

Physical Reinforcement Learning

Anonymous authors

Paper under double-blind review

Abstract

1 Digital computers are power-hungry and largely intolerant of damaged components,
2 making them potentially difficult tools for energy-limited autonomous agents in uncer-
3 tain environments. Recently developed *Contrastive Local Learning Networks* (CLLNs),
4 analog networks of self-adjusting nonlinear resistors, are inherently low-power and ro-
5 bust to physical damage, but have been exclusively used to perform supervised learn-
6 ing tasks (Dillavou et al., 2024). In this work, we demonstrate success on two simple
7 Markov decision processes using tabular Q-learning adapted for simulated CLLNs. Do-
8 ing so makes explicit the components (beyond the network being trained) required to
9 enact various tools in the RL toolbox, some of which (policy function and value func-
10 tion) are more natural in this system than others (replay buffer). Limitations to external
11 memory make the big world hypothesis, and learning via tracking and streaming, an at-
12 tractive framework for reinforcement learning on CLLN-based agents (Javed & Sutton,
13 2024; Elsayed et al., 2024). We discuss assumptions such as the physical safety that
14 digital hardware requires, CLLNs can forgo, and biological systems cannot rely on, and
15 highlight secondary goals that are important in biology and trainable in CLLNs, but
16 make little sense in digital computers. We share our code for Q-learning adapted for
17 simulated CLLNs, for further consideration as a physical framework for learning in the
18 big world [anonymized link].

19 1 Introduction

20 Digital computers are powerful and versatile machines, but have significant weaknesses, especially
21 in areas relevant to reinforcement learning (RL). First, they are fault intolerant and thus susceptible
22 to damage. Disable a handful or even a single transistor in a CPU or GPU and it may crash the
23 entire system¹. Having a single transistor erroneously flip even on 1% of write operations could
24 catastrophically disable the entire machine. This is because the function of each component is
25 intrinsically tied to its location within the system, and the entire operation relies on (nearly) error-
26 free operation. Second, they are power-hungry, with laptop CPUs consuming approximately 50 W.
27 This stems from the high energy cost of maintaining ‘perfect’ operation, as well as the shuttling of
28 data between processing and memory.

29 Low-power operability and fault tolerance are areas where biological systems thrive, as they are
30 autonomous agents in an often energy-scarce and dangerous world. Brains can take significant dam-
31 age and continue to function, including destruction of single neurons (Wang et al., 2014), traumatic
32 brain injuries (Chua et al., 2007), and even removal of large brain regions (Granovetter et al., 2022).
33 This robustness stems from brains’ distributed processing and emergent, rather than linear, compu-
34 tation. The human brain is also energy efficient, using 20 W in total (Balasubramanian, 2021) while
35 performing perception, cognition, motor control, homeostatic functions, learning, and more. This
36 energy efficiency is in part due to the overlap of memory and computation, and to *natural primi-*
37 *tives* (Mead, 1990), whereby low-level computations utilize low-power analog physical or chemical
38 processes rather than digital logic.

¹<https://www.ntchip.com/electronics-news/transistors-in-cpu>

39 In spite of the potential advantages, there are few instances of artificial and non-digital hardware
 40 performing RL tasks (Mak et al., 2007; 2010). Many instances of digitally-enhanced or simulated
 41 analog systems have been used in an RL setting (as in Mikaitis et al. (2018)), but few if any hard-
 42 ware demonstrations combine the distributed memory and computation with analog signals that give
 43 brains some significant survivability advantages.

44 Recently developed Contrastive Local Learning Networks (CLLNs) (Dillavou et al., 2024) may fill
 45 this gap. A network of self-adjusting nonlinear resistors, these systems perform supervised learning
 46 in an emergent manner by utilizing physical processes to perform computation. Learning occurs in
 47 each individual element using local rules, eschewing central and digital processing. As a result, they
 48 are energy efficient (Stern et al., 2024) and fault tolerant (Dillavou et al., 2022). However, CLLNs
 49 have not been demonstrated in an RL context, even in simulation.

50 Designing an RL implementation of CLLNs provides a setting to ask a range of questions. For
 51 instance, which RL tools are natural for a self-learning network, and which require additional pre-
 52 programmed hardware? There are numerous challenges associated with digital RL acting in the real
 53 (non-simulated) world (Dulac-Arnold et al., 2021); what additional challenges stem from placing
 54 the agent’s brain outside of the pristine digital realm (Fig. 1)? A successful navigation of these
 55 challenges would open a path towards energy-efficient and fault-tolerant autonomous agents.

