

---

# Backpropagation-Free Parallel Deep Reinforcement Learning

---

**William H. Guss**

Machine Learning at Berkeley  
2650 Durant Ave, Berkeley CA, 94720  
wguss@ml.berkeley.edu

**Mike Zhong**

Machine Learning at Berkeley  
Berkeley CA, 94720  
mlyzhong@berkeley.edu

**Utkarsh S**

Machine Learning at Berkeley  
Berkeley CA, 94720  
philkuz@ml.berkeley.edu

**Max Johansen**

Machine Learning at Berkeley  
Berkeley CA, 94720  
max@ml.berkeley.edu

## Abstract

In this paper we conjecture that an agent, environment pair  $(\pi, E)$  trained using DDPG with an actor network  $\mu$  and critic network  $Q^\pi$  can be decomposed into a number of sub-agent, sub-environment pairs  $(\pi_n, E_n)$  ranging over every neuron in  $\mu$ ; that is, we show empirically that treating each neuron  $n$  as an agent  $\pi_n : \mathbb{R}^n \rightarrow \mathbb{R}$  of its inputs and optimizing a value function  $Q^{\pi_n}$  with respect to the weights of  $\pi_n$  is dual to optimizing  $Q^\pi$  with respect to the weights of  $\mu$ . Finally we propose a learning rule which simultaneously optimizes each  $\pi_n$  without error backpropagation achieving state of the art performance and speed across a variety of OpenAI Gym environments.

## Todo list

Introduction to DDPG and recent advances in deep RL. . . . .	2
Biological diffusion of dopamine in the brain $\implies$ error backpropagation is not biologically feasible. . . . .	2
Synthetic gradients are a step in the right direction, but still require eventual back propagation. . . . .	2
Therefore it is feasible that each neuron is maximizing the expectation on his future dopamine intake, and so we propose the following theorem. . . . .	2
A high level description of the section. . . . .	3
Fix the time thing. Remove $\ell$ . . . . .	3
FIX THIS WHEN FIXING $\ell$ . The environment is NOT fully observed. . . . .	3
Insert reference and make this a footnote . . . . .	3
Cite lilicrap . . . . .	3
Make DDPG figure . . . . .	3
Cite deepmind . . . . .	3
fix this equation, $\epsilon$ doesn't directly commute . . . . .	5
Prove the Decomposition . . . . .	5

Emperical justification of the iff using the following experiment (s). . . . .	6
1. Training a network on Atari using DDPG and plotting average critic functions for neurons using window. . . . .	6
Introduce the notion of similary used to compare the q-functions. . . . .	7
Incorporate metrics regarding the performance $Q_n$ w.r.t $Q^\mu$ 's weights . . . . .	7
Therefore we propose the following learning rule in aims to evidence the reverse, training $\mu$ using simultaneous optimization on all $Q_n$ w.r.t $\pi_n$ 's weights. . . . .	8
Proposal of the rule. Linear approximation of the $Q$ function for every neuron is good enough, (experimentally). . . . .	9
Implications of the rule to DDPG . . . . .	9
Implications of the rule to entirely recurrent networks (infinite time horizion and NO unrolling since the environment the local actions of the neuron which globally recur to that neuron again are <i>encoded</i> into $Q_n$ ; large time horizion probably implies that better regresser needed for $Q_n$ .) . . . . .	9
Parallelism, no error backprop, and only 2x operations, but no locking on GPU, so all can be run sumultaneously if we cache! . . . . .	9
To validate the new learning rule we throw a fuck ton of experiments together on the following list (or better using OpenAI Gym). . . . .	9
2. Show that training decentralized policy gradient $\implies$ total policy optimization . . . . .	9
3. Show speed improvements on update step through parallelism (samples per second vs DDPG). . . . .	9
4. Show results are comparable with the state of the art. . . . .	9
We wrecked deep reinforcement learning using biological inspiration. . . . .	10
Would like to try the method with full recurrent networks and purely asynchronous imple- mentation of leaky integration networks. . . . .	10
Would like to prove the conjecture. List possible methods of proof. . . . .	10

## 1 Introduction

Introduction to DDPG and recent advances in deep RL.

Biological diffusion of dopamine in the brain  $\implies$  error backpropagation is not biologically feasible.

