
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
lol@gmail.com

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 (π, E) trained using DDPG with an actor network μ and critic network Q^π can be decomposed into a number of sub-agent, sub-environment pairs (π_n, E_n) ranging over every neuron in μ ; 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^{π_n} with respect to the weights of π_n is dual to optimizing Q^π with respect to the weights of μ . Finally we propose a learning rule which simultaneously optimizes each π_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.	2
Insert reference and make this a footnote	3
Cite lilicrap	3
Make DDPG figure	3
Cite deepmind	3
fix this equation, ϵ doesn't directly commute	4
IN PROGRESS: Write conjecture on decomposition which is free of neural configuration.	
Subject to change in later versions of ArXiv paper	5
Emperical justification of the iff using the following experiment (s).	5

1. Training a network on Atari using DDPG and plotting average critic functions for neurons using window.	5
2. Possibly others.	5
Therefore we propose the following learning rule in aims to evidence the reverse, training μ using simultaneous optimization on all Q_n w.r.t π_n 's weights.	5
Proposal of the rule. Linear approximation of the Q function for every neuron is good enough, (experimentally).	5
Implications of the rule to DDPG	5
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 <i>encoded</i> into Q_n ; large time horizon probably implies that better regressor needed for Q_n .)	5
Parallelism, no error backprop, and only 2x operations, but no locking on GPU, so all can be run simultaneously if we cache!	5
To validate the new learning rule we throw a fuck ton of experiments together on the following list (or better using OpenAI Gym).	5
1. Show that training decentralized policy gradient \implies total policy optimization	6
2. Show speed improvements on update step through parallelism (samples per second vs DDPG).	6
3. Show results are comparable with the state of the art.	6
We wrecked deep reinforcement learning using biological inspiration.	6
Would like to try the method with full recurrent networks and purely asynchronous implementation of leaky integration networks.	6
Would like to prove the conjecture. List possible methods of proof.	6

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 propagation.

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.

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 di-

mensional real vector spaces. As in DDPG we assume that the environment E is *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 σ -algebra $\mathcal{M} \subset \mathcal{P}(\mathcal{A})$ conditioned on \mathcal{S} so that $\pi(\mathcal{A} \mid s \in \mathcal{S}) = 1$. However, we will deal only with *deterministic* policies where for every s_t there is unique a_t so that $\pi(\{a_t\} \mid s = s_t) = 1$ 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 π by

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

For a policy π the action-value function is the expected future reward under π 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 π is deterministic. [Some survey] provides an extensive exposition into a justification of this equation and choice for the action-value of π , 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 $Q^\mu(s_t, \mu(s_t))$ with respect to the parameters θ^μ of μ . This method is provably the true policy gradient of μ if Q^μ is known. Recently (DDPG) utilizes the universality of DNNs in order to approximate both μ and Q^μ 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^μ

In order to decompose the Q^μ 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:

- $\epsilon : \mathcal{S} \rightarrow N_I \subset \mathcal{V}$ encodes the state into the voltages of *input neurons*, a subspace N_I of the voltages $\mathcal{V} \subset \mathbb{R}^N$ of every neuron in the network. Specifically $\epsilon(s_t) = \text{proj}_{N_I}(s_t)$.
- $\delta : \mathcal{V} \rightarrow \mathcal{A}$ decodes the voltages of the *output neurons* $N_O \subset \mathcal{V}$ into an action.
- $K : \mathcal{V} \rightarrow \mathcal{V}$ is the linear voltage graph transition function of the graph representing the topology of \mathcal{N} , parameterized by θ .
- $\Theta : \mathcal{V} \rightarrow \mathcal{V}$ 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

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

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 τ is not in sync with the discrete time step t of E . Therefore for every t there is an evaluation delay ℓ so that

$$\begin{aligned} V(s_t) &= \underbrace{D \circ D \circ \dots \circ D}_{\ell \text{ times}}(V(s_{t-1}), \epsilon(s_t)) & V(s_0) &= D(0, \epsilon(s_0)) \\ \mu(s_t) &= \delta(V(s_t)) \end{aligned} \quad (2.2.2)$$

For example if \mathcal{N} standard ℓ layer DNN, then the policy decodes a voltage after an ℓ step delay.