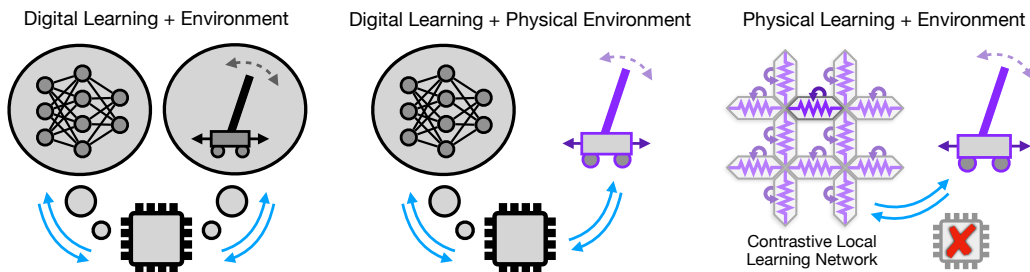


Figure 1: A schematic of reinforcement learning in three scenarios. Left: A digital agent learns in a digital environment, both simulated by a computer. Middle: A digital agent interacts with a physical environment, and the learning process is controlled by a computer. Right (proposed): A Contrastive Local Learning Network interacts with a physical environment, and learns based on those interactions. There is no digital processing; learning is done in a distributed, analog fashion. One of the many identical self-adjusting components is highlighted. Physical components are highlighted in purple, and simulated processes and digital components are shown in grey.

56 In the following text, we first present the dynamics of CLLNs. Following this, we highlight aspects
 57 of CLLNs that make the big world hypothesis an attractive framework for algorithm development in
 58 this setting (Javed & Sutton, 2024). We then illustrate the adaptation of Q-Learning to the dynamics
 59 of a CLLN in the context of two simple Markov decision processes (MDP), and discuss the broader
 60 implications of using CLLNs as a substrate for reinforcement learning methods.

61 2 Reinforcement learning with a simulated CLLN

62 2.1 Background: Contrastive Local Learning Networks

63 Contrastive Local Learning Networks (CLLNs) are networks of self-adjusting resistive elements
 64 whose individual dynamics approximate gradient descent on a global loss function (Dillavou et al.,
 65 2024). The standard operation of these systems is for supervised learning and is briefly summarized
 66 here.

67 We consider a network of resistors as shown schematically in Fig. 2. We will treat four of the nodes
 68 (yellow) as inputs, and encode data by enforcing the voltage values of these nodes ($U_1^{\text{in}}, \dots, U_4^{\text{in}}$).
 69 We will also choose four nodes (blue), and treat their equilibrium voltage values, $U_1^{\text{out}}, \dots, U_4^{\text{out}}$,

70 as outputs. That is, we treat our network as a physical function, $(U_1^{\text{out}}, U_2^{\text{out}}, U_3^{\text{out}}, U_4^{\text{out}}) \equiv$
 71 $\mathcal{F}(U_1^{\text{in}}, U_2^{\text{in}}, U_3^{\text{in}}, U_4^{\text{in}})$. It is a consequence of our eventual choice of task that the number of in-
 72 puts and outputs is equal; they are independently specified in general. When input voltages are
 73 applied, output values are ‘calculated’ by the physical processes, in this case the flow of current in
 74 the network. The result of this calculation is a consequence of the network structure and the con-
 75 ductance (inverse of resistance) values of each edge of the network, G_i . We note that at equilibrium,
 76 the voltages at unconstrained nodes have minimized power dissipation with respect to the boundary
 77 conditions².

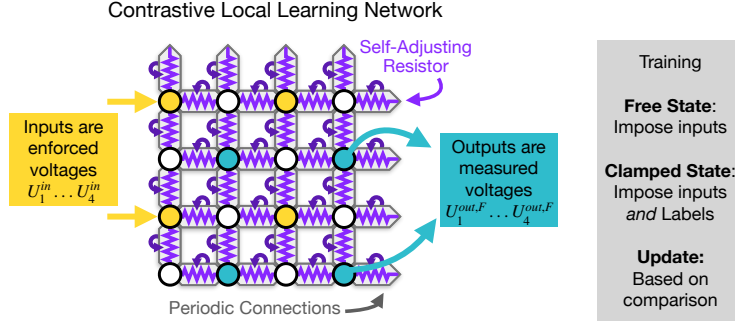


Figure 2: Schematic of a Contrastive Local Learning Network (CLLN). The configuration shown is used for the Markov decision process with four states and four actions. A modified network architecture and input-output position is used for the 9-state navigation task (see Fig. 4). A high-level description of the contrastive training protocol is outlined in the grey box. Note that each element in a CLLN requires only local measurements of the two states (free and clamped) to update itself, decentralizing the training process.

