviews Neuroscience*, 9(3):206, 2008.
- [3] João Sacramento, Rui Ponte Costa, Yoshua Bengio, and Walter Senn. Dendritic cortical microcircuits approximate the backpropagation algorithm. In *Advances in Neural Information Processing Systems*, pages 8721–8732, 2018.
- [4] Larry F Abbott and Sacha B Nelson. Synaptic plasticity: taming the beast. *Nature neuroscience*, 3(11s):1178, 2000.
- [5] Christoph Börgers and Nancy Kopell. Synchronization in networks of excitatory and inhibitory neurons with sparse, random connectivity. *Neural computation*, 15(3):509–538, 2003.
- [6] Juergen Fell and Nikolai Axmacher. The role of phase synchronization in memory processes. *Nature reviews neuroscience*, 12(2):105, 2011.
- [7] Larry F Abbott. Lapiques introduction of the integrate-and-fire model neuron (1907). *Brain research bulletin*, 50(5-6):303–304, 1999.
- [8] Anil K Jain, Jianchang Mao, and K Moidin Mohiuddin. Artificial neural networks: A tutorial. *Computer*, 29(3):31–44, 1996.
- [9] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.
- [10] William A Catterall. Structure and function of voltage-gated ion channels. *Annual review of biochemistry*, 64(1):493–531, 1995.
- [11] Shelly B Flagel, Jeremy J Clark, Terry E Robinson, Leah Mayo, Alayna Czuj, Ingo Willuhn, Christina A Akers, Sarah M Clinton, Paul EM Phillips, and Huda Akil. A selective role for dopamine in stimulus–reward learning. *Nature*, 469(7328):53, 2011.
- [12] DAVID A McCormick. Gaba as an inhibitory neurotransmitter in human cerebral cortex. *Journal of neurophysiology*, 62(5):1018–1027, 1989.
- [13] Irwin B Levitan and Leonard K Kaczmarek. *The neuron: cell and molecular biology*. Oxford University Press, USA, 2015.
- [14] Maiken Nedergaard, Takahiro Takano, and Anker J Hansen. Beyond the role of glutamate as a neurotransmitter. *Nature Reviews Neuroscience*, 3(9):748, 2002.
- [15] Vernon B Mountcastle. The columnar organization of the neocortex. *Brain: a journal of neurology*, 120(4):701–722, 1997.
- [16] Ronald D Vale. The molecular motor toolbox for intracellular transport. *Cell*, 112(4):467–480, 2003.

- [17] Sabine Hilfiker, Vincent A Pieribone, Andrew J Czernik, Hung-Teh Kao, George J Augustine, and Paul Greengard. Synapsins as regulators of neurotransmitter release. *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences*, 354(1381):269–279, 1999.
- [18] Peter J Hollenbeck and William M Saxton. The axonal transport of mitochondria. *Journal of cell science*, 118(23):5411–5419, 2005.
- [19] Graham L Collingridge, John TR Isaac, and Yu Tian Wang. Receptor trafficking and synaptic plasticity. *Nature Reviews Neuroscience*, 5(12):952, 2004.
- [20] Adam Kepecs, Xiao-Jing Wang, and John Lisman. Bursting neurons signal input slope. *Journal of Neuroscience*, 22(20):9053–9062, 2002.
- [21] Adam Kepecs and John Lisman. Information encoding and computation with spikes and bursts. *Network: Computation in neural systems*, 14(1):103–118, 2003.
- [22] Alan L Hodgkin and Andrew F Huxley. A quantitative description of membrane current and its application to conduction and excitation in nerve. *The Journal of physiology*, 117(4):500–544, 1952.
- [23] Donald Olding Hebb and DO Hebb. *The organization of behavior*, volume 65. Wiley New York, 1949.
- [24] Wolfgang Maass. Networks of spiking neurons: the third generation of neural network models. *Neural networks*, 10(9):1659–1671, 1997.
- [25] Guillaume Bellec, Darjan Salaj, Anand Subramoney, Robert Legenstein, and Wolfgang Maass. Long short-term memory and learning-to-learn in networks of spiking neurons. In *Advances in Neural Information Processing Systems*, pages 787–797, 2018.
- [26] Kurt Hornik. Approximation capabilities of multilayer feedforward networks. *Neural networks*, 4(2):251–257, 1991.
- [27] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [28] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- [29] Thomas Back. *Evolutionary algorithms in theory and practice: evolution strategies, evolutionary programming, genetic algorithms*. Oxford university press, 1996.
- [30] Kenneth O Stanley, Jeff Clune, Joel Lehman, and Risto Miikkulainen. Designing neural networks through neuroevolution. *Nature Machine Intelligence*, 1(1):24–35, 2019.
- [31] Tim Salimans, Jonathan Ho, Xi Chen, Szymon Sidor, and Ilya Sutskever. Evolution strategies as a scalable alternative to reinforcement learning. *arXiv preprint arXiv:1703.03864*, 2017.

