# Accepted Manuscript

Sigmoid-weighted linear units for neural network function approximation in reinforcement learning

Stefan Elfwing, Eiji Uchibe, Kenji Doya

# Sigmoid-Weighted Linear Units for Neural Network Function Approximation in Reinforcement Learning

Stefan Elfwing[a,*], Eiji Uchibe[a,b], Kenji Doya[b]

[a]*Dept. of Brain Robot Interface, ATR Computational Neuroscience Laboratories, 2-2-2 Hikaridai, Seikacho, Soraku-gun, Kyoto 619-0288, Japan*
[b]*Okinawa Institute of Science and Technology Graduate University, 1919-1 Tancha, Onna-son, Okinawa 904-0495, Japan*

**Abstract**

In recent years, neural networks have enjoyed a renaissance as function approximators in reinforcement learning. Two decades after Tesauro's TD-Gammon achieved near top-level human performance in backgammon, the deep reinforcement learning algorithm DQN achieved human-level performance in many Atari 2600 games. The purpose of this study is twofold. First, we propose two activation functions for neural network function approximation in reinforcement learning: the sigmoid-weighted linear unit (SiLU) and its derivative function (dSiLU). The activation of the SiLU is computed by the sigmoid function multiplied by its input. Second, we suggest that the more traditional approach of using on-policy learning with eligibility traces, instead of experience replay, and softmax action selection can be competitive with DQN, without the need for a separate target network. We validate our proposed approach by, first, achieving new state-of-the-art results in both stochastic SZ-Tetris and Tetris with a small $10 \times 10$ board, using TD($\lambda$) learning and shallow dSiLU network agents, and, then, by outperforming DQN in the Atari 2600 domain by using a deep Sarsa($\lambda$) agent with SiLU and dSiLU hidden units.

*Keywords:* reinforcement learning, sigmoid-weighted linear unit, function approximation, tetris, atari 2600, deep learning

*Corresponding author
*Email addresses:* elfwing@atr.jp (Stefan Elfwing), uchibe@atr.jp (Eiji Uchibe), doya@oist.jp (Kenji Doya)

## 1. Introduction

Neural networks have enjoyed a renaissance as function approximators in reinforcement learning (Sutton and Barto, 1998) in recent years. The DQN algorithm (Mnih et al., 2015), which combines Q-learning with a deep neural network, experience replay, and a separate target network, achieved human-level performance in many Atari 2600 games. Since the development of the DQN algorithm, there have been several proposed improvements, both to DQN specifically and deep reinforcement learning in general. Van Hasselt et al. (2015) proposed double DQN to reduce overestimation of the action values in DQN and Schaul et al. (2016) developed a framework for more efficient replay by prioritizing experiences of more important state transitions. Wang et al. (2016) proposed the dueling network architecture for more efficient learning of the action value function by separately estimating the state value function and the advantages of each action. Mnih et al. (2016) proposed a framework for asynchronous learning by multiple agents in parallel, both for value-based and actor-critic methods. To date, the most impressive application of using deep reinforcement learning is AlphaGo (Silver et al., 2016, 2017), which has achieved superhuman performance in the ancient board game Go.

The purpose of this study is twofold. First, motivated by the high performance of the expected energy restricted Boltzmann machine (EE-RBM) in our earlier studies (Elfwing et al., 2015, 2016), we propose two activation functions for neural network function approximation in reinforcement learning: the sigmoid-weighted linear unit (SiLU) and its derivative function (dSiLU). The activation of the SiLU is computed by the sigmoid function multiplied by its input. After we first proposed the SiLU (Elfwing et al., 2017), Ramachandran et al. (2017) recently performed a comprehensive comparison between the SiLU, the rectifier linear unit (ReLU; Hahnloser et al., 2000), and 6 other activation functions in the supervised learning domain. They found that the SiLU consistently outperformed the other activation functions when tested in 3 deep architectures on CIFAR-10/100 (Krizhevsky, 2009), in 5 deep architectures on ImageNet (Deng et al., 2009), and on 4 test sets for English-to-German machine translation.

Second, we suggest that the more traditional approach of using on-policy learning

2

with eligibility traces, instead of experience replay, and softmax action selection with simple annealing can be competitive with DQN, without the need for a separate target network. Our approach is something of a throwback to the approach used by Tesauro (1994) to develop TD-Gammon more than two decades ago. Using a neural network function approximator and TD($\lambda$) learning (Sutton, 1988), TD-Gammon reached near top-level human performance in backgammon, which to this day remains one of the most impressive application of reinforcement learning.

35

To evaluate our proposed approach, we first test the performance of shallow network agents with SiLU, ReLU, dSiLU, and sigmoid hidden units in stochastic SZ-Tetris, which is a simplified but difficult version of Tetris. The best agent, the dSiLU network agent, improves the average state-of-the-art score by 20 %. In stochastic SZ-Tetris, we also train deep network agents using raw board configurations as states. An agent with SiLUs in the convolutional layers and dSiLUs in the fully-connected layer (SiLU-dSiLU) outperforms the previous state-of-the-art average final score. We thereafter train a dSiLU network agent in standard Tetris with a smaller, 10×10, board size, achieving a state-of-the-art score in this more competitive version of Tetris as well. We then test a deep SiLU-dSiLU network agent in the Atari 2600 domain. It improves the mean DQN normalized scores achieved by DQN and double DQN by 232 % and 161 %, respectively, in 12 unbiasedly selected games. We finally analyze the ability of on-policy value-based reinforcement learning to accurate estimate the expected discounted returns and the importance of softmax action selection for the games where our proposed agents performed particularly well.

40

45

50

## 2. Method

### 2.1. TD($\lambda$) and Sarsa($\lambda$)

55

In this study, we use two reinforcement learning algorithms: TD($\lambda$) (Sutton, 1988) and Sarsa($\lambda$) (Rummery and Niranjan, 1994; Sutton, 1996). TD($\lambda$) learns an estimate of the state-value function, $V^{\pi}$, and Sarsa($\lambda$) learns an estimate of the action-value function, $Q^{\pi}$, while the agent follows policy $\pi$. If the approximated value functions, $V_t \approx V^{\pi}$ and $Q_t \approx Q^{\pi}$, are parameterized by the parameter vector $\boldsymbol{\theta}_t$, then the

3

gradient-descent learning update of the parameters is computed by

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha\delta_t\boldsymbol{e}_t, \tag{1}$$

where the TD-error, $\delta_t$, is

$$\delta_t = r_t + \gamma V_t(s_{t+1}) - V_t(s_t) \tag{2}$$

for TD($\lambda$) and

$$\delta_t = r_t + \gamma Q_t(s_{t+1}, a_{t+1}) - Q_t(s_t, a_t) \tag{3}$$

for Sarsa($\lambda$). The eligibility trace vector, $\boldsymbol{e}_t$, is

$$\boldsymbol{e}_t = \gamma\lambda\boldsymbol{e}_{t-1} + \nabla_{\boldsymbol{\theta}_t}V_t(s_t), \ \boldsymbol{e}_0 = \boldsymbol{0}, \tag{4}$$

for TD($\lambda$) and

$$\boldsymbol{e}_t = \gamma\lambda\boldsymbol{e}_{t-1} + \nabla_{\boldsymbol{\theta}_t}Q_t(s_t, a_t), \ \boldsymbol{e}_0 = \boldsymbol{0}, \tag{5}$$

for Sarsa($\lambda$). Here, $s_t$ is the state at time $t$, $a_t$ is the action selected at time $t$, $r_t$ is the reward for taking action $a_t$ in state $s_t$, $\alpha$ is the learning rate, $\gamma$ is the discount factor of future rewards, $\lambda$ is the trace-decay rate, and $\nabla_{\boldsymbol{\theta}_t}V_t$ and $\nabla_{\boldsymbol{\theta}_t}Q_t$ are the vectors of partial derivatives of the function approximators with respect to each component of $\boldsymbol{\theta}_t$.