78 The evolution of the gate voltages (the trainable parameters) in this network is *contrastive* in that
 79 it requires comparing two distinct states that we can impose. Given a single data point (inputs
 80 $U_1^{\text{in}}, \dots, U_4^{\text{in}}$) and label, (desired outputs L_1, \dots, L_4), the two states are as follows. First, the *free*
 81 state, wherein only inputs $U_1^{\text{in}}, \dots, U_4^{\text{in}}$ are imposed. In this state, each resistor experiences voltage
 82 drop ΔU_i^F and the outputs find values $U_n^{\text{out},F}$. We use subscript i to denote quantities on network
 83 edges, and n to denote quantities on nodes (inputs or outputs). Next, the *clamped* state, in which
 84 inputs as well as the desired outputs (labels) are enforced. Now each resistor experiences voltage
 85 drop ΔU_i^C , and the (enforced) outputs are equal to the labels $U_1^{\text{out},C} = L_1$, etc. In practice the
 86 clamped output is actually nudged *towards* the label: $U_n^{\text{out},C} = U_n^{\text{out},F}(1 - \eta) + \eta L_n$, where η is a
 87 hyperparameter, and $\eta = 1$ simply clamps the label ($U_n^{\text{out},C} = L_n$). In this work we use $\eta = 0.1$.

88 The system will perform gradient descent on a *contrastive* function, with respect to its conductance
 89 values, following the Coupled Learning framework (Stern et al., 2021). This framework is closely
 90 related to Equilibrium Propagation (Scellier & Bengio, 2017) and Contrastive Hebbian Learning
 91 (Movellan, 1991). Specifically,

$$\delta G_i = -\alpha \frac{d}{dG_i} [\mathcal{P}^C - \mathcal{P}^F], \quad (1)$$

92 where \mathcal{P}^C and \mathcal{P}^F are the total dissipated power in the clamped and free states respectively, and
 93 α is a learning rate. We note that because the electronic network minimizes power with respect to
 94 its boundary conditions, and the clamped state is equivalent to the free state with added boundary
 95 conditions (the clamped outputs), the contrastive function must be non-negative. Further, it is only
 96 zero-valued when the clamping procedure does not change the outputs, that is, the labels are already

²In a nonlinear network like the one we consider, the true minimized quantity is *co-content*, which reduces to power in the linear case. As this is a less-commonly known quantity we use power throughout. This does not change the intuition or derivations.

97 satisfied. Thus, it has the same minima as a mean-squared error cost function, and in practice its
 98 minimization works well to minimize error (Stern et al., 2021; Dillavou et al., 2024; 2022).

99 Calculating the derivative on the RHS of eq. 1 for the clamped power with the chain rule, we find
 100 that

$$\frac{d\mathcal{P}^C}{dG_i} = \sum_j^{\text{edges}} \frac{\partial \mathcal{P}^C}{\partial G_j} \frac{dG_j}{dG_i} + \frac{\partial \mathcal{P}^C}{\partial \Delta U_j^C} \frac{d\Delta U_j^C}{dG_i} = \frac{\partial \mathcal{P}^C}{\partial G_i}. \quad (2)$$

101 The first term cancels for $i \neq j$ as the conductances are independently set, and the second cancels
 102 as power is (as noted) at a minimum with respect to the voltages. As a result, the total derivative
 103 becomes a partial derivative, which we can further simplify by noting that power is a sum over local
 104 powers:

$$\mathcal{P}^C = \sum_j^{\text{edges}} (\Delta U_j^C)^2 G_j \rightarrow \frac{d\mathcal{P}^C}{dG_i} = \frac{\partial \mathcal{P}^C}{\partial G_i} = (\Delta U_i^C)^2 \quad (3)$$

105 Thus we complete the same calculation for the free power, and the result is a local learning rule,

$$\delta G_i = -\alpha \frac{d}{dG_i} [\mathcal{P}^C - \mathcal{P}^F]_{\text{contrastive fn}} = \alpha [(\Delta U_i^F)^2 - (\Delta U_i^C)^2], \quad (4)$$

106 wherein each element needs only to measure its own voltage drop in both states to enact this gradient
 107 descent. In practice, this is done by creating twin, co-evolving networks, one holding the free state,
 108 the other the clamped state (Dillavou et al., 2022; 2024).

109 We model our simulations after Dillavou et al. (2024), where each conductive element is in fact a
 110 MOSFET transistor in the triode (passive) regime, with conductance equation

$$G_i = S(U_{G,i} - U_T - \bar{U}_i) \quad (5)$$

111 where in our simulations $S = 1$ and $U_T = 0.7$ are constants, $U_{G,i}$ is the (adjustable) gate voltage
 112 that will act as weights in our system, and \bar{U} is the average of the node voltages on either side of
 113 the edge. This final term allows the system to perform nonlinear input to output transformations.
 114 We note that because $\partial G_i / \partial U_{G,i} = S$, the gradient with respect to G_i and $U_{G,i}$ is co-linear, and we
 115 may enact eq 4 by changing $U_{G,i}$, which will be our learning procedure. We restrict the range of our
 116 parameters such that $1.0 < U_{G,i} < 5.5$ to mimic experimental conditions.