Synthetic gradients are a step in the right direction, but still require eventual back propogation.

Recent work has attempted to address the issue of decoupling connections between layers in the network using decoupled neural interfaces [?]. The initial coupling is due to the reliance on the forwards and backwards computations to compute the backpropagated error gradient for a given layer. Synthetic gradient modules model the error gradient using only local information, allowing immediate feedback that is later corrected when the backpropogated error is finally computed. This method has had success in both feedforward and recurrent architectures, where the network is able to model much greater time horizons.

Although synthetic gradients decouple the network modules to some degree, this technique still relies on synchronous backpropagation to learn the local regressors. In addition, they model at a lower granularity, learning layer-level gradients instead of those of the individual neurons.

Therefore it is feasible that each neuron is maximizing the expectation on his future dopamine intake, and so we propose the following theorem.

## 2 Agent-Environment Value Decomposition

A high level description of the section.

Fix the time thing. Remove  $\ell$ .

### 2.1 Background

Recall the standard reinforcement learning setup. We say  $E$  is an *environment* if  $E \stackrel{\text{def}}{=} (\mathcal{S}, \mathcal{A}, \mathcal{R}, T, r)$  where  $T$  describes transition probability measure  $T(s_{t+1} | s_t, a_t)$  and  $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{R}$  is a reward function. Furthermore  $\mathcal{S}, \mathcal{A}, \mathcal{R}$  are the *state space*, *action space*, and *reward space* respectively. We restrict  $\mathcal{R}$  to a compact subset of  $\mathbb{R}$  and action space and state space to finite dimensional real vector spaces. As in DDPG we assume that the environment  $E$  is *fully observed*

**FIX THIS WHEN FIXING  $\ell$ . The environment is NOT fully observed.**

; that is, at any time step the state  $s_t$  is fully described by the observation presented,  $x_t$ , and not by the history  $(x_1, a_1, \dots, a_{t-1})$ .

We define the policy for an agent to be  $\pi : \mathcal{P}(\mathcal{A}) \times \mathcal{S} \rightarrow [0, 1]$ . In general the policy is a probability measure on some  $\sigma$ -algebra  $\mathcal{M} \subset \mathcal{P}(\mathcal{A})$  conditioned on  $\mathcal{S}$  so that  $\pi(\mathcal{A} | s \in \mathcal{S}) = 1$ . However, we will deal only with *deterministic* policies where for every  $s_t$  there is unique  $a_t$  so that  $\pi(\{a_t\} | s = s_t) = 1$  and the measure is 0 otherwise. Thus we will abuse notation and define a *deterministic agent* by a policy function  $\pi : \mathcal{S} \rightarrow \mathcal{A}$ . Additionally we denote the state-space trajectories of  $\pi$  by

$$\Gamma_\pi(\mathcal{S}) = \{(s_1, s_2, \dots) | s_1 \sim T(s_0), s_{t+1} \sim T(s | s_t, \pi(s_t))\}. \quad (2.1.1)$$

For a policy  $\pi$  the action-value function is the expected future reward under  $\pi$  by performing  $a_t$  at state  $s_t$  using the Bellman equation

$$Q^\pi(s_t, a_t) = \mathbb{E}_{s_{t+1} \sim E} [r(s_t, a_t) + \gamma Q^\pi(s_{t+1}, \pi(s_{t+1}))] \quad (2.1.2)$$

with  $\gamma \in (0, 1)$  a discount factor, and the second expectation removed because  $\pi$  is deterministic. [Some survey] provides an extensive exposition into a justification of this equation and choice for the action-value of  $\pi$ , so we will assume such a choice is a valid measure of performance.

In deterministic policy gradient methods, we define an actor  $\mu : \mathcal{S} \rightarrow \mathcal{A}$  and a critic  $Q^\mu : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  and optimize  $Q^\mu(s_t, \mu(s_t | \theta^\mu))$  with respect to the parameters  $\theta^\mu$  of  $\mu$ . This method is provably the true policy gradient of  $\mu$  if  $Q^\mu$  is known. Recently (DDPG) utilizes the universality of DNNs in order to approximate both  $\mu$  and  $Q^\mu$  along with delayed weight-transfer networks to stabilize learning and prevent divergence as depicted in Figure 1. In order to decompose the action-value function we will make heavy use of this methodology at a scale local to each neuron in the flavor of (Synthetic gradients.)