### 2.2. Sigmoid-weighted Linear Units

We proposed the EE-RBM as a function approximator in reinforcement learning (Elfwing et al., 2016). In the case of state-value based learning, given a state vector $\mathbf{s}$, an EE-RBM approximates the state-value function $V$ by the negative expected energy of an RBM (Smolensky, 1986; Freund and Haussler, 1992; Hinton, 2002) network:

$$
\begin{align}
V(\mathbf{s}) &= \sum_k z_k\sigma(z_k) + \sum_i b_i s_i, \tag{6} \\
z_k &= \sum_i w_{ik}s_i + b_k, \tag{7} \\
\sigma(x) &= \frac{1}{1+e^{-x}}. \tag{8}
\end{align}
$$

Here, $z_k$ is the input to hidden unit $k$, $\sigma(\cdot)$ is the sigmoid function, $b_i$ is the bias weight for input unit $s_i$, $w_{ik}$ is the weight connecting state $s_i$ and hidden unit $k$, and $b_k$ is the

4

bias weight for hidden unit $k$. Note that Equation 6 can be regarded as the output of a one-hidden layer feedforward neural network with hidden unit activations computed by $z_k \sigma(z_k)$ and with uniform output weights of one. In this study, motivated by the high performance of the EE-RBM in both the classification (Elfwing et al., 2015) and the reinforcement learning (Elfwing et al., 2016) domains, we propose the SiLU as an activation function for neural network function approximation in reinforcement learning. The activation $a_k$ of the $k^{\text{th}}$ SiLU for input $z_k$ is computed by the sigmoid function multiplied by its input (i.e., equal to the contribution from a hidden node to the value function in an EE-RBM):

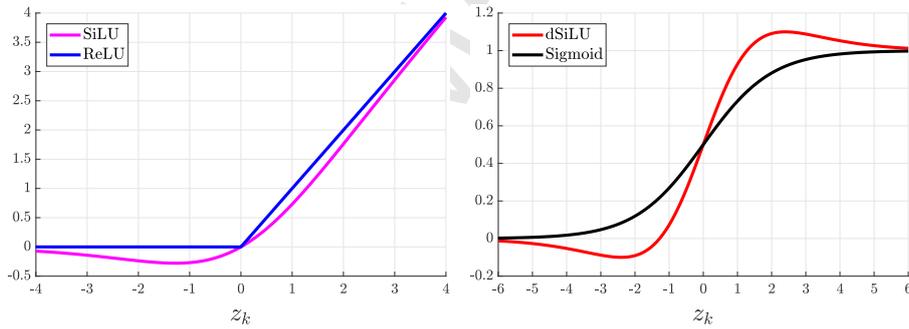$$a_k(z_k) = z_k \sigma(z_k). \tag{9}$$



Figure 1: The activation functions of the SiLU and the ReLU (left panel), and the dSiLU and the sigmoid unit (right panel).

For $z_k$-values of large magnitude, the activation of the SiLU is approximately equal to the activation of the ReLU (see left panel in Figure 1), i.e., the activation is approximately equal to zero for large negative $z_k$-values and approximately equal to $z_k$ for large positive $z_k$-values. Unlike the ReLU (and other commonly used activation units such as sigmoid and tanh units), the activation of the SiLU is not monotonically increasing. Instead, it has a global minimum value of approximately $-0.28$ for $z_k \approx -1.28$. An attractive feature of the SiLU is that it has a self-stabilizing property, which we demonstrated experimentally in Elfwing et al. (2015). The global minimum, where the

5

derivative is zero, functions as a "soft floor" on the weights that serves as an implicit regularizer that inhibits the learning of weights of large magnitudes.

In Elfwing et al. (2015), we discovered that the derivative function of the SiLU (i.e., the derivative of the contribution from a hidden node to the output in an EE-RBM) looks like a steeper and "overshooting" version of the sigmoid function. In this study, we call this function the dSiLU and we propose it as a competitive alternative to the sigmoid function in neural network function approximation in reinforcement learning. The activation of the dSiLU is computed by the derivative of the SiLU (see right panel in Figure 1):

$$a_k(z_k) = \sigma(z_k)\left(1 + z_k(1 - \sigma(z_k))\right). \tag{10}$$

The dSiLU has a maximum value of approximately $1.1$ and a minimum value of approximately $-0.1$ for $z_k \approx \pm 2.4$, i.e., the solutions to the equation $z_k = -\log\left((z_k - 2)/(z_k + 2)\right)$.

*2.3. Action selection*

We use softmax action selection with a Boltzmann distribution in all experiments. For Sarsa($\lambda$), the probability to select action $a$ in state $s$ is defined as

$$\pi(a|s) = \frac{\exp(Q(s,a)/\tau)}{\sum_b \exp(Q(s,b)/\tau)}. \tag{11}$$

For the model-based TD($\lambda$) algorithm, we select an action $a$ in state $s$ that leads to the next state $s'$ with a probability defined as

$$\pi(a|s) = \frac{\exp(V(f(s,a))/\tau)}{\sum_b \exp(V(f(s,b))/\tau)}. \tag{12}$$

Here, $f(s,a)$ returns the next state $s'$ according to the deterministic state transition dynamics and $\tau$ is the temperature that controls the trade-off between exploration and exploitation. We used hyperbolic annealing of the temperature and the temperature was decreased after every episode $i$:

$$\tau(i) = \frac{\tau_0}{1 + \tau_k i}. \tag{13}$$

Here, $\tau_0$ is the initial temperature and $\tau_k$ controls the rate of annealing.

## 3. Experiments

### 3.1. SZ-Tetris

Szita and Szepesvári (2010) proposed stochastic SZ-Tetris (Burgiel, 1997) as a benchmark for reinforcement learning that preserves the core challenges of standard Tetris but allows faster evaluation of different strategies due to shorter episodes by removing easier tetrominos. Stochastic SZ-Tetris is played on a board of standard Tetris size with a width of 10 and a height of 20. In each time step, either an S-shaped tetromino or a Z-shaped tetromino appears with equal probability. The agent selects a rotation (lying or standing) and a horizontal position within the board. In total, there are 17 possible actions for each tetromino (9 standing and 8 lying horizontal positions). After the action selection, the tetromino drops down the board, stopping when it hits another tetromino or the bottom of the board. If a row is completed, then it disappears. The agent gets a score of +1 point for each completed row. An episode ends when a tetromino does not fit within the board.