117 2.2 Connections to the big world hypothesis

118 Several properties of CLLNs make the big world hypothesis an attractive framework for learning
 119 directly with this physical substrate. Given that parameters and outputs in the CLLN are bounded,
 120 the learning response to a signal will naturally be limited in the larger external environment. Perhaps
 121 most pertinently, common deep RL approaches incorporate the storage of history, through eligibility
 122 traces and replay buffers (Haarnoja et al., 2018; Schulman et al., 2017; Mnih et al., 2013). While
 123 function approximation is naturally part of a CLLN, it is not obvious how to naturally introduce
 124 history storage into a CLLN, and methods like tracking and streaming deep RL (Elsayed et al.,
 125 2024) are attractive as they better support the vision of ‘a network in the wild’ (Fig. 1). We share
 126 our code, for consideration as a simulated physical substrate for reinforcement learning in big world
 127 tasks [anonymized link].

128 2.3 Example 1: A toy Markov decision process with four states and four actions

129 In this section, we adapt the update rule for the CLLN to Q-Learning on a Markov decision pro-
 130 cess (MDP) with four states and four actions. We will treat our CLLN simulation as enacting
 131 the Q matrix. At each training step t , we encode the environmental state S_t as the input; states
 132 $S_1, \dots, S_4 \equiv [U_1^{\text{in}}, U_2^{\text{in}}, U_3^{\text{in}}, U_4^{\text{in}}] = [1 \ 0 \ 1 \ 0], [0 \ 1 \ 0 \ 1], [1 \ 1 \ 0 \ 0],$ and $[0 \ 0 \ 1 \ 1]$ V respectively.

133 Note that all voltages are analog, despite the input appearing to be digital, and the outputs will fall
 134 between 0 and 1. We then use ϵ -greedy action selection where ϵ decays from 0.05 to 0 linearly over
 135 the training run. We select the maximum output of the four (with probability $1 - \epsilon$) as our action A_t .

136 The environment then issues a state-dependent reward $R(S_t, A_t) + \mathcal{N}(0, 0.01)$ and the state changes
 137 to S_{t+1} . $R(S_t, A_t)$ was sampled from $\mathcal{N}(0.1, 0.1)$, and the selected values are shown in Fig. 3A.
 138 Our environment is extremely simple, with action 0 inducing state 0, action 1 inducing state 1 and
 139 so on. We then impose the inputs for state S_{t+1} , and take the maximum output minus the mean of
 140 all four outputs as our predicted future reward. We use a discount factor $\gamma = 0.5$ and all of this
 141 information to create future-weighted score L_t :

$$L_t = R(S_t, A_t) + \gamma [\max(\mathcal{F}(S_{t+1})) - \text{mean}(\mathcal{F}(S_{t+1}))] \quad (6)$$

142 We find that subtracting the mean term improves performance, as it removes a constraint on the
 143 average set of outputs, permitting our small network additional flexibility. We allow the system to
 144 evolve (eq 4) while imposing encoded S_t as inputs, U_i^{out} on all clamped-state outputs corresponding
 145 to actions the system *did not* take, and L_t on the clamped-state output whose action was taken. Our
 146 approach is summarized in Algorithm 1.

147 As an example: the system is in state $S_t = 1$, and U_2^{out} is the biggest (we assume we do not take the
 148 rare random action). The agent takes action $A_t = 2$ as a result, gets reward R_t , and the state moves
 149 to $S_{t+1} = A_t = 2$. L_t is calculated as above, and then we allow our system to evolve (eq 4) while
 150 S_t is imposed as input in both free and clamped states, and $\mathcal{F}(S_{t+1})$ is imposed as the clamped label
 151 with L_t inserted at position 2 (the action taken). That is, we train the system to keep the outputs
 152 of the untouched actions the same, while moving the chosen action’s output value towards L_t . The
 153 result is evolution similar to the Bellman equation. Updates are batched and imposed on the system
 154 every 50 steps.

155 2.3.1 Results

156 We perform 10 trials of 100,000 training steps, each time randomly initializing $U_{G,i}$ using
 157 $\mathcal{N}(1.5, 0.1)$. Rewards are tallied over time, and averaged in logarithmically spaced intervals (min-
 158 imum of 10 steps) for each trial, shown as purple lines in Fig. 3B, with their average overlaid in
 159 black. We simulate rewards for our same protocol (with $\epsilon = 0$ but including reward noise) using
 160 10,000 steps with each of the $4^4 = 256$ possible strategies, and denote the best and worst cases with
 161 black dashed lines in Fig. 3B. Our system always learns a near-optimal strategy, including ending at
 162 the best possible strategy in 8 of 10 trials.

