

# Learning to Run challenge solutions: Adapting reinforcement learning methods for neuromusculoskeletal environments

Łukasz Kidziński, Sharada Prasanna Mohanty, Carmichael Ong, Zhewei Huang, Shuchang Zhou, Anton Pechenko, Adam Stelmaszczyk, Piotr Jarosik, Mikhail Pavlov, Sergey Kolesnikov, Sergey Plis, Zhibo Chen, Zhizheng Zhang, Jiale Chen, Jun Shi, Zhuobin Zheng, Chun Yuan, Zihui Lin, Henryk Michalewski, Piotr Miłoś, Błażej Osiński, Andrew Melnik, Malte Schilling, Helge Ritter, Sean Carroll, Jennifer Hicks, Sergey Levine, Marcel Salathé, Scott Delp

**Abstract** In the NIPS 2017 *Learning to Run* challenge, participants were tasked with building a controller for a musculoskeletal model to make it run as fast as possible through an obstacle course. Top participants were invited to describe their algorithms. In this work, we present eight solutions that used deep reinforcement learning approaches, based on algorithms such as Deep Deterministic Policy Gradient, Proximal Policy Optimization, and Trust Region Policy Optimization. Many solutions use similar relaxations and heuristics, such as reward shaping, frame skipping, discretization of the action space, symmetry, and policy blending. However, each of the eight teams implemented different modifications of the known algorithms.

## 1 Introduction

In the Learning to Run challenge participants were tasked to build a controller for a human musculoskeletal model, optimizing muscle activity such that the model travels as far as possible within 10 seconds [16]. Participants were solving a control problem with a continuous space of 41 input and 18 output parameters with high order relations between actuations and actions, simulating human musculoskeletal system. Expensive computational cost of the musculoskeletal simulations encouraged participants to develop new techniques tailored for this control problem.

All participants whose models traveled at least 15 meters in 10 seconds of the simulator time were invited to share their solutions in this manuscript. Nine teams agreed to contribute. The winning algorithm is published separately [14], while the remaining eight are collected in this manuscript. Each section in the remainder of this document describes an approach taken by one team. Sections are self-contained, they can be read independently, and each of them starts with an introduction summarizing the approach. For information on compositions of teams, affiliations and acknowledgments please refer to Section 10.

## 2 Learning to Run with Actor-Critic Ensemble

Zhewei Huang and Shuchang Zhou

We introduce an Actor-Critic Ensemble (ACE) method for improving the performance of Deep Deterministic Policy Gradient (DDPG) algorithm[19, 34]. At inference time, our method uses a critic ensemble to select the best action from proposals of multiple actors running in parallel. By having a larger candidate set, our method can avoid actions that have fatal consequences, while staying deterministic. Using ACE, we have won the 2nd place in NIPS’17 Learning to Run competition.

### 2.1 Methods

#### 2.1.1 Dooming Actions Problem of DDPG

We found that in the *Learning to Run* challenge environment legs of a fast running skeleton can easily be tripped up by obstacles. This caused the skeleton to enter an unstable state with limbs swinging and falling down after a few frames. We observed that it was almost impossible to recover from the unstable states. We call the action causing the skeleton to enter unstable state a “dooming action”.

To investigate dooming actions, we let the critic network inspect the actions at inference time. We found that most of the time, the critic could recognize dooming actions by anticipating low scores. However, as there was only one action proposed by the actor network in DDPG at every step, the dooming actions could not be avoided. This observation led us to use an actor ensemble to allow the agent to avoid dooming actions by having a critic ensemble to pick the best action from the proposed ones, as shown in Fig. 1(a).

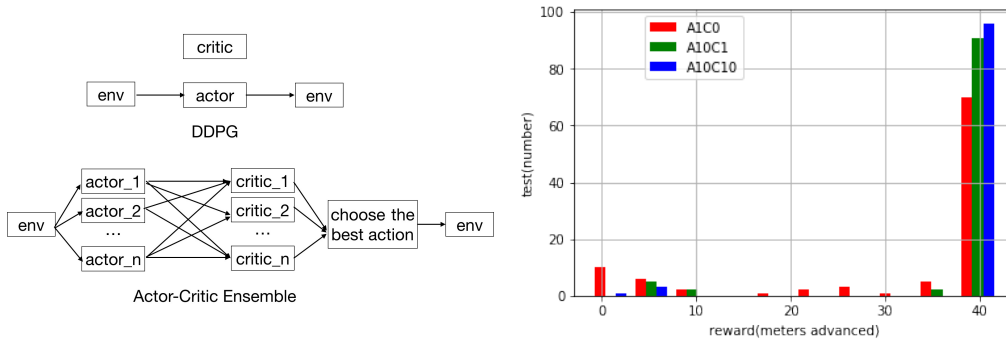


Fig. 1: Schema for DDPG and ACE

### 2.1.2 Inference-Time ACE

We first trained multiple actor-critic pairs separately, using the standard DDPG method. Then we built a new agent with many actor networks proposing actions at every step. Given multiple actions, a critic network was used to select the best action. The actor picked the action with the highest score in a greedy manner.