The standard learning approach for Tetris has been to use a model-based setting and define the evaluation function or state-value function as the linear combination of hand-coded features. Value-based reinforcement learning algorithms have a lousy track record using this approach. In regular Tetris, their reported performance levels are many magnitudes lower than black-box methods such as the cross-entropy (CE) method and evolutionary approaches. In stochastic SZ-Tetris, the reported scores for a wide variety of reinforcement learning algorithms are either approximately zero (Szita and Szepesvári, 2010) or in the single digits [1].

Value-based reinforcement learning has had better success in stochastic SZ-Tetris when using non-linear neural network based function approximators. Faußer and Schwenker (2013) achieved a score of about 130 points using a shallow neural network function approximator with sigmoid hidden units. They improved the result to about 150 points by using an ensemble approach consisting of ten neural networks. We achieved an average score of about 200 points using three different neural network function ap-

---

[1]http://barbados2011.rl-community.org/program/SzitaTalk.pdf

proximators: an EE-RBM, a free energy RBM, and a standard neural network with sigmoid hidden units (Elfwing et al., 2016). Jaskowski et al. (2015) achieved the current state-of-the-art results using systematic n-tuple networks as function approximators: average scores of 220 and 218 points achieved by the evolutionary VD-CMA-ES method and TD-learning, respectively, and the best mean score in a single run of 295 points achieved by TD-learning.

In this study, we use the TD($\lambda$) algorithm and softmax action selection to compare the performance of different hidden activation units in two learning settings: 1) shallow network agents with one hidden layer using hand-coded state features and 2) deep network agents using raw board configurations as states, i.e., a state node is set to one if the corresponding board cell was occupied by a tetromino and set to zero otherwise.
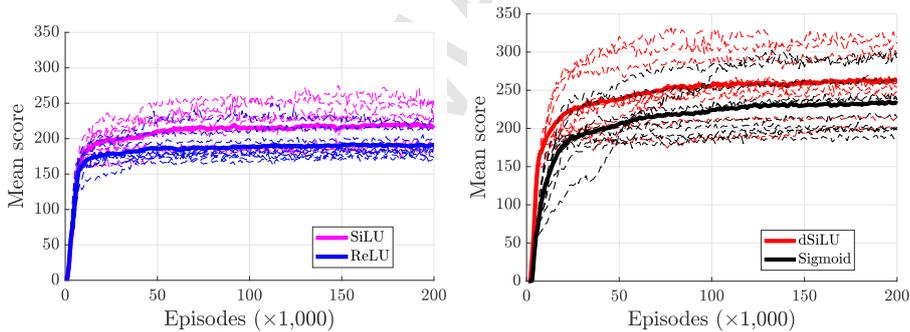


Figure 2: Learning curves in stochastic SZ-Tetris for the four types of shallow neural network agents. The figure shows the average scores over ten separate runs (tick solid lines) and the scores of individual runs (thin dashed lines). The mean scores were computed over every 1,000 episodes.

In the setting with state features, we trained shallow network agents with SiLU, ReLU, dSiLU, and sigmoid hidden units. We used the same experimental setup as used in our earlier work (Elfwing et al., 2016). The networks consisted of one hidden layer with 50 hidden units and a linear output layer. The features were similar to the original 21 features proposed by Bertsekas and Ioffe (1996), except for not including the maximum column height and using the differences in column heights instead of the absolute differences. The length of the binary state vector was 460. The shallow network agents were trained for 200,000 episodes and the experiments were repeated

8

for ten separate runs for each type of activation unit.

165     In the deep reinforcement learning setting, we used a deep network architecture consisting of two convolutional layers with 15 and 50 filters of size $5 \times 5$ using a stride of 1, a fully-connected layer with 250 units, and a linear output layer. Both convolutional layers were followed by max-pooling layers with pooling windows of size $3 \times 3$ using a stride of 2. We compared networks with SiLUs in both the con-

170 volutional and fully-connected layers (SiLU-SiLU) with networks with ReLUs in all hidden layers (ReLU-ReLU). Based on the high performance of the shallow networks with dSiLUs in the hidden layer, we also tested a deep network with dSiLUs in the last, fully-connected, hidden layer, combined with SiLUs in the convolutional layers (SiLU-dSiLU). The deep network agents were trained for 200,000 episodes and the

175 experiments were repeated for five separate runs for each type of network.
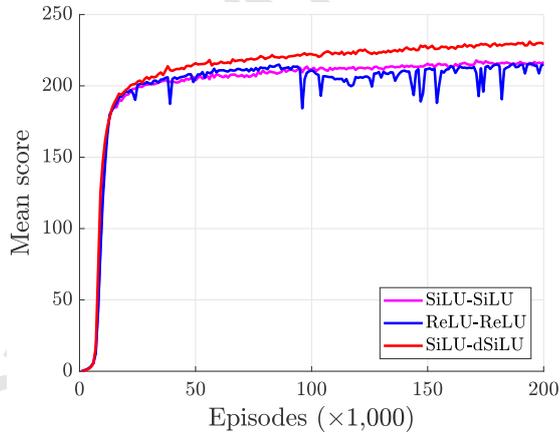


Figure 3: Average Learning curves in stochastic SZ-Tetris for the three types of deep neural network agents. The figure shows the average scores over five separate runs, computed over every 1,000 episodes.

We used the following reward function (proposed by Faußer and Schwenker (2013)):

$$r(\mathbf{s}) = e^{-(\text{number of holes in } \mathbf{s})/33}. \tag{14}$$

We set $\gamma$ to 0.99, $\lambda$ to 0.55, $\tau_0$ to 0.5, and $\tau_k$ to 0.00025. We used a rough grid-like search to find appropriate values of the learning rate $\alpha$ and it was determined to

9

<sub>180</sub> be 0.001 for the four shallow network agents and 0.0001 for the three deep network agents.