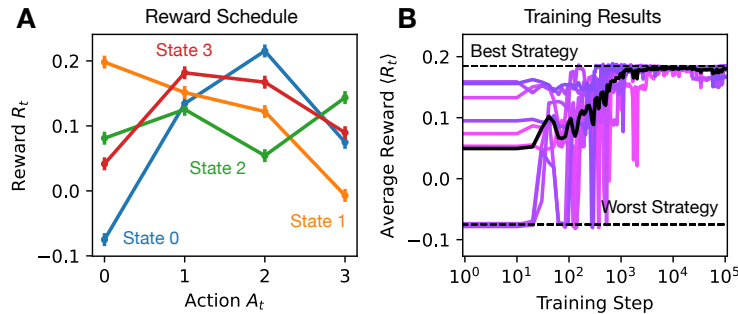


Figure 3: Q Learning with CLLN. (A) Reward schedule for each state (noise effect shown as small error bars). The optimal strategy involves a cycle through all four states. (B) Average reward over training for 10 trials (purples) overlaid with their average (black).

163 **2.4 Example 2: Navigation on a grid with nine states and four actions**

164 We use a modified network architecture to attempt a slightly more difficult task, navigating towards
 165 a goal location in discretized 2D space. The reward schedule is shown in Fig. 4A, with a reward
 166 gradient (shaping) 5000 times smaller than the large reward (top left), invisible on the heatmap. The
 167 optimal strategy (overlaid) is to always move up or left.

168 We now simulate a larger, 44-edge network to perform this task. The architecture is equivalent to
 169 a densely connected neural network, with 6, 4, 4, and 1 node in each layer. Note however that our
 170 network is *not* feed forward, and thus we place our four output nodes in the third layer, and have
 171 each one represent a move (up, down, left, right). The row and state values are rescaled [0, 0.5, 1]
 172 and imposed as two inputs in the first layer. Their inverses [1, 0.5, 0] are also imposed in the first
 173 layer, and the final two nodes in this layer are always set to 0 and 1, respectively. The single node
 174 in the final layer is set to 0.5, making 7 total inputs, only 4 of which change (with the state). When
 175 training (and testing) the system is randomly placed in a new state every 5 steps.

176 **2.4.1 Results**

177 As before, we randomly initialize our network ($U_{G,i}$ from $\mathcal{N}(1.5, 0.1)$) and batch updates every 50
 178 training steps. We train 10 trials for 300,000 steps, reduce ϵ from 0.1 to 0 linearly in each trial,
 179 and report the average rewards for each section of training in purple in Fig. 4B, with their average
 180 overlaid in black. In 8 of our 10 trials, our system finds one of the optimal strategies (several are
 181 equivalent). Note that the system cannot sit on the 'best strategy' line while $\epsilon > 0$, which is the case
 182 until the very end of training.

183 After training, we simulate each trained agent operating in the environment for 10,000 steps, still
 184 randomly setting their location every 5 steps, but with no random actions beyond this ($\epsilon = 0$). We
 185 find the agents spend most of their time in the high-reward square, as shown in Fig. 4C.

186 **3 Discussion**

187 Our experiments illustrate that the reinforcement learning (RL) scheme developed for a physical net-
 188 work performs well on an MDP comprising four states and four actions, and on a navigational task
 189 in a two-dimensional grid, with nine states and four actions. Trained agents find optimal strategies
 190 with few exceptions. While the above algorithm is quite simple, its development required engaging
 191 with two atypical features when compared to digital RL:

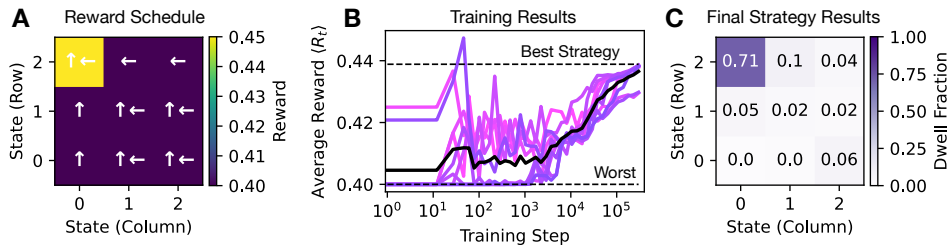


Figure 4: Navigation Task. (A) Reward schedule, with all states approximately equal except for the target (upper left) state. There is no reward noise. The optimal action in each grid is indicated via arrows. (B) Average reward over training for 10 trials (purples) is overlaid with their average (black). The single purple spike above the best average strategy is a result of a few lucky runs, and that binning is narrower early in training. (C) The 10 strategies at the end of training are each simulated for 10,000 steps (reset randomly every 5 steps), and the fraction of total time spent in each state is shown as a heatmap.

- 192 • **Parameters and outputs in our system are bounded.** As a result, we must take care not to define
193 rewards or discount factors such that the system cannot produce the required outputs.
- 194 • **CLLNs are not ‘feed forward’.** A voltage anywhere in the network (or a change in resistance
195 anywhere) potentially impacts every output. Thus, our protocol required us to train unselected
196 outputs to remain stationary. The limitations this poses for the complexity of agent strategies is
197 a subject for future study, as are the advantages of a network capable of sensing inputs from any
198 location.