Empirically, we found that actors of heterogeneous nature, e.g. trained with different hyper-parameters, perform better than actors from different epochs of the same training setting. This was in agreement with the observations in the original work on Ensemble Learning [7].

To further improve critic’s prediction quality, we built an ensemble of critics, by picking the pairing critics of actors. We combined the outputs of the critic networks by averaging them.

### 2.1.3 Training with ACE

If we put actor networks together to train, all the actor networks are updated at every step, even if a certain action was not used. The modified Bellman equation takes form

$$i_{t+1} = \arg \max_j Q(s_{t+1}, \mu_j(s_{t+1}))$$

$$Q(s_t, a_t) = r(s_t, a_t) + \gamma Q(s_{t+1}, \mu_{i_{t+1}}(s_{t+1})).$$

## 2.2 Experiments and results

### 2.2.1 Baseline Implementation

We used the DDPG as our baseline. To describe the state of the agent, we collected three consecutive frames of observations from the environment. We performed feature engineering as proposed by Yongliang Qin<sup>1</sup>, enriching the observation before we feeding into the network.

As the agent was expected to run 1000 steps to finish a successful episode, we found the vanishing gradient problem (i.e. too small magnitude of the update step in the learning process) to be critical. We made several attempts to deal with this difficulty. First, we found that with the original simulation timestep, the DDPG converges slowly. In contrast, using four times larger simulation timestep, which was equivalent to changing the action only every four frames, was found to speedup convergence significantly.

We also tried unrolling DDPG as in  $TD(\lambda)$  with  $\lambda = 4$  [2], but found it inferior to simply increasing simulation timestep. Second, we tried several activation functions and found that the activation function of Scaled Exponential Linear Units(SELU)[18] is superior to ReLU, as shown in Fig. 4. SELU also outperformed Leaky ReLU, Tanh and Sigmoid.

---

<sup>1</sup> <https://github.com/ctmakro/stanford-osrl>

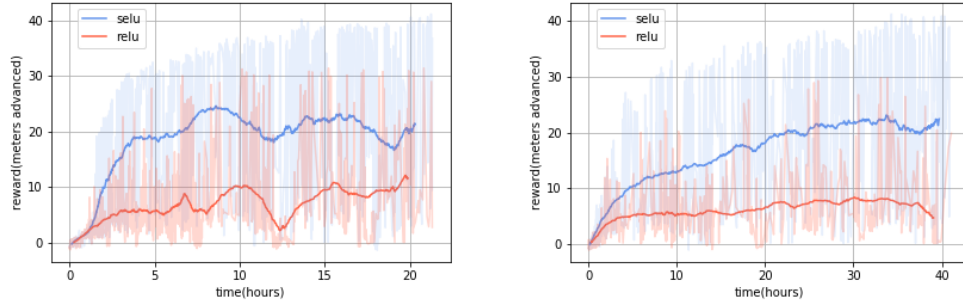


Fig. 2: Training with different activation functions and different number of processes for generating training data, by DDPG

Table 1: Hyper-parameters used in the experiments

Actor network architecture	$[FC800, FC400]$ , Tanh for output layer and SELU for other layers
Critic network architecture	$[FC800, FC400]$ , linear for output layer and SELU for other layers
Actor learning rate	$3e-4$
Critic learning rate	$3e-4$
Batch size	128
$\gamma$	0.96
replay buffer size	$2e6$

### 2.2.2 ACE experiments

For all models we used an identical architecture of actor and critic networks, with hyper-parameters listed in Table 1. Our code used for competition can be found online<sup>2</sup>.

We built the ensemble by drawing models trained with settings of the last section. Fig. 1(b) gives the distribution of rewards when using ACE, where AXCY stands for X number of actors and Y number of critics. It can be seen that A10C10 (having 10 critics and 10 actors) has a much smaller chance of falling (rewards below 30) compared to A1C0, which is equivalent to DDPG. The maximum rewards also get improved, as shown in Tab. 2.

Training with ACE was found to have similar performance as Inference-Time ACE.

Table 2: Performance of ACE

Experiment #	Test #	Actor #	Critic #	Average reward	Max reward	# Fall off
A1C0	100	1	0	32.0789	41.4203	25
A10C1	100	10	1	37.7578	41.4445	7
A10C10	100	10	10	39.2579	41.9507	4

<sup>2</sup> <https://github.com/hzwer/NIPS2017-LearningToRun>

## 2.3 Discussion

We propose Actor-Critic Ensemble, a deterministic method that avoids dooming actions at inference time by asking an ensemble of critics to pick actions proposed by an ensemble of actors. Our experiments found that ACE can significantly improve performance of DDPG, by reducing of the number of falls and increasing the speed of running skeletons.

## 3 Deep Deterministic Policy Gradient and improvements

Mikhail Pavlov, Sergey Kolesnikov and Sergey Plis

We benchmarked state of the art policy-gradient methods and found that Deep Deterministic Policy Gradient (DDPG) method is the most efficient method for this environment. We also applied several improvements to DDPG method, such as layer normalization, parameter noise, action and state reflecting. All this improvements helped to stabilize training and improve its sample-efficiency.