Table 1: Average scores ($\pm$ standard deviations) achieved in stochastic SZ-Tetris, computed over the final 1,000 episodes for all runs and the best single runs.

| Network | Final average score | Final best score |
|---------|---------------------|------------------|
| Shallow networks | | |
| SiLU | $214 \pm 74$ | $253 \pm 83$ |
| ReLU | $191 \pm 58$ | $227 \pm 76$ |
| dSiLU | $\mathbf{263 \pm 80}$ | $\mathbf{320 \pm 87}$ |
| Sigmoid | $232 \pm 75$ | $293 \pm 73$ |
| Deep networks | | |
| SiLU-SiLU | $217 \pm 53$ | $219 \pm 54$ |
| ReLU-ReLU | $215 \pm 54$ | $217 \pm 52$ |
| SiLU-dSiLU | $\mathbf{229 \pm 55}$ | $\mathbf{235 \pm 54}$ |

Figure 2 shows the average learning curves as well as learning curves for the individual runs for the shallow networks, Figure 3 shows the average learning curves for the deep networks, and the final results are summarized in Table 1. The results show
<sub>185</sub> significant differences ($p < 0.0001$) in final average score between all four shallow agents. The networks with bounded hidden units (dSiLU and sigmoid) outperformed the networks with unbounded units (SiLU and ReLU), the SiLU network outperformed the ReLU network, and the dSiLU network outperformed the sigmoid network. The final average score (best score) of 263 (320) points achieved by the dSiLU network
<sub>190</sub> agent is a new state-of-the-art score, improving the previous best performance by 43 (25) points or 20 % (8 %). In the deep learning setting, the SiLU-dSiLU network significantly ($p < 0.0001$) outperformed the other two networks and the average final score of 229 points is better than the previous state-of-the-art of 220 points. There were no significant difference ($p = 0.32$) between the final performance of the SiLU-SiLU
<sub>195</sub> network and the ReLU-ReLU network.

### 3.2. 10×10 Tetris

The result achieved by the dSiLU network agent in stochastic SZ-Tetris is impressive, but we cannot compare the result with the methods that have achieved the highest performance levels in standard Tetris, because those methods have not been applied
200 to stochastic SZ-Tetris. Furthermore, it is not feasible to apply our method to Tetris with a standard board height of 20, because of the prohibitively long learning time. The current state-of-the-art for a single run of an algorithm, achieved by the CBMPI algorithm (Gabillon et al., 2013; Scherrer et al., 2015), is a mean score of 51 million cleared lines. However, for the best methods applied to Tetris, there are reported results
205 for a smaller, 10×10, Tetris board, and in this case the learning time for our method is long, but not prohibitively so.
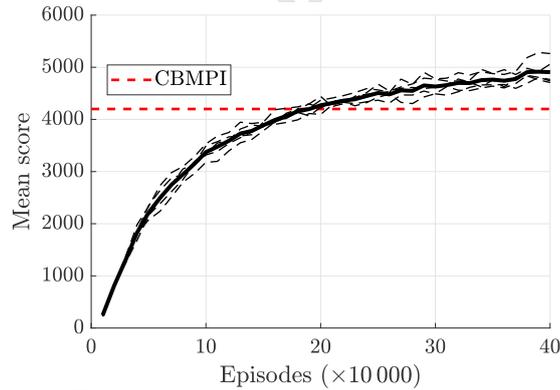


Figure 4: Learning curves for a dSiLU network agent with 250 hidden nodes in 10×10 Tetris. The figure shows the average score over five separate runs (tick solid lines) and the scores of individual runs (thin dashed lines). The red dashed line show the previous best average score of 4,200 points achieved by the CBMPI algorithm.

10×10 Tetris is played with the standard seven tetrominos and the numbers of actions are 9 for the block-shaped tetromino, 17 for the S-, Z-, and stick-shaped tetrominos, and 34 for the J-, L- and T-shaped tetrominos. In each time step, the agent gets
210 a score equal to the number of completed rows, with a maximum of +4 points that can only be achieved by the stick-shaped tetromino.

We trained a shallow neural network agent with dSiLU units in the hidden layer.

11

To handle the more complex learning task, we increased the number of hidden units to 250 and the number of episodes to 400,000. We repeated the experiment for five

215  separate runs. We used the same 20 state features as in the SZ-Tetris experiment, but the length of the binary state vector was reduced to 260 due to the smaller board size. The reward function was changed as follows for the same reason:

$$r(\mathbf{s}) = e^{-(\text{number of holes in } \mathbf{s})/(33/2)}. \tag{15}$$

We used the same values of the meta-parameters as in the stochastic SZ-Tetris experiment.

220  The average learning curve as well as learning curves for the five separate runs in $10 \times 10$ Tetris are shown in Figure 4. The dSiLU network agent reached an average score of 4,900 points over the final 10,000 episodes and the five separate runs, which is a new state-of-the-art in $10 \times 10$ Tetris. The previous best average scores are 4,200 points achieved by the CBMPI algorithm, 3,400 points achieved by the DPI algorithm,

225  and 3,000 points achieved by the CE method (Gabillon et al., 2013). The best individual run achieved a final mean score of 5,300 points, which is also a new state-of-the-art, improving on the score of 5,000 points achieved by the CBMPI algorithm.

It is particularly impressive that the dSiLU network agent achieved its result using features similar to the original Bertsekas features. Using only the Bertsekas features,

230  the CBMPI algorithm, the DPI algorithm, and the CE method could only achieve average scores of about 500 points (Gabillon et al., 2013). The CE method has achieved its best score by combining the Bertsekas features, the Dellacherie features (Fahey, 2003), and three original features (Thiery and Scherrer, 2009). The CBMPI algorithm achieved its best score using the same features as the CE method, except for using five

235  original RBF height features instead of the Bertsekas features.

### 3.3. Atari 2600 games

To further evaluate the use of value-based on-policy reinforcement learning with eligibility traces and softmax action selection in high-dimensional state space domains we applied Sarsa($\lambda$) with a deep convolution neural network function approximator

240  in the Atari 2600 domain using the Arcade Learning Environment (Bellemare et al.,

12

2013). Based on the results for the deep networks in SZ-Tetris, we used a SiLU-dSiLU network with SiLU units in the convolutional layers and dSiLU units in the fully-connected layer. To limit the number of games and prevent a biased selection of the games, we selected the 12 games played by DQN (Mnih et al., 2015) that begin

245 with the letters 'A' and 'B': Alien, Amidar, Assault, Asterix, Asteroids, Atlantis, Bank Heist, Battle Zone, Beam Rider, Bowling, Boxing, and Breakout.

We used a similar experimental setup as Mnih et al. (2015). We pre-processed the raw 210×160 Atari 2600 RGB frames by extracting the luminance channel, taking the maximum pixel values over consecutive frames to prevent flickering, and then

250 downsampling the grayscale images to 105×80. For computational reasons, we used a smaller network architecture. Instead of three convolutional layers, we used two with half the number of filters, each followed by a max-pooling layer. The input to the network was a 105×80×2 image consisting of the current and the fourth previous pre-processed frame. As we used frame skipping where actions were selected every fourth

255 frame and repeated for the next four frames, we only needed to apply pre-processing to every fourth frame. The first convolutional layer had 16 filters of size 8×8 with a stride of 4. The second convolutional layer had 32 filters of size 4×4 with a stride of 2. The max-pooling layers had pooling windows of size 3×3 with a stride of 2. The convolutional layers were followed by a fully-connected hidden layer with 512 dSiLU

260 units and a fully-connected linear output layer with 4 to 18 output (or action-value) units, depending on the number of valid actions in the considered game. We selected meta-parameters by a preliminary search in the Alien, Amidar and Assault games and used the same values for all 12 games: $\alpha$: 0.001, $\gamma$: 0.99, $\lambda$: 0.8, $\tau_0$: 0.5, and $\tau_k$: 0.0005. As in Mnih et al. (2015), we clipped the rewards to be between $-1$ and $+1$,

265 but we did not clip the values of the TD-errors.