### 2.2 Towards Neurocomputational Decomposition of $Q^\mu$

In order to decompose the  $Q^\mu$  algorithm we will abstractly define a neurocomputational agent in terms of an operator on voltages with no restrictions on the topology of the network, and then relate the action-value function of the whole agent to those which are defined for each individual neuron in the network.

If  $\mathcal{V}$  is an  $N$ -dimensional vector space then a *neurocomputational agent* is a tuple  $\mathcal{N} = (\mu, \epsilon, \delta, K, \Theta, \sigma, D)$  such that:

Insert reference and make this a footnote

Cite lil-crap

Make DDPG figure

Cite deep-mind

- $\epsilon : \mathcal{S} \rightarrow \mathcal{V}$  encodes the state into the voltages. Realistically, only a subset of all neurons are input neurons, denoted as  $N_I \subset \mathcal{V}$ , so  $\epsilon(s_t) = \text{proj}_{N_I}(\epsilon(s_t))$ .
- $\delta : \mathcal{V} \rightarrow \mathcal{A}$  decodes the voltages of the *output neurons*  $N_O \subset \mathcal{V}$  into an action, so that  $\delta(v_t) = \delta(\text{proj}_{N_O}(v_t))$ .
- $K : \mathcal{V} \rightarrow \mathcal{V}$  is the linear voltage graph transition function of the graph representing the topology of  $\mathcal{N}$ , parameterized by  $\theta$ .
- $\Theta : \mathcal{V} \rightarrow \mathcal{V}$  is a nonlinear inhibition function.
- $\sigma : \mathcal{V} \rightarrow \mathcal{V}$  is the elementwise application of some activation function to the voltage vector.
- $D : \mathcal{V} \times \mathcal{V} \rightarrow \mathcal{V}$  is called voltage dynamic of  $\mathcal{N}$  such that

$$v_{cur+1} \stackrel{\text{def}}{=} D(v_{cur}, v_{in}) \stackrel{\text{def}}{=} \sigma(\Theta K[v_{cur}]) + v_{in} \quad (2.2.1)$$

where  $v_{cur}$  is the internal voltage vector of  $\mathcal{N}$  and  $v_{in}$  is an input voltage. We will occasionally abuse notation and say that  $D(v_{cur}) = D(v_{cur}, 0)$  when  $v_{in}$  is 0.

- $\mu : \mathcal{S} \rightarrow \mathcal{A}$  is the deterministic policy for  $\mathcal{N}$ . For some agents, the internal time  $\tau$  is not in sync with the discrete time step  $t$  of  $E$ . Therefore for every  $t$  there is an evaluation delay  $\ell$  so that

$$\begin{aligned} v_{t+1} &= D(v_t) + \epsilon(s_t) \\ &= D(v_{t-1}, \epsilon(s_t)) \\ \mu(s_t) &= \delta(v_t) \end{aligned} \quad (2.2.2)$$

It is not hard to see that this definition encompasses any DQN or DDPG network with either recurrent or non recurrent layers. Additionally other paradigms such as the leaky integrator are neuro-computational agents. (See appendix.)

If  $n$  is some neuron in  $\mathcal{N}$ , we say  $E^n = (\mathcal{S}^n, \mathcal{A}^n, \mathcal{R}^n, T^n, r^n)$  is deterministic *sub-environment* of  $E$  with respect to  $\mathcal{N}$  if it has the following properties:

- an environment  $\mathcal{S}^n = \mathcal{V}$ , in that the environment that neuron  $n$  sees is the voltages of all other neurons. We assume that the graph is fully connected so that all neurons are connected to all other neurons. Realistically, each neuron only sees the voltages of a subset of all neurons. We will deal with this in a further iteration of the paper.
- an action  $\mathcal{A}^n = \mathbb{R}$  denotes the output voltage of neuron  $n$  one time step later, i.e.  $\alpha_t^n = v_{t+1}^n$ .
- a reward  $\mathcal{R}^n = \mathbb{R}$ , denoting the reward that neuron  $n$  receives. It is the same as the reward that the entire agent sees.
- a transition function  $T^n : \mathcal{S}^n \times \mathcal{A}^n \rightarrow \mathcal{S}^n$  such that  $T^n(v_t, \alpha_t^n) = (I - \delta_{n,n})D(v_t, \epsilon(s_t)) + e_n \alpha_t^n$
- and a reward function  $r^n(v_\tau, \alpha) = r(s_t, \mu(s_t))$  if  $I(\tau) \neq 0$  and 0 otherwise. Essentially the state space of  $E^n$  at time step  $\tau + 1$  is just the normal dynamics of  $\mathcal{N}$  applied to the previous state along with a possible encoded input state  $\epsilon(s_t)$  from  $E$  except for at neuron  $n$ . Lastly an agent  $\mu^n : \mathcal{V} \rightarrow \mathbb{R}$  is called *neuromorphically local* to  $\mathcal{N}$  if  $v_\tau \mapsto \langle D(v_\tau), e_n \rangle$ .

We now can think of every neuron in  $\mathcal{N}$  as an agent in its own environment, acting on its inputs, and we can extend the action-value definition to  $\mu^n$  as follows

$$Q^{\mu^n}(v_\tau, \alpha_\tau) = \mathbb{E}_{v_{\tau+1} \sim E^n} \left[ r^n(v_\tau, \alpha_\tau) + \gamma Q^{\mu^n}(v_{\tau+1}, \mu^n(v_{\tau+1})) \right]. \quad (2.2.3)$$

### 2.2.1 Results

Provided with the previous definitions, the following question arises: does deterministic policy gradient learning on  $\mathcal{N}$ , specifically  $\mu$  on  $E$ , *commute* with performing the same operation simultaneously on every neuromorphically local agent  $\mu^n$  comprising  $\mathcal{N}$  and their respective sub-environments  $E^n$ ? Supposing that we have the true  $Q^\mu$  function and  $\mu$  is optimal with respect to  $Q^\mu$ , then it is intuitive, but not obvious, that every  $\mu^n$  should behave optimally with respect to an infinite time horizon – but will the reverse hold? We give the following results:

**Theorem 2.2.1.** Let  $E$  and  $\mathcal{N}$  be defined as before. Then for every  $n \in \mathcal{N}$ , it follows that  $\Gamma_\mu(\mathcal{S})$  is equal to  $\Gamma_{\mu^n}(\mathcal{V})$  up to bijection and the following diagram commutes.

$$\begin{array}{ccc}
 \mathcal{V} \times \mathcal{S} & \xrightarrow{\mu \circ \pi_2} & \mathcal{A} \\
 \downarrow \text{id}_{\mathcal{V}} \times \epsilon & & \uparrow \delta \\
 \underbrace{(\mathcal{V} \times \mathcal{V})}_{(v_\tau, \epsilon(s_t))} & \xrightarrow{D} & \underbrace{\mathcal{V}}_{v_{\tau+1}} & \xrightarrow{D} & \underbrace{\mathcal{V} \rightarrow \dots \rightarrow \mathcal{V}}_{v_{\tau+2}, \dots, v_{\tau+\ell-1}} & \xrightarrow{D} & \mathcal{V} \\
 & \searrow \mu^n \circ \pi_1 + \pi_n \circ \pi_2 & \downarrow \pi_n & \searrow \mu^n & \downarrow \pi_n & \searrow \mu^n & \downarrow \pi_n \\
 & \mathbb{R} & \mathbb{R} \dots \mathbb{R} & & \mathbb{R}
 \end{array} \tag{2.2.4}$$

*Proof.* We first show that (2.2.4) commutes. Let  $v \in \mathcal{V}$  and  $s \in \mathcal{S}$ . Observe that  $\mu(\mathcal{S}) =$  is clear that  $V[b, c] \circ V[a, b] = V[a, c]$  by (??), so the upper part of the diagram is equivalent to

$$\begin{array}{ccc}
 \mathcal{S} & \xrightarrow{\mu} & \mathcal{A} \\
 \downarrow \epsilon & & \uparrow \delta \\
 \mathcal{V} & \xrightarrow{V[\tau, \ell]} & \mathcal{V}
 \end{array} \tag{2.2.5}$$

and by definition  $\mu$  with an evaluation time  $\ell$  we have that

$$\begin{aligned}
 \mu(s_t) &= \delta(V[\tau + \ell](\sigma(K\Theta[V(\tau)]) + \epsilon(s_t))) \\
 &= \delta(V[\tau + \ell](V(\tau))).
 \end{aligned} \tag{2.2.6}$$