### 3.1 Methods

#### 3.1.1 DDPG improvements

We used standard reinforcement learning techniques: action repeat (the agent selects action every 5th state and selected action is repeated on skipped steps) and reward scaling. After several attempts, we choose a scale factor of 10 (i.e. multiply reward by ten) for remaining experiments. For exploration we used Ornstein-Uhlenbeck (OU) process [37] to generate temporally correlated noise, considered efficient in exploration of physical environments. Our DDPG implementation was parallelized as follows:  $n$  processes collected samples with fixed weights all of which were processed by the learning process at the end of an episode, which updated their weights. Since DDPG is an off-policy method, the stale weights of the samples only improved the performance providing each sampling process with its own weights and thus improving exploration.

#### 3.1.2 Parameter noise

Another improvement is the recently proposed parameters noise [25] that perturbs network weights encouraging state dependent exploration. We used parameter noise only for the actor network. Standard deviation  $\sigma$  for the Gaussian noise was chosen according to the original work [25] so that measure where  $\tilde{\pi}$  is the policy with noise, equals to  $\sigma$  in OU. For each training episode we switched between the action noise and the parameter noise choosing them with 0.7 and 0.3 probability respectively.

### 3.1.3 Layer norm

Henderson et al. showed that layer normalization [3] stabilizes the learning process in a wide range of reward scaling. We have investigated this claim in our settings. Additionally, layer normalization allowed us to use same perturbation scale across all layers despite the use of parameters noise [25]. We normalized the output of each layer except the last for critic and actor by standardizing the activations of each sample. We applied layer normalization before the nonlinearity.

### 3.1.4 Actions and states reflection symmetry

The musculoskeletal model to control in the challenge has bilateral body symmetry. State components and actions can be reflected to increase sample size by factor of 2. We sampled transitions from replay memory, reflected states and actions and used original states and actions as well as reflected as batch in training step. This procedure improves stability of learned policy. When we did not use this technique our model learned suboptimal policies, when for example muscles for only one leg are active and other leg just follows the first leg.

## 3.2 Experiments and results

For all experiments we used environment with 3 obstacles and random strengths of the psoas muscles. We tested models on setups running 8 and 20 threads. For comparing different PPO, TRPO and DDPG settings we used 20 threads per model configuration. We have compared various combinations of improvements of DDPG in two identical settings that only differed in the number of threads used per configuration: 8 and 20. The goal was to determine whether the model rankings are consistent when the number of threads changes. For  $n$  threads (where  $n$  is either 8 or 20) we used  $n - 2$  threads for sampling transitions, 1 thread for training, and 1 thread for testing. For all models we used identical architecture of actor and critic networks. All hyperparameters are listed in Table 3. Our code used for competition and described experiments can be found in a github repo.<sup>3</sup> Experimental evaluation is based on the non-discounted return.

### 3.2.1 Benchmarking different models

Comparison of our winning model with the baseline approaches is presented in Figure 3. Among all methods the DDPG significantly outperformed PPO and TRPO. The environment is time expensive and method should utilized experience as effectively as possible. DDPG due to experience replay (re)uses each sample from environment many times making it the most effective method for this environment.

---

<sup>3</sup> Theano: [https://github.com/fgvbrt/nips\\_rl](https://github.com/fgvbrt/nips_rl) and PyTorch: <https://github.com/Scitator/Run-Skeleton-Run>

Table 3: Hyperparameters used in the experiments.

parameters	Value
Actor network architecture	[64, 64], elu activation
Critic network architecture	[64, 32], tanh activation
Actor learning rate	linear decay from $1e-3$ to $5e-5$ in $10e6$ steps with Adam optimizer
Critic learning rate	linear decay from $2e-3$ to $5e-5$ in $10e6$ steps with Adam optimizer
Batch size	200
$\gamma$	0.9
replay buffer size	$5e6$
rewards scaling	10
parameter noise probability	0.3
OU exploration parameters	$\theta = 0.1, \mu = 0, \sigma = 0.2, \sigma_{min} = 0.05, dt = 1e-2, n_{steps}$ annealing $\sigma_{decay} 1e6$ per thread

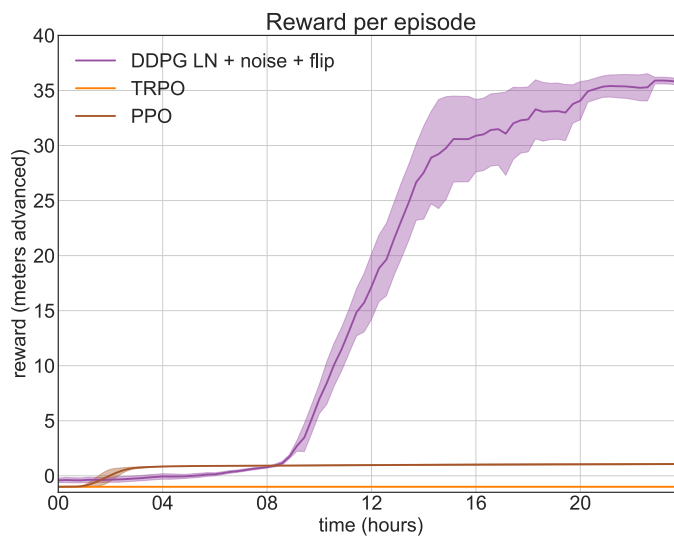


Fig. 3: Comparing test reward of the baseline models and the best performing model that we have used in the “Learning to run” competition.

