# Persistence Length Based Exploration for Continuous Control

Riashat Islam

(Joint work with Maziar Gomrokchi, Susan Amin & Doina Precup)

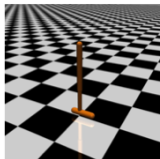Reasoning and Learning Lab



McGill University

20th April 2017

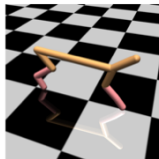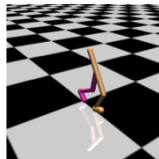# Deep Reinforcement Learning

Locomotion Tasks



Swimmer

Hopper

Half Cheetah

Walker

Ant

Simplified Humanoid

Full Humanoid
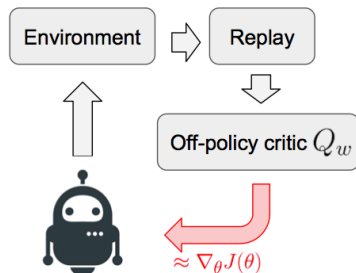
# Exploration in Continuous Control

- Exploring environment $\longleftrightarrow$ Exploiting good behaviour
- In continuous control :
  default exploration is through random control noise
- High dimensional continous actions
    - Many directed exploration methods ($\epsilon$-greedy, Boltzmann) are limited to discrete action spaces
    - Current exploration strategies are insufficient

**We propose trajectory based exploration method suited for continuous control tasks**

# Motivation

Off-Policy Actor-Critic

- ▶ DDPG in continuous control [Lillicrap et. al., 2016, Silver et. al., 2014 ]



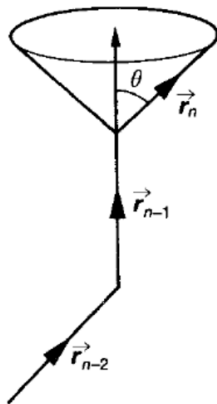However, no good exploration strategy to collect off-policy samples

- ▶ this talk : **propose exploration method for off-policy actor-critic for continuous control**
- ▶ Related current benchmark : VIME in on-policy TRPO [Houthooft et. al., 2016]

# Persistence Length Exploration

Intuition :

- Choice of next exploratory action should dependent on the trajectory so far
- Trajectories should fill up the entire state space

# Persistence Length Exploration

- Mechanism of locally self avoiding random walk
- Adopted from physics literature to describe behaviour of polymer chains
- Consider trajectory upto current state to decide next action
- Pure exploration $\rightarrow$ plan trajectory to fill up entire environment

# Persistence Length Exploration

- Self avoiding chains in d-dimensional action space
- Self avoiding trajectory
- Travel quickly around environment depending on parameterization
- Persistence length $L_p$ quantifies stiffness of the chain

# PolyRL + DDPG

**Algorithm 1:** PolyRL Algorithm (2D Action Space) on top of DDPG

1   Randomly initialize critic network Q(s,a | $\theta^Q$) and actor network $\mu(s|\theta^\mu$ with weights $\theta^Q$ and $\theta^\mu$;
2   Initialise target network $Q'$ and policy network $\mu'$;
3   Initialise two replay buffers $B^e$ and $B^d$;
4   **for** *episode=1, 2, ... M* **do**
5     PolyRL pure exploration phase → **for** *expl epoch until e = E* **do**
6       **if** *e == 0* **then**
7        Sample $\mathbf{A}_0$ and $S_0$ w.r.t $\rho$;
8       **else if** *e == 1* **then**
9        Initialize $\mathbf{H}_1$ s.t. $||\mathbf{H}_1|| = b_o$;
10        $\mathbf{A}_1 \leftarrow \mathbf{A}_o + \mathbf{H}_1$;
11       **else**
12        Draw a sample $\theta$ from $\mathcal{N}(\mu, \sigma)$;
13        $\theta_t \leftarrow$ toss a coin and choose between $\theta$ and $-\theta$;
14        $\mathbf{A}_e \leftarrow \mathbf{A}_{e-1} +$ apply $\prod_2^{\theta_e}$ on $\mathbf{H}_{e-1}$;
15       **if** *$A_e$ is not valid* **then**
16        Terminate the episode;
17       **else**
18        Apply step function on action $\mathbf{A}_e$ and observe $S_{e+1}$ and $R_{e+1}$;
19        **if** *$S_{e+1}$ is valid* **then**
20         Continue;
21        **else**
22         End the episode and re-start the chain;
23       Sample a random minibatch of transitions from buffer $B^e$;
24       Update the Q critic network using off-policy exploration samples;
25   Return trajectory of states and actions;
26   Return end of trajectory state and action;
27   Return updated Q critic network from PolyRL exploration phase;