Next for each  $V[\tau + k, 1]$ ,  $k \in \mathbb{N} \cup \{0\}$  observe the cooresponding triangle in the diagram. When  $\pi_n$  is the canonical projection, we have

$$(\pi_n \circ V[\tau + k, 1])(v_\tau) = \langle V[\tau + k + 1](v), e_n \rangle \tag{2.2.7}$$

and by (??)

□

**Theorem 2.2.2.** If  $\mathcal{N}$  is a nuerocomputational agent in  $E$  and for every  $n \in \mathcal{N}$  there is a nueromorphically local agent  $\mu^n$  in sub-environment  $E^n$ , then policy gradient for  $\mu$  agrees with the simultaneous policy gradients of every neuromorphically local agent; that is

$$\prod_{n=1}^N \nabla_{K^n} Q^{\mu^n}(v, a) \Big|_{v=v_t, a=\mu^n(v_t)} = \nabla_K Q^\mu(s, a) \Big|_{s=s_t, a=\mu(s_t)} \tag{2.2.8}$$

for every time step  $t$ , where  $K^n$  represents the  $n$ th column of the linear voltage graph transition function, i.e. the weights of the connections from all neurons to neuron  $n$ .

*Suggestion (from Mike):*

If  $\mathcal{N}$  is a nuerocomputational agent in  $E$  and for every  $n \in \mathcal{N}$  there is a nueromorphically local agent  $\mu^n$  in sub-environment  $E^n$ , then policy gradient for  $\mu$  agrees with the simultaneous policy gradients of every neuromorphically local agent; that is

$$\forall n, \nabla_{K^n} Q^{\mu^n}(v, a) \Big|_{v=v_t, a=\mu^n(v_t)} = \nabla_{K^n} Q^\mu(s, a) \Big|_{s=s_t, a=\mu(s_t)} \tag{2.2.9}$$

for every time step  $t$ , where  $K^n$  represents the  $n$ th column of the linear voltage graph transition function, i.e. the weights of the connections from all neurons to neuron  $n$ .

*Proof.*

Prove the Decomposition

□

Emperical justification of the iff using the following experiment (s).

Intuitively, it should be the case that by treating each neuron in the neural network as its own Q-learner, the policy gradient of an individual agent should be the same as for the policy gradient entire agent, if each neuron sees the same reward  $r_t$  as the entire agent. In other words, the optimal way of updating the entire connection matrix  $K$  should be the optimal way for updating the weights connected to each neuron w.r.t. each neuron as its own Q-learner, given the same rewards.

3. Plot the output of both  $Q^\mu$  and  $Q_n$  using TENSORFLOW summaries.

From Experiment 1, we see that there is a correlation between the Q-gradient for neuron  $n$ , and the  $n$ th column of the gradient of the agent's Q function:

from strong to weak:

$$\begin{aligned}\nabla_{K^n} Q^{\mu^n}(v, a) &= \nabla_{K^n} Q^\mu(s, a) \\ \nabla_{K^n} Q^{\mu^n}(v, a) &\propto \nabla_{K^n} Q^\mu(s, a) \\ \text{corr}[\nabla_{K^n} Q^{\mu^n}(v, a), \nabla_{K^n} Q^\mu(s, a)] &\approx \text{EMPERICALVALUE}\end{aligned}$$

although (assuming we get a weak experimental result), we propose that it is due to our neurons not being the ideal Q-learner (convergence issues, etc.).