### 3.2.2 Testing improvements of DDPG

To evaluate each component we used an ablation study as it was done in the rainbow article [12]. In each ablation, we removed one component from the full combination. Results of experiments are presented in Figure 4a and Figure 4b for 8 and 20 threads respectively. The figures demonstrate that each modification leads to a statistically significant performance increase. The model containing all modifications scores the highest reward. Note, the substantially lower reward in the case, when parameter noise was employed without the layer norm. One of the reasons is our use of the same perturbation scale across all layers, which does not work that well without normalization. Also note, the behavior is quite stable across number of threads, as well as the model ranking. As expected, increasing the number of threads improves the result.

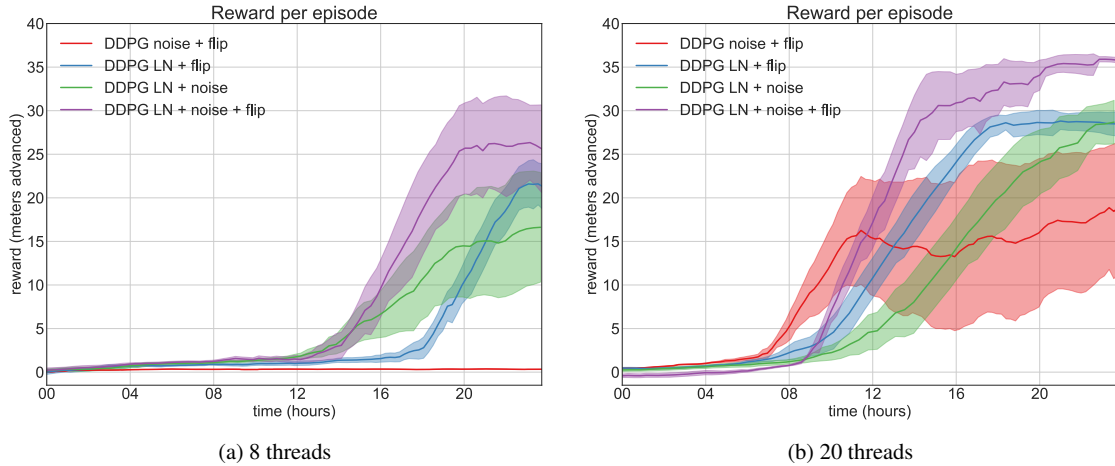


Fig. 4: Comparing test reward for various modifications of the DDPG algorithm with 8 threads per configuration (Figure 4a) and 20 threads per configuration (Figure 4b). Although the number of threads significantly affects performance, the model ranking approximately stays the same.

Maximal rewards achieved in the given time for 8 and 20 threads cases for each of the combinations of the modifications is summarized in Table 4. The main things to observe is a substantial improvement effect of the number of threads, and stability in the best and worst model rankings, although the models in the middle are ready to trade places.

Table 4: Best achieved reward for each DDPG modification.

agent \ # threads	8	20
	DDPG + noise + flip	0.39
DDPG + LN + flip	25.29	31.91
DDPG + LN + noise	25.57	30.90
DDPG + LN + noise + flip	<b>31.25</b>	<b>38.46</b>

### 3.3 Discussion

Our results in OpenSim experiments indicate that in a computationally expensive stochastic environments that have high-dimensional continuous action space the best performing method is off-policy DDPG. We have tested 3 modifications to DDPG and each turned out to be important for learning. Action states reflection doubles the size of the training data and improves stability of learning and encourages the agent to learn to use left and right muscles equally well. With this approach the agent truly learns to run. Examples of the learned policies with and without the reflection are present at this URL <https://tinyurl.com/ycvfq8cv>.



Parameter and Layer noise additionally improves stability of learning due to introduction of state dependent exploration.

## 4 Asynchronous DDPG with Deep Residual Network for Learning to Run

Zhibo Chen, Zhizheng Zhang, Jiale Chen and Jun Shi

For improving the training effectiveness of DDPG on this physics-based simulation environment which has high computational complexity, we designed a parallel architecture with deep residual network for the asynchronous training of DDPG. In this work, we describe our approach and we introduce supporting implementation details.

### 4.1 Methods

#### 4.1.1 Asynchronous DDPG

In our framework, the agent could collect interactive experiences and update its network parameters asynchronously. For the collection of experiences, the *Learning to Run* environments with different seeds and same difficulty-level settings were wrapped by multi-process programming. All step-by-step interactive experiences in every wrapped environment would be stored in a specific storage until this episode finished. Then we decided which step experiences to put into the experience replay memory according to their corresponding step rewards and episode rewards. In terms of the updating of networks' parameters, the updating process would sample a batch from the replay memory after each interaction with the RL environments no matter which specific environment process this interaction takes place in.