In each of the 12 Atari games, we trained a SiLU-dSiLU agent for 200,000 episodes and the experiments were repeated for two separate runs. An episode started with up to 30 'do nothing' actions (*no-op condition*) and it was played until the end of the game or for a maximum of 18,000 frames (i.e., 5 minutes). We evaluated the agents

270 by computing the mean scores over every 100 consecutive episodes in each run, which we call the *mScore*. Figures 5 and 6 show the average learning curves, as well as
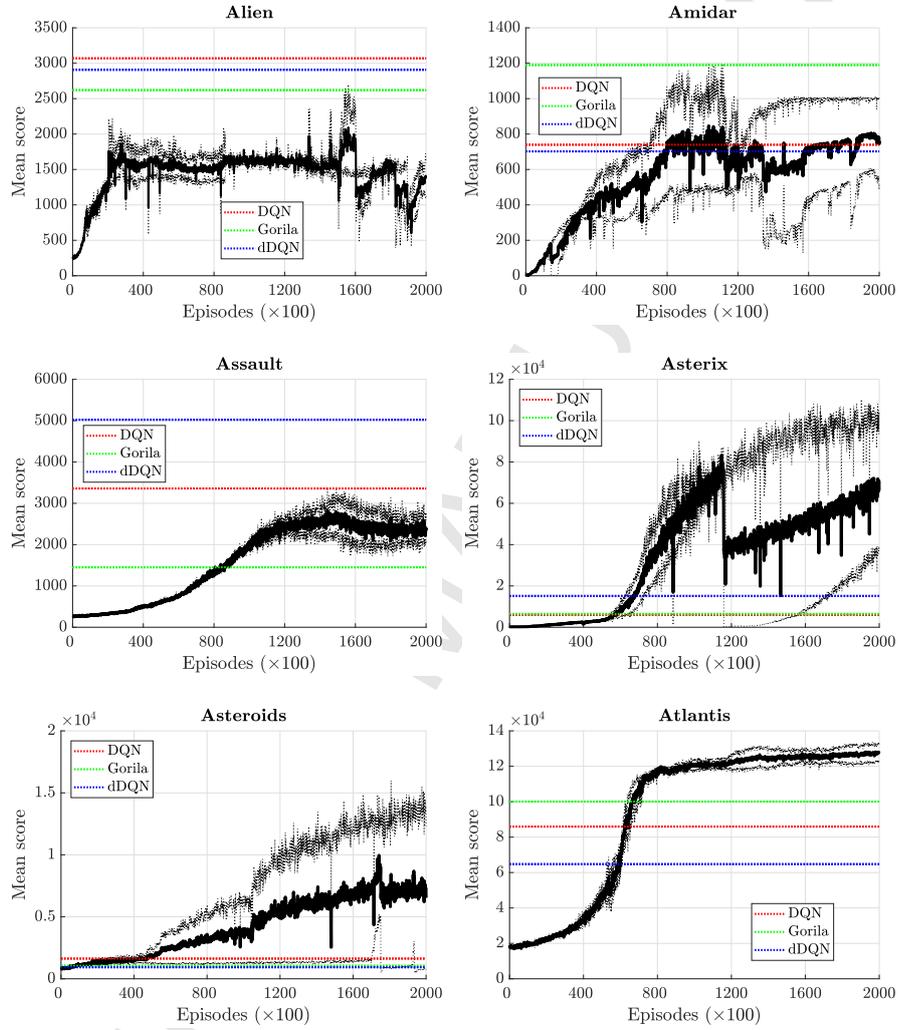
13

Figure 5: Average learning curves (solid lines) over two separate runs (dashed lines) for the SiLU-dSiLU agents in the 6 Atari games that begin with the letter 'A'. The dotted lines show the reported results for DQN (red), the Gorila implementation of DQN (green), and double DQN (blue).

the learning curves for the two separate runs, in the 12 Atari 2600 games. Table 2 summarizes our results using two metrics:

- **Final score**: the average score computed over the final mScores.
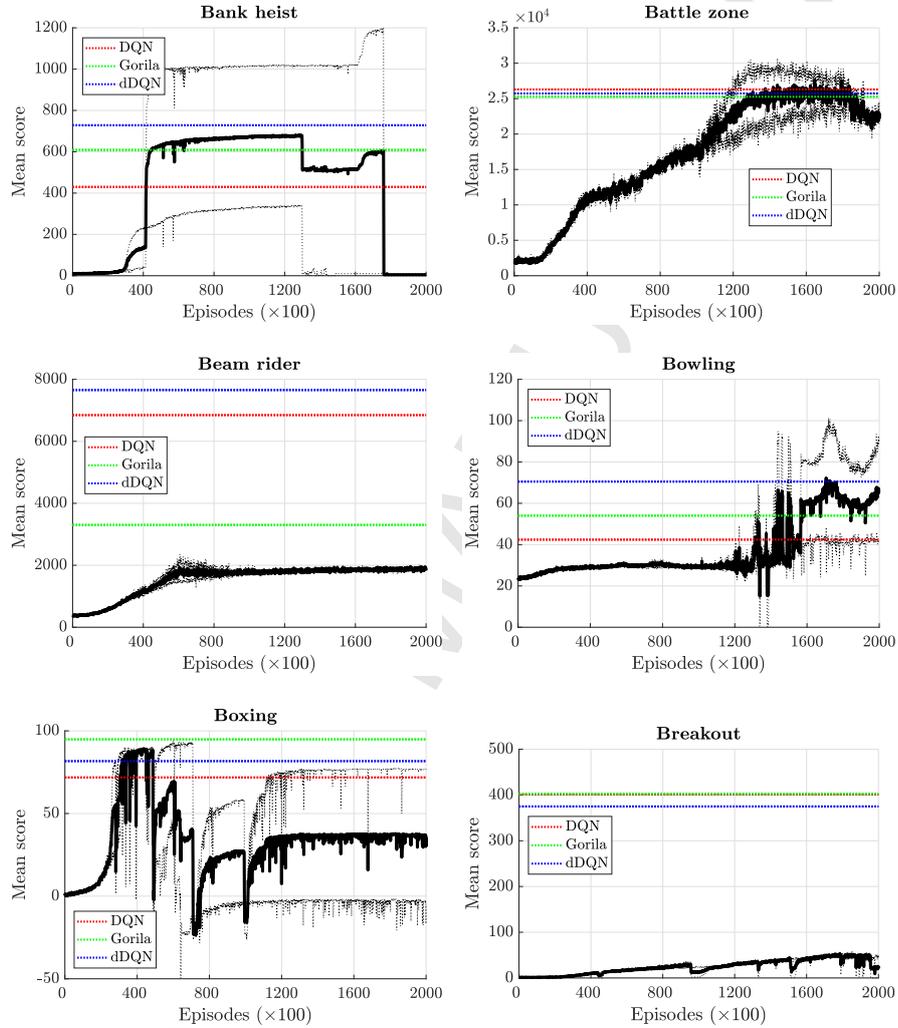
14

Figure 6: Average learning curves (solid lines) over two separate runs (dashed lines) for the SiLU-dSiLU agents in the 6 Atari games that begin with the letter 'B'. The dotted lines show the reported results for DQN (red), the Gorila implementation of DQN (green), and double DQN (blue).