199 Further, envisioning a physical system performing our scheme necessitates a somewhat ‘unnatural’
200 step backwards in time:

- 201 1. System evaluates state S_t and selects action A_t
- 202 2. Environment returns reward R_t , and switches to S_{t+1}
- 203 3. System evaluates state S_{t+1} and calculates L_t
- 204 4. System updates based on **previous** state S_t and L_t

205 This requires storing and reapplying S_t after the environment has moved on (S_{t+1}). It is possible
206 that some constructed hybrid state (involving S_{t+1} and L) could allow the system to avoid such
207 ‘backtracking’ and still train successfully.

208 Even our current simple implementation required some external control, to implement the random-
209 ization and the max function in the epsilon-greedy algorithm, and to hold and swap back S_t as
210 discussed above. While a future CLLN could certainly be trained to perform these functions, the
211 system is then being asked to learn the learning algorithm before succeeding at the task. Other
212 schemes in the RL literature are far more natural to include. For instance, implementing both a
213 policy function and value function could be done with a single network, or by dividing the network
214 in two. Our adaptation of tabular Q-learning to a CLLN substrate is illustrated in an open-source
215 codebase.

216 By virtue of being a simulation, our CLLN has several features of the digital world. However when
217 physically constructed, these aspects change. First, in hardware, imperfections in components and
218 measurements can manifest as biased elements, hamstringing the learning process and preventing
219 capture of subtle learning signals (Dillavou et al., 2025). Such effects are broadly avoided (for good
220 reason) in RL and machine learning, but spotlight a distinct advantage that these algorithms have
221 over physically instantiated learning systems like the brain.

222 Second, the notion of ‘damage’ is somewhat meaningless in simulation, but has very real conse-
223 quences for biological systems that rely on a physical substrate. As CLLNs are collections of self-
224 adjusting elements and can retrain after damage (Dillavou et al., 2022), developing architectures and
225 algorithms to allow on-the-fly recovery in RL tasks is an exciting research direction, and one that
226 makes little sense in the digital domain.

227 Third, the energy used by the processor in simulating our CLLN is independent of the resistances
228 and currents *in silico*. However a physically constructed version burns power directly proportional
229 to these values, and small modifications to the learning have been demonstrated to bias physically
230 instantiated CLLNs towards low power solutions (Stern et al., 2024). Perhaps networks could like-
231 wise be biased into states more robust to physical damage, or with faster input-output transmission.
232 Modifications to learning algorithms that account for such secondary goals are an exciting avenue
233 for future work, and as above, make little sense in digital systems.

234 4 Conclusion

235 This work highlights a new type of analog, distributed system, a Contrastive Local Learning Net-
236 work (CLLN), whose features make it a natural fit for reinforcement learning (RL), particularly in
237 the big world framework. We simulated two CLLNs, and developed a training scheme to leverage
238 their self-learning dynamics to enact Q Learning on one Markov decision process each, the first

239 with four states and four actions, the second with a nine states (3×3 grid) and four actions (cardinal
240 direction moves). The CLLN agents were largely successful, achieving near optimal rewards across
241 trials. Our scheme highlighted several aspects of typical RL algorithms that are less natural to apply
242 to physical self-learning networks, most including some form of additional nontrivial memory. The
243 limitations associated with external memory make the big world an attractive framework for devel-
244 oping RL methods entirely outside the pristine digital realm, towards advantages like low-power
245 operability and fault tolerance. To support further consideration of CLLNs as a physical substrate
246 for RL, we shared the codebase demonstrating tabular RL on a simulated CLLN [[anonymized link](#)].
247 Finally, we discussed disadvantages like imperfect components and secondary goals such as low-
248 power operation and robustness to damage, all of which may be improvable via modified training
249 methods in CLLNs, a process that makes little sense in the digital domain.