#### 4.1.2 The Neural Network Structure

Whether for the human body in real-world or the musculoskeletal model used in this simulation, the accurate physical motions are determined by multiple joints and implemented by the excitations of multiple different muscles. Taking it naturally, we applied 1D convolution modules in the neural networks for both actor and critic networks with the expectation of capturing the correlation among 41 values of the observation. And our experimental results indicated that 1D convolution neural networks were better able to prevent converging to the local optimal solution than fully connected networks. In order to improve the efficiency and stability of training, we added the residual blocks (see Figure 5) to make our model easier to train and converge.

We also tried to take advantage of 2D convolution to process the 1D observation information and learn the features from historical actions, inspired by the work on 3D convolution neural network for human action recognition [15]. However, the performance of the RL agent with 2D convolution was less likely to converge steadily.

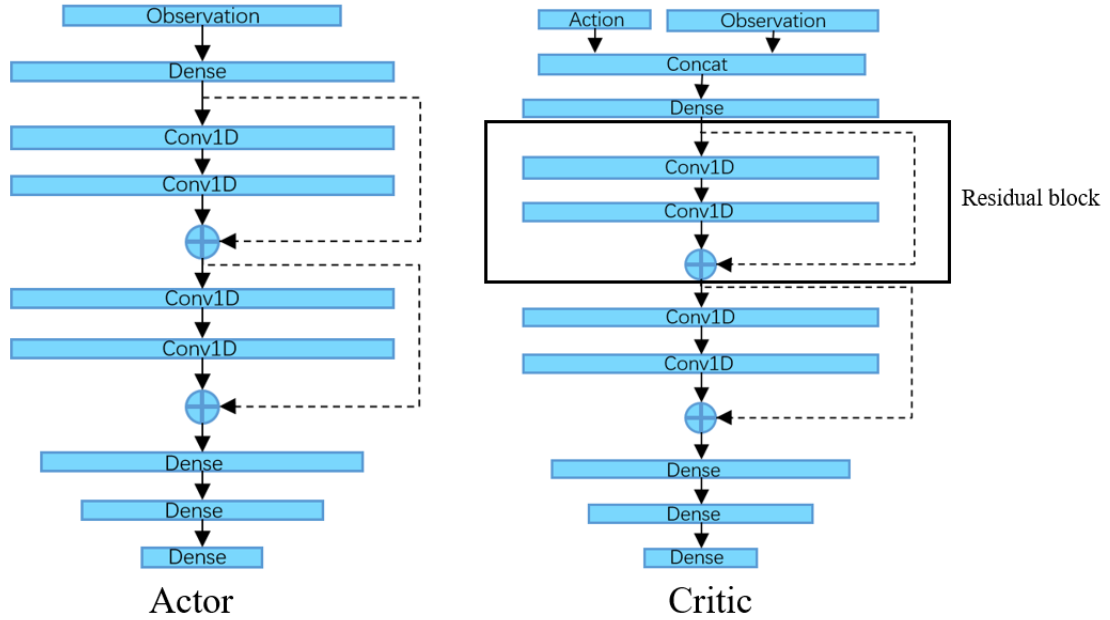


Fig. 5: The diagram of our network structure with residual blocks

#### 4.1.3 Noise for Exploration

We tried both parameter space noise [25] and action space noise for the exploration in this task. For parameter space noise, it was really hard to fine-tuning and get an optimal solution in *Learning to Run* environment, which might be caused by the structure of our neural network. In terms of action space noise, we found that Ornstein-Uhlenbeck noise would lead to inefficient exploration and convergence to local optimal solution. Instead, correlated Gaussian noise was more suitable in this task. Additionally, considering there should be the continuity among actions as the outputs of the actor, we designed a so-called random walk Gaussian noise for this continuous task, which brought us the highest grade but with a little large variance. Hence, we thought that normal correlated Gaussian noise and the random walk Gaussian noise were both effective for exploration in this environment, but each had its advantages and disadvantages.

#### 4.1.4 Detailed Implementation

We would like to discuss the stability in this section, which included the stability of training and the stability of the policy for obstacles crossing. For the stability of training, we applied some common techniques in our model, such as layer normalization, reward scaling, prioritized experience replay [31] and action repetition. Additionally, we applied a training trick with a small learning rate, named “trot”. In detail, we sampled one batch and used it to implement back-propagation for multiple times with a small learning rate. For the stability of the policy for obstacles crossing, a ceil option for the radius of the obstacle turned out to be significant to improve the agent’s performance. The obstacles for the agent would be slightly larger than their real sizes and

the mathematic space to be fitted by neural network will be also reduced by this method. Hence, we could make it easier for the learning by neural network and improve the stability of obstacles crossing meanwhile.

## 4.2 Experiments and results

### 4.2.1 The Number of Parallel Process

Based on the experimental verification, we found the number of environment process would have a large impact on the learning performance of the agent. We tested our model with setting 12, 24, and 64 processes. The experimental results indicated that the more processes we used, the more samples we could get but not inevitably the better for the learning performance. In our settings, the model with 24 environment process would get the highest grade. Excessive number of the parallel environment processes could cause that too much similar transitions were pushed into the experience replay memory, which adversely affected the training of agent.