275    • **Best score**: the average score computed over the maximum mScores.

We included the Best score to be able to compare our results with those achieved by DQN (single run scores, where the mScores were computed over 30 episodes; Mnih et al., 2015), the Gorila implementation of DQN (average scores over 5 runs, where the

15

mScores were computed over 30 episodes; Nair et al., 2015), and double DQN (singles
run scores, where the mScores were computed over 100 episodes; van Hasselt et al.,
2015). The last two rows of the table shows summary statistics over the 12 games,
which were obtained by computing the mean and the median of the DQN normalized
scores:

$$\text{Score}_{\text{DQN\_normalized}} = \frac{\text{Score}_{\text{agent}} - \text{Score}_{\text{random}}}{\text{Score}_{\text{DQN}} - \text{Score}_{\text{random}}}$$

Here, $\text{Score}_{\text{random}}$ is the score achieved by a random agent in Mnih et al. (2015).

Table 2: The Best and Final scores achieved by our SiLU-dSiLU agents in 12 Atari 2600 games, and the
reported Best scores achieved by DQN, the Gorila implementation of DQN, and double DQN (dDQN) in the
no-op condition with 5 minutes of evaluation time.

| | DQN | Gorila | dDQN | SiLU-dSiLU | |
| :--- | :---: | :---: | :---: | :---: | :---: |
| **Game** | **Best** | **Best** | **Best** | **Final** | **Best** |
| Alien | **3,069** | 2,621 | 2,907 | 1,370 | 2,246 |
| Amidar | 740 | **1,190** | 702 | 762 | 904 |
| Assault | 3,359 | 1,450 | **5,023** | 2,415 | 2,944 |
| Asterix | 6,012 | 6,433 | 15,150 | 70,942 | **100,322** |
| Asteroids | 1,629 | 1,048 | 931 | 6,537 | **10,614** |
| Atlantis | 85,950 | 100,069 | 64,758 | 127,651 | **128,983** |
| Bank Heist | 430 | 609 | 728 | 5 | **770** |
| Battle Zone | 26,300 | 25,267 | 25,730 | 22,930 | **29,115** |
| Beam Rider | 6,846 | 3,303 | **7,654** | 1,829 | 2,176 |
| Bowling | 42 | 54 | 71 | 67 | **75** |
| Boxing | 72 | **95** | 82 | 36 | 92 |
| Breakout | 401 | **402** | 375 | 25 | 55 |
| **Mean** (DQN Normalized) | 100 % | 102 % | 127 % | 218 % | **332** % |
| **Median** (DQN Normalized) | 100 % | 104 % | 105 % | 78 % | **125** % |

The results clearly show that our SiLU-dSiLU agent outperformed the other agents,
improving the mean (median) DQN normalized Best score score from 127 % (105 %)

achieved by double DQN to $332\%$ ($125\%$). The SiLU-dSiLU agents achieved the highest Best score in 6 out of the 12 games and only performed much worse than the other 3 agents in one game, Breakout, where the learning never took off during the 200,000 episodes of training (see Figure 6). The performance was especially impressive in the Asterix (Best score of 100,322) and Asteroids (Best score of 10,614) games, which improved the Best scores achieved by the second-best agent by $562\%$ and $552\%$, respectively.

## 4. Analysis

### 4.1. Value estimation

First, we investigate the ability of TD($\lambda$) and Sarsa($\lambda$) to accurately estimate discounted returns:

$$R_t = \sum_{k=0}^{T-t} \gamma^k r_{t+k}.$$

Here $T$ is the length of an episode. The reason for doing this is that van Hasselt et al. (2015) showed that the double DQN algorithm improved the performance of DQN in Atari 2600 games by reducing the overestimation of the action values. It is known (Thrun and Schwartz, 1993; van Hasselt, 2010) that Q-learning based algorithms, such as DQN, can overestimate action values due to the max operator, which is used in the computation of the learning targets. TD($\lambda$) and Sarsa($\lambda$) do not use the max operator to compute the learning targets and they should therefore not suffer from this problem.

Figure 7 shows that for SZ-Tetris episodes of average (or expected) length, the best dSiLU network agent at the end of the learning ($\tau = 0.0098$) learned good estimates of the discounted returns, both along the episodes (left panel) and as measured by the normalized sum of differences between $V(s_t)$ and $R_t$ (right panel):

$$\frac{1}{T} \sum_{t=1}^{T} \left( V(s_t) - R_t \right).$$

The linear fit of the normalized sum of differences data for 1,000 episodes gives a small underestimation (-0.43) for an episode of average length (866 time steps). The $V(s_t)$-values overestimated the discounted returns for short episodes and underestimated the
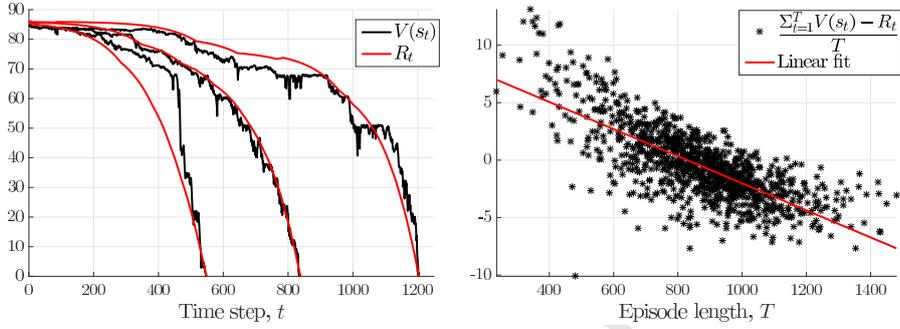
17

Figure 7: The left panel shows learned $V(s_t)$-values and $R_t$-values, for examples of short, medium-long, and long episodes in SZ-Tetris. The right panel shows the normalized sum of differences between $V(s_t)$ and $R_t$ for 1,000 episodes and the best linear fit of the data ($-0.012T + 9.8$).

discounted returns for long episodes (especially in the middle part of the episodes), which is accurate since the episodes ended earlier and later, respectively, than were
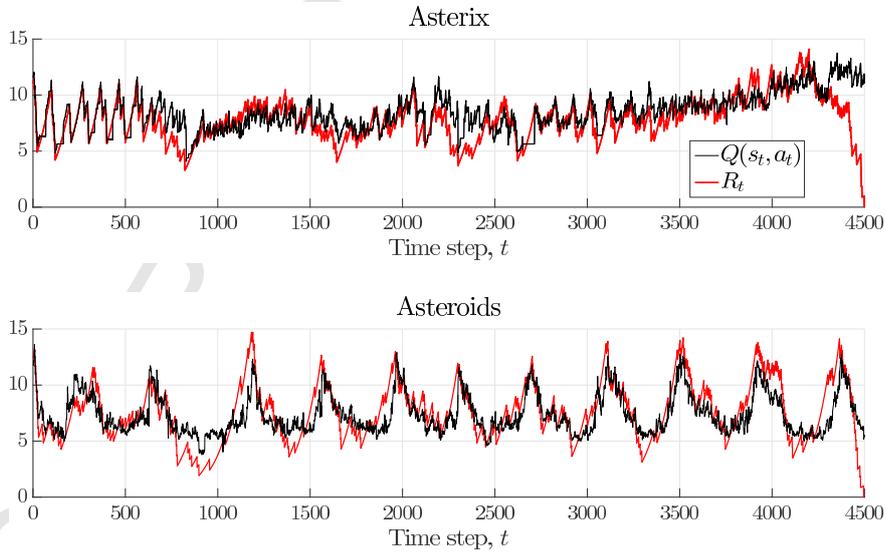315 expected.