1. Training a network on Atari using DDPG and plotting average critic functions for neurons using window.

**EXPERIMENT 1 SPECIFICATION.** 1. Set up a standard DDPG to play the set of atari games in OpenAI Gym using TENSORFLOW (this will be mac,linux, or windows bash only). If on Windows bash install Xming (its an X server) and run all OpenAI Gym commands with `DISPLAY=localhost:0.0 python3 src/experiment1/some_script_in_src.py`. This will pipe the visual output fo the OpenAI Gym simulators to the display. If you cannot get this to work on your screen, do not do `env.render`. We can also stop using the atari games, since this works for a fact on the basic Box2d versions. **We are going to write all of this in Python3, make sure to install gym in python3.**

2. DDPG has a  $Q^\mu$  network which we use to optimize  $\mu(s_t | \theta)$  with respect to  $\theta$ . The goal of this experiment is to train a standard DDPG network to play one of these OpenAI Gym simulations, whilst concurrently estimating and viewing the  $Q$  functions for every single neuron. THEREFORE, we need to select a subset of neurons in the fully connected layers (for example) of the  $\mu$  network (actor) and concurrently train a network  $Q^n(s, a)$  to estimate the  $Q$  function of the neuron based on its inputs  $s$  and its SINGLE output voltage  $s$ . This can be a 3 layer fully connected network with  $|s| + 1$  inputs (one for each input the neuron  $n$  and 1 for the voltage of the neuron after receiving that output.) Tensorflow is a dataflow language so the output of a layer looks like  $O2 = \sigma(W * O1)$ . Therefore you just need to make another "network" whose dataflow could be like  $Q_n = \sigma(W_3^n * \sigma(W_2^n * \sigma(W_1^n * \text{concat}(O1, O2[i])))$  where  $n$  is the  $i$ th neuron on layer  $O2$ . And  $O1, O2, O3$  are the outputs of the neurons on those layers in the network. Then as in standard DQN you train  $Q_n$  with a lag network  $W(Q_n') = W(Q_n)(1 - \tau) + \tau W(Q_n)$  where  $W$  denotes the weights of  $Q_n$ , say  $W_3^n, W_2^n, W_1^n, \dots$ . And then actually do gradient decent on the weights of  $Q_n$  not  $Q_n'$  by minimizing the following bellman equation

$$L(s_t, a_t, r_t, s_{t+1}, a_{t+1}) = (Q_n(O1(s_t), O2[i](s_t)) - r_t - Q_n'(O1(s_{t+1}), O2[i](s_{t+1})))^2$$

with respect to the parameters  $W(Q_n)$ . Note I didn't actually use  $a_t$  above since really the  $Q_n$  function takes in the input of the previous layer (to  $n$ ) as its input, say  $O1(s_t)$  and the action for that same time step which is just the output of the neuron  $n$ , say  $O2(s_t)[i]$ . The same goes for  $Q_n'$  but at the next time step.

### 3 Experimentation

#### 3.1 Experiment 1: Learning Q functions on components of the Actor network.

**Motivation:** Our first batch of experiments aimed to establish empirically that the Q-functions of each layer, denoted  $Q_1 \dots Q_L$  are similar to the Q-function of the overall network, denoted  $Q_\mu$ . To

that end, we extended the actor-critic methodology used in Lillicrap et al (2016) as follows. In addition to training  $Q^\mu$ , to estimate the Q functions of the actor network  $\mu$ , sub-critic networks  $Q^n$  were initialized for each individual layer. We then compare the  $Q$  values estimated by the subcritics to the  $Q$  estimation provided by the main critic,  $Q^\mu$ .

The inputs used in the calculation of the Q-functions for each layer follow the conceptual framework laid out in Section 2.2. The state and action parameters are therefore voltages represented by the outputs of  $\epsilon : \mathcal{S} \rightarrow \mathcal{V}$  and  $\delta : \mathcal{V} \rightarrow \mathcal{A}$  respectively, where the changes in voltage are described by a voltage transition function  $K : \mathcal{V} \rightarrow \mathcal{V}$ .

To determine and compare the rate at which the subcritics and the main critic learn, the experiment was run in two phases. First, each of the subcritics and the main critic were trained using the standard DDPG algorithm on some actor  $\mu$ . In the second phase, a new actor  $\mu'$  was initialized and its Q-function set to  $Q_\mu$  as determined above in phase 1. The subcritics were trained, and the values of  $Q_1 \dots Q_L$  were plotted as training occurred.

Introduce the notion of similarity used to compare the q-functions.