### 4.2.2 The Neural Network Structure

The neural network structure seriously affects the learning ability of the agent. According to our training results (see Figure 6), the 1D convolution neural network with residual blocks was the best fit for both actor and critic in DDPG whether for the maximum learning ability or the stability. Moreover, the 1D convolution neural network without residual blocks was also better than the fully connection network. Moreover, widening the network was more effective than deepening the network in this task, because it was easier to get the gradient in a suitable range for wide networks to implement policy updating.

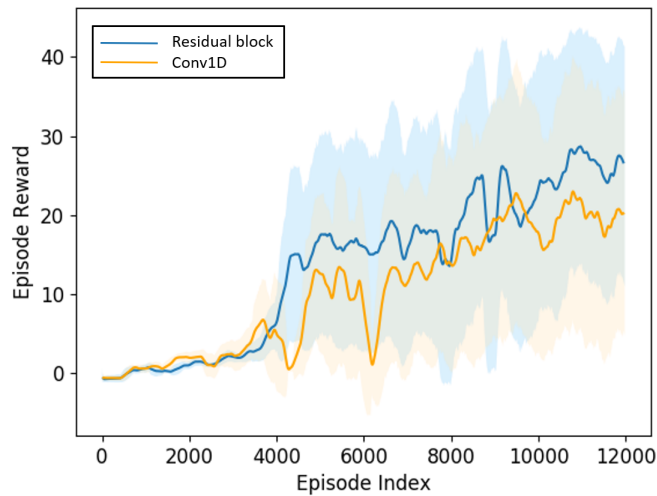


Fig. 6: The results of the comparison of network structure with and without residual module

### 4.3 Discussion

In this section, we discuss the difference between the first and the second round of the competition. Due to the setting of that the current episode will be finished if you fall down, the requirements on stability of obstacles crossing were significantly higher than it in the first round. More precisely, the obstacles only appeared in the beginning of the journey, it was easier for an agent to cross the obstacles in the set time, because the agent had not accelerated yet to a great speed. In the second round, the agent should cross the obstacles when it runs with a high speed, which made it easy to fall down. As our previous description, a ceiling constraint for the radius of obstacles could solve this problem effectively to some extent.

Interestingly, our agent discovered a really smart trick by itself in the training process. Instead of avoiding the obstacles, it would try to deliberately step on a large obstacle and then used it as a stepping stone. Although it didn't master this trick due to the limitation of training time, we were really surprised by this exploitation of the environment.

## 5 Proximal Policy Optimization with Policy Blending

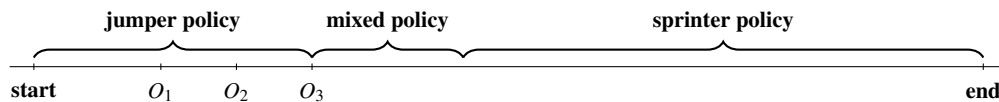
Henryk Michalewski, Piotr Miłoś and Błażej Osipiński

Our solution was based on the distributed Proximal Policy Optimization algorithm combined with a few efficiency-improving techniques. We used the *frameskip* to increase exploration. We changed rewards to encourage the agent to *bend its knees*, which significantly stabilized the gait and accelerated the training. In the final stage, we found it beneficial to transfer skills from small networks (easier to train) to bigger ones (with more expressive power). For this purpose we developed *policy blending*, a general cloning/transferring technique.

### 5.1 Methods

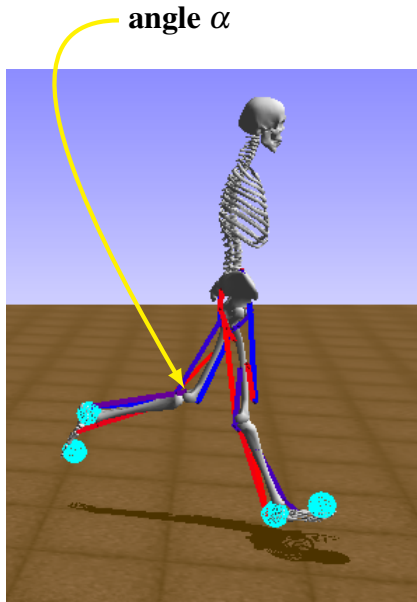
Combining cautious and aggressive strategies.

In the first stage of the competition our most successful submission was a combination of 2 agents. The first agent was cautious and rather slow. It was designed to steadily jump over the three obstacles (approx. first 200 steps of an episode). For the remaining 800 steps we switched to a 20%-faster agent, trained beforehand in an environment without obstacles.



The switching of policies was a rather delicate task. Our attempts to immediately change from the first policy left the agent in a state unknown to the second one and caused the agent to fall. Eventually, we switched the policies gradually by applying the linear combination  $(1 - \frac{k}{n})a_\eta + \frac{k}{n}a_\nu$ , where  $k$  is the transition step,

$a_\eta, a_v$  are actions of the jumper and sprinter respectively;  $n = 150$  was required to smooth the transition. A more refined version of this experiment should include learning of a macro policy which would decide on its own how the jumper and sprinter should be combined (see [36, 38] for a broader introduction to hierarchical learning).