Insert reference and make this a footnote

Cite lil-crap

Make DDPG figure

Cite deepmind

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 $E^n = (\mathcal{V}, \mathbb{R}, \mathcal{R}, T^n, r^n)$ is deterministic *sub-environment* of E with respect to \mathcal{N} if it has the following properties: a discrete time step τ not necessarily one-to-one with the time t of E ; an input function $I(\tau) = \epsilon(s_t)$ if s_t presented at time τ and 0 otherwise; a transition function T^n such that $v_\tau \in \mathcal{V}$, $\alpha_\tau \in \mathbb{R}$, and $v_{\tau+1}^j = f(v_\tau, I(\tau)) \langle D(v_\tau, I(\tau)), e_j \rangle$ for all $j \neq n$ and $v_{\tau+1}^n = \alpha_\tau$, then

$$T^n(v_{\tau+1} | v_\tau, \alpha_\tau) = 1; \quad (2.2.3)$$

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 μ^n as follows

$$Q^{\mu^n}(v_\tau, \alpha_\tau) = \mathbb{E}_{v_{\tau+1} \sim E^n} \left[r^n(v_\tau, a_\tau) + \gamma Q^{\mu^n}(v_{\tau+1}, \mu^n(v_{\tau+1})) \right]. \quad (2.2.4)$$

2.2.1 Results

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

$$\begin{array}{ccc} \mathcal{V} \times \mathcal{S} & \xrightarrow{\mu \circ \pi_2} & \mathcal{A} \\ \downarrow \text{id}_{\mathcal{V}} \times \epsilon & & \uparrow \delta \\ \underbrace{(\underbrace{v_\tau, \epsilon(s_t)}_{\mathcal{V} \times \mathcal{V}})}_{\mathcal{V} \times \mathcal{V}} & \xrightarrow{D} \underbrace{\mathcal{V}}_{v_{\tau+1}} & \xrightarrow{D} \underbrace{\mathcal{V} \rightarrow \dots \rightarrow \mathcal{V}}_{v_{\tau+2}, \dots, v_{\tau+\ell-1}} & \xrightarrow{D} \mathcal{V} \\ & \searrow \mu^n \circ \pi_1 + \pi_n \circ \pi_2 & \searrow \mu^n & \searrow \mu^n \\ & \mathbb{R} & \mathbb{R} \dots \mathbb{R} & \mathbb{R} \end{array} \quad (2.2.5)$$

Proof. We first show that (2.2.5) 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

$$\begin{array}{ccc} \mathcal{S} & \xrightarrow{\mu} & \mathcal{A} \\ \downarrow \epsilon & & \uparrow \delta \\ \mathcal{V} & \xrightarrow{V[\tau, \ell]} & \mathcal{V} \end{array} \quad (2.2.6)$$

and by definition μ with an evaluation time ℓ we have that

$$\begin{aligned} \mu(s_t) &= \delta(V[\tau + \ell](\sigma(K\Theta[V(\tau)]) + \epsilon(s_t))) \\ &= \delta(V[\tau + \ell](V(\tau))). \end{aligned} \quad (2.2.7)$$

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

$$(\pi_n \circ V[\tau + k, 1])(v_\tau) = \langle V[\tau + k + 1](v), e_n \rangle \quad (2.2.8)$$

and by (??)

fix this equation, ϵ doesn't directly commute

□

IN PROGRESS: Write conjecture on decomposition which is free of neural configuration. Subject to change in later versions of ArXiv paper

Conjecture 2.2.2. *The following diagram commutes*

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

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

2. Possibly others.

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

3 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 sumultaneously if we cache!

4 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
```

1. Show that training decentralized policy gradient \implies total policy optimization
2. Show speed improvements on update step through parallelism (samples per second vs DDPG).
3. Show results are comparable with the state of the art.

5 Conclusion

We wrecked deep reinforcement learning using biological inspiration.

5.1 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.