After [number of iterations], each  $Q^n(s, a)$  networks effectively learned the same cost function as  $Q^\mu(s, a)$ .

Incorporate metrics regarding the performance  $Q_n$  w.r.t  $Q^\mu$ 's weights

. Furthermore, the  $Q$ -value plots for each  $Q^n$  correlate directly with the  $Q^\mu$ 's  $Q$  plot. This confirms that we are able to teach the  $Q^n$  model on the level of subcomponents of  $\mu$  supporting the argument that we can use  $Q^n$  nets to approximate  $Q^\mu$ . We can now use each  $Q^n$  to calculate the gradient of the specific neuron component that the  $Q^n$  is critiquing. We proceed with this in Experiment 2.

### 3.2 Experiment 2: Treating each neuron as its Actor-Critic network using linear approximators

Now that we have established that the Critic networks for each of the individual neuron learn the same Q-function as does the entire agent, we now treat each neuron as its own actor. Each neuron changes the weights of its presynaptic neurons' connections, as parameters for its actor – its voltage on the next timestep – to optimize its learned approximated Q function. We use the linear approximation:

$$Q^n(v, a) \approx \theta_{n,v}^T v + \theta_{n,a} a = \theta_n^T(v, a)^T$$

$$\mu^n(v) = \sigma(K^n v)$$

GOT RID OF  $\Theta$  in that it's not very necessary ...

The algorithm for learning (very much INSPIRED by [CONTINUOUS CONTROL WITH DEEP REINFORCEMENT LEARNING]):

For each neuron n:

Initialize random weights  $\theta_{n,v}$ ,  $\theta_{n,a}$ , and  $K^n$

Initialize target network  $Q^{n'}$  and  $\mu^{n'}$  with same weights,  $\theta'_{n,v}$ ,  $\theta'_{n,a}$ , and  $K'^n$

Initialize replay buffer R

For each episode:

Receive initial observation voltages  $v_1$

for  $t=1, T$ , do:

follow dynamics,  $a_{t+1} = \sigma(K^n v_t)$

Observe reward  $r_t$ , observe new voltages  $v_{t+1}$

Store transition  $(v_t, a_t, r_t, v_{t+1})$  in R.

Sample a random minibatch of  $N$  transitions  $(v_i, a_i, r_i, v_{i+1})$  from R

Set  $y_i = r_i + \gamma Q'(v_{i+1}, \mu'(v_{i+1}))$

Update critic weights  $\theta_n$  by minimizing loss  $L := \frac{1}{N} \sum_i (y_i - Q(v_i, a_i))^2$   
 Update actor weights (connections) using the sample policy gradient:

$$\nabla_{K^n} J \approx \frac{1}{N} \sum_i \nabla_a Q(v, a)|_{v=v_i, a=\mu^n(v_i)} \nabla_{K^n} \mu^n(v)|_{v_i}$$

update the target networks:

$$\theta'_n \leftarrow \tau \theta_n + (1 - \tau) \theta'_n$$

$$K^{n'} \leftarrow \tau K^n + (1 - \tau) K^{n'}$$

end for

end for

We train the critic weights  $(\theta_{n,v}, \theta_{n,a})$  by minimizing the loss:

$$L = \frac{1}{2} ((r_i + \gamma Q^{n'}(v_{i+1}, a_{i+1})) - Q(v_i, a_i))^2 = \frac{1}{2} D_i^2$$

where  $D_i := (r_i + \gamma \theta_n^{T'}(v_{i+1}, a_{i+1})) - \theta_n(v_i, a_i)$

We thus have:

$$\nabla_{\theta_n} L = -D_i(v_i, a_i)$$

or, using gradient descent with learning rate  $\eta_Q$ , we have:

$$\theta_n \leftarrow \theta_n - \eta_Q \nabla_{\theta_n} L = \theta_n + \eta_Q D_i(v_i, a_i)$$

or:

$$\theta_{n,v} \leftarrow \theta_{n,v} + \eta_Q D_i v_i$$

and

$$\theta_{n,a} \leftarrow \theta_{n,a} + \eta_Q D_i a_i$$

As for updating the actor policy, with learning rate  $\eta_A$ , we have:

$$K^n \leftarrow K^n + \eta_A \nabla_a Q(v, a)|_{v=v_i, a=\mu^n(v_i)} \nabla_{K^n} \mu^n(v)|_{v_i} = K^n + \eta_A \theta_{n,a} \sigma'(K^n v_i) v_i$$