#### Frameskip.

Our initial experiments led to suboptimal behaviors, such as kangaroo-like jumping (two legs moving together). We conjectured that the reason was a poor exploration and thus applied frameskip. In this way we obtained our first truly bipedal walkers. In particular, frameskip set to 4 led to a slow but steadily walking agent capable of scoring approximately 20 points. Reward shaping.

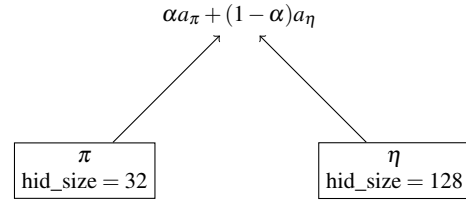
In the process of training we obtained various agents walking or running, but still never bending one of the legs. We decided to *shape agent's behavior* [8] through adding of an extra reward for bending of knees. More specifically, we added a reward for getting the angle  $\alpha$  into a prescribed interval, see the figure on the left. This adjustment resulted in significant improvements in the walking style. Moreover, training with the knee reward caused our policy to train significantly faster, see Figure 8. We experimented with various other rewards. Perhaps the most interesting was explicit penalization of falling over. As a consequence we created an ultra-cautious agent, who always concluded the whole episode but at the cost of a big drop of the average speed.

#### Final tuning - policy blending.

We found it easier to train a reasonable bipedal walking policy  $\pi$  when using small nets. However, small nets suffered from unsatisfactory final performance in environment with many obstacles. We used pretrained  $\pi$  and a method we called *policy blending* to softly guide the learning process. Namely, we kept  $\pi$  fixed and trained new  $\eta$ ; the agent was issuing actions  $\alpha a_\pi + (1 - \alpha)a_\eta$ ,  $\alpha \in (0, 1)$ . One can see blending as a simplified version of progressive networks considered in [27].

Even with  $\alpha \approx 0.1$  the walking style of  $\pi$  was coarsely preserved, while the input from bigger net of  $\eta$  led to significant improvements in harder environments. In some cases blended policies could successfully deal with obstacles even though  $\pi$  was trained in an obstacle-free environment. In some experiments after an initial period of training with  $\alpha > 0$ , we continued with  $\alpha = 0$ , which can be seen as knowledge transfer from  $\pi$  to  $\eta$ . In our experience, such knowledge transfer worked better than direct behavioral cloning.

To test cloning we followed a procedure of [4], first copying the original policy to a new neural net through supervised learning followed by further training of the new policy on its own. However, in our experiments policies obtained through cloning showed at best the same performance as the original policies. Conversely, as can be seen in Figure 9 transferring a policy from obstacle-free environment using policy blending performed better than simple retraining.

Fig. 7: Blending of  $\pi$  and  $\eta$ .

## 5.2 Experiments and results

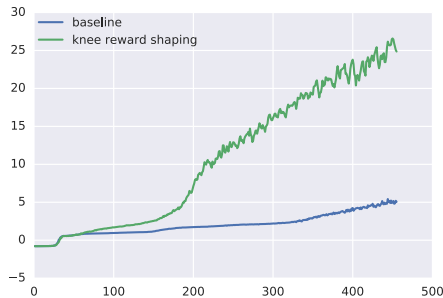


Fig. 8: Reward shaping

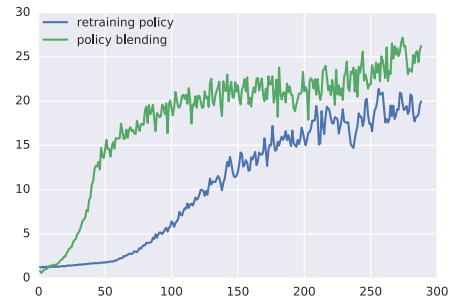


Fig. 9: Policy blending

Overall, we performed approximately 5000 experiments of lengths up to 72 hours. We used the PPO optimization algorithm [32]. In Figure 8 we present how the knee reward helps in the training. Figure 9 we compare retraining with blending of policies.

## 6 Double Bootstrapped DDPG for Continuous Deep Reinforcement Learning

Zhuobin Zheng, Chun Yuan and Zhihui Lin

Deep Deterministic Policy Gradient (DDPG) provides substantial gains in sample efficiency on off-policy experience data over on-policy methods. However, vanilla DDPG may explore inefficiently and be easily trapped into local optima in the Learning to Run Challenge. We proposed *Double Bootstrapped DDPG*, an algorithm that combines efficient exploration with promising stability in continuous control tasks.

## 6.1 Methods

### 6.1.1 Double Bootstrapped DDPG

Inspired by Bootstrapped DQN [22], Double Bootstrapped DDPG, abbreviated with DB-DDPG, extends the actor-critic architecture to completely bootstrapped networks (see Figure 10 for an overview of the approach). Both actor and critic networks have a shared body for feature extraction, followed by multiple heads with different random initialization.