250 References

- 251 Vijay Balasubramanian. Brain power. *Proceedings of the National Academy of Sciences*, 118(32):
252 e2107022118, August 2021. DOI: 10.1073/pnas.2107022118.
- 253 Karen Sg Chua, Yee-Sien Ng, Samantha Gm Yap, and Chek-Wai Bok. A Brief Review of Traumatic
254 Brain Injury Rehabilitation. *Annals of the Academy of Medicine, Singapore*, 36(1):31–42, January
255 2007. ISSN 0304-4602. DOI: 10.47102/annals-acadmedsg.V36N1p31.
- 256 Sam Dillavou, Menachem Stern, Andrea J. Liu, and Douglas J. Durian. Demonstration of De-
257 centralized Physics-Driven Learning. *Physical Review Applied*, 18(1):014040, July 2022. DOI:
258 10.1103/PhysRevApplied.18.014040.
- 259 Sam Dillavou, Benjamin D. Beyer, Menachem Stern, Andrea J. Liu, Marc Z. Miskin, and Douglas J.
260 Durian. Machine learning without a processor: Emergent learning in a nonlinear analog network.
261 *Proceedings of the National Academy of Sciences*, 121(28):e2319718121, July 2024. DOI: 10.
262 1073/pnas.2319718121.
- 263 Sam Dillavou, Marcelo Guzman, Andrea J. Liu, and Douglas J. Durian. Understanding and Em-
264 bracing Imperfection in Physical Learning Networks, May 2025.
- 265 Gabriel Dulac-Arnold, Nir Levine, Daniel J. Mankowitz, Jerry Li, Cosmin Paduraru, Sven Goyal,
266 and Todd Hester. Challenges of real-world reinforcement learning: Definitions, benchmarks and
267 analysis. *Machine Learning*, 110(9):2419–2468, September 2021. ISSN 1573-0565. DOI: 10.
268 1007/s10994-021-05961-4.
- 269 Mohamed Elsayed, Gautham Vasan, and A Rupam Mahmood. Streaming deep reinforcement learn-
270 ing finally works. *arXiv preprint arXiv:2410.14606*, 2024.
- 271 Michael C. Granovetter, Sophia Robert, Leah Ettensohn, and Marlene Behrmann. With childhood
272 hemispherectomy, one hemisphere can support—but is suboptimal for—word and face recogni-
273 tion. *Proceedings of the National Academy of Sciences*, 119(44):e2212936119, November 2022.
274 DOI: 10.1073/pnas.2212936119.
- 275 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
276 maximum entropy deep reinforcement learning with a stochastic actor. In *International confer-
277 ence on machine learning*, pp. 1861–1870. Pmlr, 2018.
- 278 Khurram Javed and Richard S Sutton. The big world hypothesis and its ramifications for artificial
279 intelligence. In *Finding the Frame: An RLC Workshop for Examining Conceptual Frameworks*,
280 2024.
- 281 Terrence Mak, Kai-Pui Lam, H. S. Ng, Guy Rachmuth, and Chi-Sang Poon. A CMOS Current-Mode
282 Dynamic Programming Circuit. *IEEE Transactions on Circuits and Systems I: Regular Papers*,
283 57(12):3112–3123, December 2010. ISSN 1558-0806. DOI: 10.1109/TCSI.2010.2052661.

- 284 Terrence S.T. Mak, K.P. Lam, H.S. Ng, G. Rachmuth, and C.-S. Poon. A Current-Mode Analog
285 Circuit for Reinforcement Learning Problems. In *2007 IEEE International Symposium on Circuits
286 and Systems (ISCAS)*, pp. 1301–1304, May 2007. DOI: 10.1109/ISCAS.2007.378410.
- 287 C. Mead. Neuromorphic electronic systems. *Proceedings of the IEEE*, 78(10):1629–1636, October
288 1990. ISSN 1558-2256. DOI: 10.1109/5.58356.
- 289 Mantas Mikaitis, Garibaldi Pineda García, James C. Knight, and Steve B. Furber. Neuromodulated
290 Synaptic Plasticity on the SpiNNaker Neuromorphic System. *Frontiers in Neuroscience*, 12,
291 February 2018. ISSN 1662-453X. DOI: 10.3389/fnins.2018.00105.
- 292 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wier-
293 stra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint
294 arXiv:1312.5602*, 2013.
- 295 Javier R. Movellan. Contrastive Hebbian Learning in the Continuous Hopfield Model. In David S.
296 Touretzky, Jeffrey L. Elman, Terrence J. Sejnowski, and Geoffrey E. Hinton (eds.), *Connectionist
297 Models*, pp. 10–17. Morgan Kaufmann, January 1991. ISBN 978-1-4832-1448-1. DOI: 10.1016/
298 B978-1-4832-1448-1.50007-X.
- 299 Benjamin Scellier and Yoshua Bengio. Equilibrium Propagation: Bridging the Gap between Energy-
300 Based Models and Backpropagation. *Frontiers in Computational Neuroscience*, 11, 2017. ISSN
301 1662-5188. DOI: 10.3389/fncom.2017.00024.
- 302 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
303 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 304 Menachem Stern, Daniel Hexner, Jason W. Rocks, and Andrea J. Liu. Supervised Learning in
305 Physical Networks: From Machine Learning to Learning Machines. *Physical Review X*, 11(2):
306 021045, May 2021. DOI: 10.1103/PhysRevX.11.021045.
- 307 Menachem Stern, Sam Dillavou, Dinesh Jayaraman, Douglas J. Durian, and Andrea J. Liu. Training
308 self-learning circuits for power-efficient solutions. *APL Machine Learning*, 2(1):016114, Febru-
309 ary 2024. ISSN 2770-9019. DOI: 10.1063/5.0181382.
- 310 Xueying Wang, John A Hayes, Ann L Reville, Hanbing Song, Andrew Kottick, Nikolas C Vann,
311 M Drew LaMar, Maria Cristina D Picardo, Victoria T Akins, Gregory D Funk, and Christopher A
312 Del Negro. Laser ablation of Dbx1 neurons in the pre-Bötzinger complex stops inspiratory rhythm
313 and impairs output in neonatal mice. *eLife*, 3:e03427, July 2014. ISSN 2050-084X. DOI: 10.
314 7554/eLife.03427.