IN CONCLUSION:

$$\theta_{n,v} \leftarrow \theta_{n,v} + \frac{1}{N} \sum_i \eta_Q D_i v_i$$

$$\theta_{n,a} \leftarrow \theta_{n,a} + \frac{1}{N} \sum_i \eta_Q D_i a_i$$

$$K^n \leftarrow K^n + \frac{1}{N} \sum_i \eta_A \theta_{n,a} \sigma'(K^n v_i) v_i$$

where  $D_i := (r_i + \gamma \theta_n^{T'}(v_{i+1}, a_{i+1})) - \theta_n(v_i, a_i)$

Therefore we propose the following learning rule in aims to evidence the reverse, training  $\mu$  using simultaneous optimization on all  $Q_n$  w.r.t  $\pi_n$ 's weights.



## 4 Decentralized Deep Deterministic Policy Gradient Learning

Proposal of the rule. Linear approximation of the  $Q$  function for every neuron is good enough, (experimentally).

Implications of the rule to DDPG

Implications of the rule to entirely recurrent networks (infinite time horizon and NO unrolling since the environment the local actions of the neuron which globally recur to that neuron again are *encoded* into  $Q_n$ ; large time horizon probably implies that better regressor needed for  $Q_n$ .)

Parallelism, no error backprop, and only 2x operations, but no locking on GPU, so all can be run simultaneously if we cache!

## 5 Results

To validate the new learning rule we throw a fuck ton of experiments together on the following list (or better using OpenAI Gym).

```
blockworld1 1.156 1.511 0.466 1.299 -0.080 1.260
blockworld3da 0.340 0.705 0.889 2.225 -0.139 0.658
canada 0.303 1.735 0.176 0.688 0.125 1.157
canada2d 0.400 0.978 -0.285 0.119 -0.045 0.701
cart 0.938 1.336 1.096 1.258 0.343 1.216
cartpole 0.844 1.115 0.482 1.138 0.244 0.755
cartpoleBalance 0.951 1.000 0.335 0.996 -0.468 0.528
cartpoleParallelDouble 0.549 0.900 0.188 0.323 0.197 0.572
cartpoleSerialDouble 0.272 0.719 0.195 0.642 0.143 0.701
cartpoleSerialTriple 0.736 0.946 0.412 0.427 0.583 0.942
cheetah 0.903 1.206 0.457 0.792 -0.008 0.425
fixedReacher 0.849 1.021 0.693 0.981 0.259 0.927
fixedReacherDouble 0.924 0.996 0.872 0.943 0.290 0.995
fixedReacherSingle 0.954 1.000 0.827 0.995 0.620 0.999
gripper 0.655 0.972 0.406 0.790 0.461 0.816
gripperRandom 0.618 0.937 0.082 0.791 0.557 0.808
hardCheetah 1.311 1.990 1.204 1.431 -0.031 1.411
hopper 0.676 0.936 0.112 0.924 0.078 0.917
hyq 0.416 0.722 0.234 0.672 0.198 0.618
movingGripper 0.474 0.936 0.480 0.644 0.416 0.805
pendulum 0.946 1.021 0.663 1.055 0.099 0.951
reacher 0.720 0.987 0.194 0.878 0.231 0.953
reacher3daFixedTarget 0.585 0.943 0.453 0.922 0.204 0.631
reacher3daRandomTarget 0.467 0.739 0.374 0.735 -0.046 0.158
reacherSingle 0.981 1.102 1.000 1.083 1.010 1.083
walker2d 0.705 1.573 0.944 1.476 0.393 1.397
```

2. Show that training decentralized policy gradient  $\implies$  total policy optimization

3. Show speed improvements on update step through parallelism (samples per second vs DDPG).

4. Show results are comparable with the state of the art.

## 6 Conclusion

We wrecked deep reinforcement learning using biological inspiration.

### 6.1 Future Work

Would like to try the method with full recurrent networks and purely asynchronous implementation of leaky integration networks.

Would like to prove the conjecture. List possible methods of proof.