- [32] Jürgen Schmidhuber. *Evolutionary principles in self-referential learning, or on learning how to learn: the meta-meta-... hook*. PhD thesis, Technische Universität München, 1987.
- [33] Samy Bengio, Yoshua Bengio, Jocelyn Cloutier, and Jan Gecsei. On the optimization of a synaptic learning rule. In *Preprints Conf. Optimality in Artificial and Biological Neural Networks*, pages 6–8. Univ. of Texas, 1992.
- [34] Jürgen Schmidhuber. Learning to control fast-weight memories: An alternative to dynamic recurrent networks. *Neural Computation*, 4(1):131–139, 1992.
- [35] Marcin Andrychowicz, Misha Denil, Sergio Gomez, Matthew W Hoffman, David Pfau, Tom Schaul, Brendan Shillingford, and Nando De Freitas. Learning to learn by gradient descent by gradient descent. In *Advances in neural information processing systems*, pages 3981–3989, 2016.
- [36] David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. *arXiv preprint arXiv:1609.09106*, 2016.
- [37] Thomas Back, Frank Hoffmeister, and Hans-Paul Schwefel. A survey of evolution strategies. In *Proceedings of the fourth international conference on genetic algorithms*, volume 2. Morgan Kaufmann Publishers San Mateo, CA, 1991.
- [38] Hans-Georg Beyer and Hans-Paul Schwefel. Evolution strategies—a comprehensive introduction. *Natural computing*, 1(1):3–52, 2002.
- [39] Daan Wierstra, Tom Schaul, Tobias Glasmachers, Yi Sun, Jan Peters, and Jürgen Schmidhuber. Natural evolution strategies. *The Journal of Machine Learning Research*, 15(1):949–980, 2014.
- [40] Joel Lehman, Jay Chen, Jeff Clune, and Kenneth O Stanley. Es is more than just a traditional finite-difference approximator. In *Proceedings of the Genetic and Evolutionary Computation Conference*, pages 450–457. ACM, 2018.
- [41] Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In *International conference on machine learning*, pages 1139–1147, 2013.
- [42] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.
- [43] Travis E Oliphant. *A guide to NumPy*, volume 1. Trelgol Publishing USA, 2006.
- [44] Verena Pawlak, Jeffery R Wickens, Alfredo Kirkwood, and Jason ND Kerr. Timing is not everything: neuromodulation opens the stdp gate. *Frontiers in synaptic neuroscience*, 2:146, 2010.
- [45] Nicolas Frémaux and Wulfram Gerstner. Neuromodulated spike-timing-dependent plasticity, and theory of three-factor learning rules. *Frontiers in neural circuits*, 9:85, 2016.
- [46] Robert MJ Deacon and J Nicholas P Rawlins. T-maze alternation in the rodent. *Nature protocols*, 1(1):7, 2006.

- [47] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.
- [48] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937, 2016.
- [49] Jean M Lauder. Neurotransmitters as growth regulatory signals: role of receptors and second messengers. *Trends in neurosciences*, 16(6):233–240, 1993.
- [50] Kenji Doya. Metalearning and neuromodulation. *Neural Networks*, 15(4-6):495–506, 2002.