Algorithm 1 Q-Learning with a Contrastive Local Learning Network (CLLN)

```

1: given CLLN with nodes  $n \in \{1, \dots, M\}$ , connected by edges  $i \in \{1, \dots, E\}$  each with a
   gate voltage  $U_{G,i}$ ; environment  $(R, \text{step})$ ; hyperparameters  $\eta, \gamma, \epsilon, R_{\text{reset}}$ , total steps  $T$ , batch
   period  $B$ 
2: initialize  $U_{G,i} \sim \mathcal{N}(1.5, 0.1)$ , clipped to  $[U_G^{\min}, U_G^{\max}]$ 
3: initialize accumulated update  $\delta U_{G,i} \leftarrow 0$  for all  $i$ 
4: sample initial state  $S_0 \sim \text{Uniform}\{\text{states}\}$ 
5:
6: for  $t = 0, 1, \dots, T - 1$  do
7:   // Free state for  $S_t$ 
8:   impose input voltages encoding  $S_t$ ; equilibrate (free)  $\Rightarrow$  node voltages  $U^F(S_t)$ 
9:   record edge drops  $\Delta U_i^F = U_{\text{from}}^F - U_{\text{to}}^F$  and free outputs  $\mathcal{F}(S_t) = \{U_n^{\text{out},F}\}_{n=1}^{N_A}$ 
10:
11:   //  $\epsilon$ -greedy action selection
12:    $\epsilon_t \leftarrow \epsilon_0 (1 - t/T) + \epsilon_\infty (t/T)$ 
13:   if  $\text{rand}() < \epsilon_t$  then
14:      $A_t \leftarrow \text{Uniform}\{1, \dots, N_A\}$ 
15:   else
16:      $A_t \leftarrow \arg \max_n \mathcal{F}(S_t)$ 
17:   end if
18:
19:    $R_t \leftarrow R(S_t, A_t)$ ,  $S_{t+1} \leftarrow \text{step}(S_t, A_t)$  // Record reward and step environment
20:
21:   // Bootstrap target using evaluation of next state  $S_{t+1}$ 
22:   impose  $S_{t+1}$ ; equilibrate (free)  $\Rightarrow \mathcal{F}(S_{t+1})$ 
23:    $L_t \leftarrow R_t + \gamma[\max(\mathcal{F}(S_{t+1})) - \text{mean}(\mathcal{F}(S_{t+1}))]$ 
24:
25:   // Clamped state for  $S_t$ 
26:   labels  $L_n \leftarrow U_n^{\text{out},F}(S_t)$  for  $n \neq A_t$ ,  $L_{A_t} \leftarrow L_t$  // measured outputs for actions not taken
27:   nudge outputs:  $U_n^{\text{out},C} \leftarrow (1 - \eta) U_n^{\text{out},F} + \eta L_n$  for  $n = 1, \dots, N_A$ 
28:   impose inputs of  $S_t$  and  $U_n^{\text{out},C}$  on outputs; equilibrate  $\Rightarrow U^C(S_t)$ 
29:   record clamped edge drops  $\Delta U_i^C = U_{\text{from}(i)}^C - U_{\text{to}(i)}^C$ 
30:   // Coupled-Learning contrastive update (local rule)
31:   for  $i = 1, \dots, E$  do
32:      $\delta U_{G,i} \leftarrow \delta U_{G,i} + \alpha/\eta [(\Delta U_i^F)^2 - (\Delta U_i^C)^2]$ 
33:   end for
34:   // Apply batched update every  $B$  steps
35:   if  $(t + 1) \bmod B = 0$  then
36:      $\delta U_{G,i} \leftarrow \text{clip}(\delta U_{G,i}, -c, +c)$ 
37:      $U_{G,i} \leftarrow \text{clip}(U_{G,i} + \delta U_{G,i}, U_G^{\min}, U_G^{\max})$ 
38:      $\delta U_{G,i} \leftarrow 0$ 
39:   end if
40:   // State carry-over and periodic random reset
41:    $S_t \leftarrow S_{t+1}$ 
42:   if  $(t + 1) \bmod R_{\text{reset}} = 0$  then
43:      $S_t \leftarrow \text{Uniform}\{\text{states}\}$ 
44:   end if
45: end for

```