28   Deep Deterministic Policy Gradient (DDPG);
29   **for** *t=1, 2,...T* **do**
30     Select action $a_t$ according to current policy $\mu(s_t|\theta^{mu})$;
31     Execute action $a_t$ and observe reward $r_{t+1}$ and next state $s_{t+1}$;
32     Store transition to replay buffer $B^d$;
33     Sample random minibatch from replay buffer $B^d$;
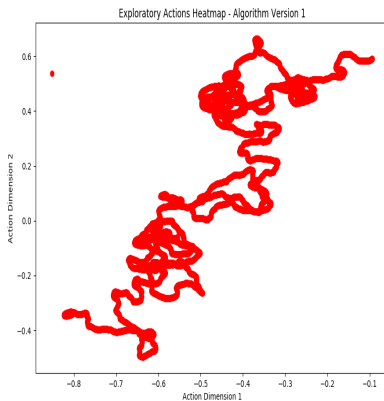34     Update the critic network by minimizing the loss;
35     Update the actor policy network using sampled policy gradient ;
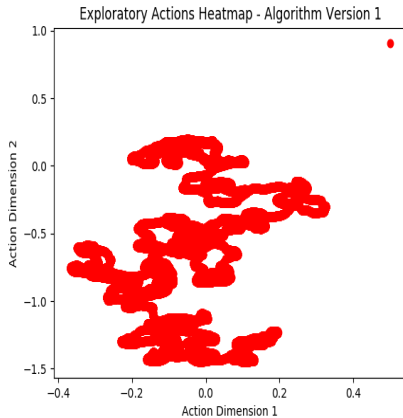36     Update the target networks;

**Persistence Length Based Exploration**

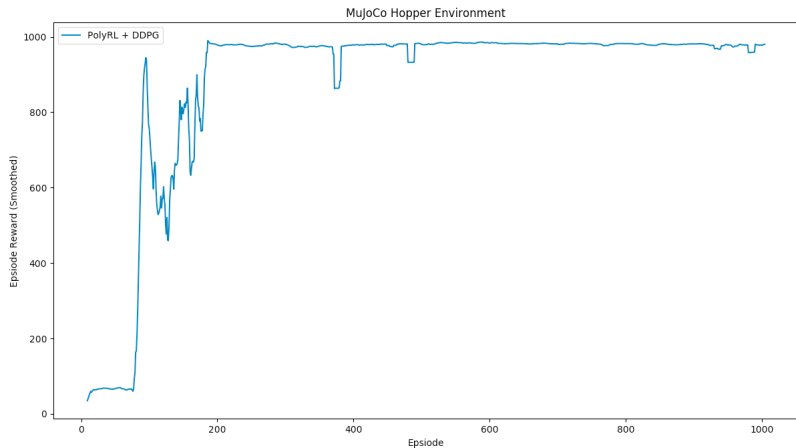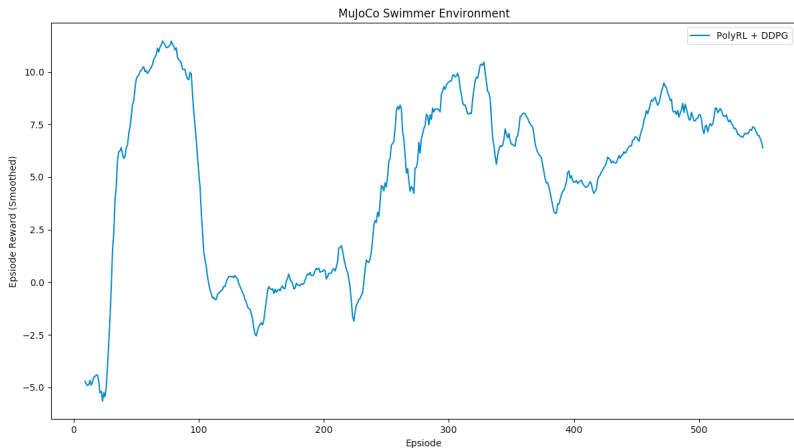**DDPG**

# PolyRL Exploration (2D Action Space)



(a) Episode 1                    (b) Episode 2

Figure: Exploratory action trajectory

# PolyRL + DDPG (MuJoCo Hopper)



MuJoCo Hopper Environment

# PolyRL + DDPG (MuJoCo Swimmer)

# Policy Gradients on MuJoCo Tasks

| Few Benchmark Results (Max Return) | | | |
|---|---|---|---|
| **Task** | **Action Dim** | **TRPO** | **DDPG** |
| Swimmer | 2D | 110 | 150 |
| Reacher | 2D | -6.7 | -6.6 |
| Hopper | 3D | 2486 | 2604 |
| HalfCheetah | 6D | 4734 | 7490 |
| Walker | 6D | 3567 | 3626 |
| Humanoid | 17D | 918 | 552 |

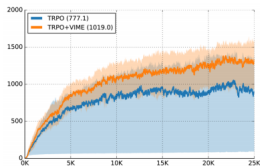# Current Benchmark - VIME

MuJoCo Walker2D, Swimmer



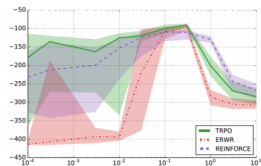Figure 3: Performance of TRPO with and without VIME on the high-dimensional Walker2D locomotion task.



Figure 4: VIME: performance over the first few iterations for TRPO, REINFORCE, and ERWR i.f.o. $\eta$ on MountainCar.
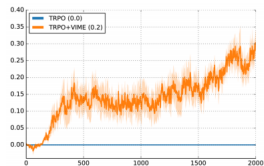


Figure 5: Performance of TRPO with and without VIME on the challenging hierarchical task SwimmerGather.

# Thank You

Questions...