Figure 8: Learned action values, $Q(s_t, a_t)$, and discounted returns, $R_t$, for the best SiLU-dSiLU agents in Asterix and Asteroids.

18

Figure 8 shows typical examples of learned action values and discounted returns along episodes in Asterix (score of 108,500) and Asteroids (score of 22,500), when the best SiLU-dSiLU agents at the end of the learning ($\tau = 0.00495$) successfully played for the full 18,000 frames (i.e., 4,500 time steps since the agents acted every

320 fourth frame). In both games, with the exception of a few smaller parts, the learned action values matched the discounted returns very well along the whole episodes. The normalized sums of differences (absolute differences) were 0.59 (1.05) in the Asterix episode and $-0.23$ (1.28) in the Asteroids episode. In both games, the agents overestimated action values at the end of the episodes. However, this is an artifact of that

325 an episode ended after a maximum of 4,500 time steps, which the agents could not predict. Videos of the corresponding learned behaviors in Asterix and Asteroids can be found at `http://www.cns.atr.jp/˜elfwing/videos/asterix_deep_ SiL.mov` and `http://www.cns.atr.jp/˜elfwing/videos/asteroids_ deep_SiL.mov`.

330 *4.2. Action selection*

Second, we investigate the importance of softmax action selection in the games where our proposed agents performed particularly well. Almost all deep reinforcement learning algorithms that have been used in the Atari 2600 domain have used $\varepsilon$-greedy action selection (one exception is the asynchronous advantage actor-critic

335 method, A3C, which used softmax output units for the actor (Mnih et al., 2016)). One drawback of $\varepsilon$-greedy selection is that it selects all actions with equal probability when exploring, which can lead to poor learning outcomes in tasks where the worst actions have very bad consequences. This is clearly the case in both Tetris games and in the Asterix and Asteroids games. In each state in Tetris, many, and often most, actions

340 will creates holes, which are difficult (especially in SZ-Tetris) to remove. In the Asterix game, random exploratory actions can kill Asterix if executed when Cacofonix's deadly lyres are passing. In the Asteroids game, one of the actions sends the spaceship into hyperspace and makes it reappear in a random location, which has the risk of the spaceship self-destructing or of destroying it by appearing on top of an asteroid.

345 We compared softmax action selection, where $\tau$ was set to the annealed values

Table 3: Mean scores and average numbers of non-greedy actions for softmax action selection and $\varepsilon$-greedy action selection.

| Game | Selection | Mean Score | Non-greedy actions |
|------|-----------|------------|--------------------|
| SZ-Tetris | $\tau_{200k} = 0.0098$ | 326 | 28.7 |
| | $\tau_{50k} = 0.0370$ | 318 | 93.0 |
| | $\tau_{10k} = 0.1429$ | 302 | 191.7 |
| | $\varepsilon = 0$ | 332 | 0 |
| | $\varepsilon = 0.001$ | 260 | 0.6 |
| | $\varepsilon = 0.01$ | 71 | 2.0 |
| | $\varepsilon = 0.05$ | 14 | 3.2 |
| Asterix | $\tau_{200k} = 0.00495$ | 104,299 | 47.6 |
| | $\tau_{50k} = 0.0192$ | 97,365 | 178.5 |
| | $\tau_{10k} = 0.0833$ | 71,505 | 576.9 |
| | $\varepsilon = 0$ | 102,890 | 0 |
| | $\varepsilon = 0.001$ | 98,264 | 3.6 |
| | $\varepsilon = 0.01$ | 66,113 | 30.0 |
| | $\varepsilon = 0.05$ | 7,152 | 56.8 |
| Asteroids | $\tau_{200k} = 0.00495$ | 15,833 | 31.3 |
| | $\tau_{50k} = 0.0192$ | 15,421 | 92.7 |
| | $\tau_{10k} = 0.0833$ | 14,219 | 355.2 |
| | $\varepsilon = 0$ | 15,091 | 0 |
| | $\varepsilon = 0.001$ | 11,105 | 2.1 |
| | $\varepsilon = 0.01$ | 3,536 | 11.7 |
| | $\varepsilon = 0.05$ | 1,521 | 47.3 |

after 10,000 episodes ($\tau_{10k}$), 50,000 episodes ($\tau_{50k}$), and 200,000 episodes ($\tau_{200k}$, i.e., the $\tau$-values at the end of learning), and $\varepsilon$-greedy action selection, where $\varepsilon$ was set to 0, 0.001, 0.01, and 0.05, for the best dSiLU network agent in SZ-Tetris and the

best SiLU-dSiLU agents in the Asterix and Asteroids games. The results (see Table 3)

350    clearly show that $\varepsilon$-greedy action selection with $\varepsilon$ set to 0.05, as used for evaluation by DQN, is not suitable for these games. The scores were only 4 % to 10 % of the scores for softmax selection using $\tau_{200k}$. The negative effects of random exploration were largest in Asteroid and SZ-Tetris. Even when $\varepsilon$ was set as low as 0.001 and the agent performed only 2.1 non-greedy actions per episode in Asteroids and 0.6 in SZ-Tetris,

355    the mean scores were reduced by 26 % and 22 %, respectively, compared with $\varepsilon = 0$. In contrast, even if the numbers of non-greedy actions in Asteroids and SZ-Tetris were increased by about an order of magnitude when using $\tau_{10k}$ compared to when using $\tau_{200k}$, the scores were only reduced by 10 % and 7 %, respectively.

## 5. Conclusions

360    In this study, we proposed SiLU and dSiLU as activation functions for neural network function approximation in reinforcement learning. We demonstrated in stochastic SZ-Tetris that SiLUs significantly outperformed ReLUs, and that dSiLUs significantly outperformed sigmoid units. The best agent, the dSiLU network agent, achieved a new state-of-the-art in stochastic SZ-Tetris and in 10×10 Tetris. In the Atari 2600 domain,

365    a deep Sarsa($\lambda$) agent with SiLUs in the convolutional layers and dSiLUs in the fully-connected hidden layer outperformed DQN and double DQN, as measured by mean and median DQN normalized scores.

An additional purpose of this study was to demonstrate that a more traditional approach of using on-policy learning with eligibility traces and softmax selection (i.e.,

370    basically a "textbook" version of a reinforcement learning agent but with non-linear neural network function approximators) can be competitive with the approach used by DQN. This means that there is a lot of room for improvements, by, e.g., using, as DQN, a separate target network, but also by using more recent advances such as the dueling architecture (Wang et al., 2016) for more accurate estimates of the action values and

375    asynchronous learning by multiple agents in parallel (Mnih et al., 2016).

## 6. Acknowledgments

380 **References**

Bellemare, M. G., Naddaf, Y., Veness, J., Bowling, M., 2013. The arcade learning environment: An evaluation platform for general agents. Journal of Artificial Intelligence Research 47, 253–279.

Bertsekas, D. P., Ioffe, S., 1996. Temporal differences based policy iteration and applications in neuro-dynamic programming. Technical Report LIDS-P-2349, MIT.

Burgiel, H., 1997. How to lose at Tetris. Mathematical Gazette 81, 194–200.

Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., Fei-Fei, L., 2009. ImageNet: A Large-Scale Hierarchical Image Database. In: CVPR09.

Elfwing, S., Uchibe, E., Doya, K., 2015. Expected energy-based restricted boltzmann machine for classification. Neural Networks 64 (3), 29–38.

Elfwing, S., Uchibe, E., Doya, K., 2016. From free energy to expected energy: Improving energy-based value function approximation in reinforcement learning. Neural Networks 84, 17–27.

Elfwing, S., Uchibe, E., Doya, K., 2017. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning. arXiv:1702.03118 [cs.LG].

Fahey, C., 2003. Tetris AI, computer plays tetris. `colinfahey.com/tetris/tetris.html` [Online; accessed 22-February-2017].

Faußer, S., Schwenker, F., 2013. Neural network ensembles in reinforcement learning. Neural Processing Letters, 1–15.

22

Freund, Y., Haussler, D., 1992. Unsupervised learning of distributions on binary vectors using two layer networks. In: Moody, J. E., Hanson, S. J., Lippmann, R. P. (Eds.), Proceedings of Advances in Neural Information Processing Systems (NIPS1992). Morgan Kaufmann.

Gabillon, V., Ghavamzadeh, M., Scherrer, B., 2013. Approximate dynamic programming finally performs well in the game of tetris. In: Proceedings of Advances in Neural Information Processing Systems (NIPS2013). pp. 1754–1762.

Hahnloser, R. H. R., Sarpeshka, R., Mahowald, M. A., Douglas, R. J., Seung, H. S., 2000. Digital selection and analogue amplification coexist in a cortex-inspired silicon circuit. Nature 405, 947–951.

Hinton, G. E., 2002. Training products of experts by minimizing contrastive divergence. Neural Computation 12 (8), 1771–1800.

Jaskowski, W., Szubert, M. G., Liskowski, P., Krawiec, K., 2015. High-dimensional function approximation for knowledge-free reinforcement learning: a case study in SZ-Tetris. In: Proceedings of the Genetic and Evolutionary Computation Conference (GECCO2015). pp. 567–573.

Krizhevsky, A., 2009. Learning multiple layers of features from tiny images. Tech. rep., University of Toronto.

Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T. P., Harley, T., Silver, D., Kavukcuoglu, K., 2016. Asynchronous methods for deep reinforcement learning. In: Proceedings of the International Conference on Machine Learning (ICML2016). pp. 1928–1937.

Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., Hassabis, D., 2015. Human-level control through deep reinforcement learning. Nature 518 (7540), 529–533.

Nair, A., Srinivasan, P., Blackwell, S., Alcicek, C., Fearon, R., Maria, A. D., Panneershelvam, V., Suleyman, M., Beattie, C., Petersen, S., Legg, S., Mnih, V., Kavukcuoglu, K., Silver, D., 2015. Massively parallel methods for deep reinforcement learning. arXiv:1507.04296 [cs.LG].

Ramachandran, P., Zoph, B., Le, Q. V., 2017. Searching for activation functions. arXiv:1710.05941 [cs.NE].

Rummery, G. A., Niranjan, M., 1994. On-line Q-learning using connectionist systems. Tech. Rep. CUED/F-INFENG/TR 166, Cambridge University Engineering Department.

Schaul, T., Quan, J., Antonoglou, I., Silver, D., 2016. Prioritized experience replay. In: International Conference on Learning Representations (ICLR2016). Puerto Rico.

Scherrer, B., Ghavamzadeh, M., Gabillon, V., Lesner, B., Geist, M., 2015. Approximate modified policy iteration and its application to the game of tetris. Journal of Machine Learning Research 16, 1629–1676.

Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., van den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., Dieleman, S., Grewe, D., Nham, J., Kalchbrenner, N., Sutskever, I., Lillicrap, T., Leach, M., Kavukcuoglu, K., Graepel, T., Hassabis, D., 2016. Mastering the game of go with deep neural networks and tree search. Nature 529, 484–503.

Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., Chen, Y., Lillicrap, T., Hui, F., Sifre, L., van den Driessche, G., Graepel, T., Hassabis, D., 10 2017. Mastering the game of go without human knowledge 550, 354–359.

Smolensky, P., 1986. Information processing in dynamical systems: Foundations of harmony theory. In: Rumelhart, D. E., McClelland, J. L. (Eds.), Parallel Distributed Processing: Explorations in the Microstructure of Cognition. Volume 1: Foundations. MIT Press.

Sutton, R. S., 1988. Learning to predict by the method of temporal differences. Machine Learning 3, 9–44.

Sutton, R. S., 1996. Generalization in reinforcement learning: Successful examples using sparse coarse coding. In: Proceedings of Advances in Neural Information Processing Systems (NIPS1996). MIT Press, pp. 1038–1044.

Sutton, R. S., Barto, A., 1998. Reinforcement Learning: An Introduction. MIT Press.

Szita, I., Szepesvári, C., 2010. SZ-Tetris as a benchmark for studying key problems of reinforcement learning. In: ICML 2010 workshop on machine learning and games.

Tesauro, G., 1994. Td-gammon, a self-teaching backgammon program, achieves master-level play. Neural Computation 6 (2), 215–219.

Thiery, C., Scherrer, B., 2009. Improvements on learning tetris with cross entropy. International Computer Games Association Journal 32.

Thrun, S., Schwartz, A., 1993. Issues in using function approximation for reinforcement learning. In: Proceedings of the 1993 Connectionist Models Summer School. pp. 255–263.

van Hasselt, H., 2010. Double q-learning. In: Proceedings of Advances in Neural Information Processing Systems (NIPS2010). pp. 2613–2621.

van Hasselt, H., Guez, A., Silver, D., 2015. Deep reinforcement learning with double q-learning. arXiv:1509.06461 [cs.LG].

Wang, Z., Schaul, T., Hessel, M., van Hasselt, H., Lanctot, M., de Freitas, N., 2016. Dueling network architectures for deep reinforcement learning. In: Proceedings of the International Conference on Machine Learning (ICML2016). pp. 1995–2003.