A simple warm-up is applied before training: the actor heads are randomly selected to interact with the environment and pre-train during every episode, together with their paired critic heads as vanilla DDPGs.

During one single training episode, the  $k_{th}$  pair of actor and critic heads are randomly activated to learn. Given a state  $s_t$ , multiple actor heads output candidate actions  $\mathbf{a}_t = (a^1, a^2, \dots, a^K)_t$ , which are concatenated to the critic network. Multiple critic heads output an  $E$ - $Q$  value matrix ( $\mathbb{R}^{K \times K}$ ) for the actions  $\mathbf{a}_t$  according to the state  $s_t$ . The final ensemble action with highest  $Q$ -value sum which is determined by the  $E$ - $Q$  value matrix as Equation (1),

$$a_t = \arg \max_a \left\{ \sum_{i=1}^K Q_i(s_t, a) \mid_{a=\mu_k(s_t)} \right\}_{k=1}^K, \quad (1)$$

is chosen to execute for the state receiving a new state  $s_{t+1}$  and a reward  $r_t$ .

Moreover, a random mask  $\mathbf{m}_t = (m^1, m^2, \dots, m^K)_t$  is generated with Bernoulli distribution simultaneously. A new transition  $(s_t, a_t, r_t, s_{t+1}, \mathbf{m}_t)$  is stored in experience memory. The selected  $k_{th}$  heads together with respective shared bodies and target networks are trained as a DDPG given a minibatch of samples. The  $i_{th}$  experience with mask  $m_i^k = 0$  is ignored when training for bootstrapping [22]. A more detailed procedure can be found in Algorithm 1.

Besides, we also used common reinforcement learning techniques, such as: frame skipping [20](the agent performs the same action every  $k$  consecutive observations. we set  $k = 4$ ), prioritized experience replay [31] and a reward trick [9] which utilizes velocity instead of distance as reward encouraging the agent to make more forward process along the track.

### 6.1.2 Observation Preprocessing

Since the original observation vector given by the environment does not satisfy the Markov property, we extended the observation by calculating more velocities and acceleration of the remaining body parts. During the experiments, we found that obstacles may be extremely close to each other. In this case, when an agent overstepped the first obstacle with one single leg, a new observation about the next obstacle was immediately updated. As a result, it probably fell down since another leg hit the previously invisible obstacle. To solve this problem, information about the previous obstacle was appended to the observation vector.

### 6.1.3 Noise Schedule

Ornstein-Uhlenbeck (OU) process [37] was used to generate noise for exploration. DB-DDPG utilized a noise decay rate for action noise for balancing exploration and exploitation. Once the noise decreased below a threshold, DB-DDPG acted without noise injection and focused on exploiting the environment. Exploration

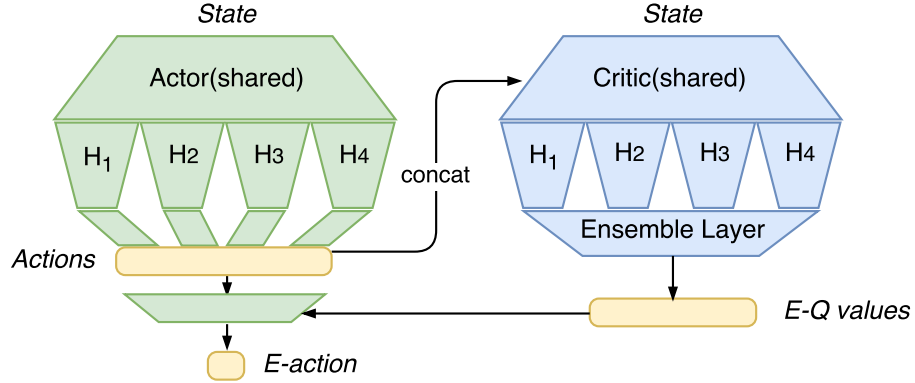


Fig. 10: Structure of Double Bootstrapped DDPG. When the actor(green) receives a state, each head outputs an action vector( $\mathbb{R}^{18}$ ). Given the same state, the critic(blue) concatenates these actions( $\mathbb{R}^{K \times 18}$ ) in the hidden layer and outputs an *Ensemble-Q* value matrix( $\mathbb{R}^{K \times K}$ ) by  $K$  heads. The actor chooses the final *Ensemble-action* determined by the *Ensemble-Q* values.

noise recovered when the uncertainty of the agent decreased. This process could be switched iteratively and fine-tuned manually for stability until convergence.

## 6.2 Experiments

### 6.2.1 Details

We used Adam [17] for learning the parameters with a learning rate of  $1e^{-4}$  and  $3e^{-4}$  for the actor and the critic network respectively. For the  $Q$  networks we set a discount factor of  $\gamma = 0.99$  and  $\tau = 1e^{-3}$  for soft target updates. All hidden layers utilized exponential linear units(ELU) [5] while the output layer of the actor utilized a  $\tanh$  layer to bound the actions followed by scale and shift operations. The shared network of the actor had two hidden layers with 128 and 64 units respectively. This was followed by multiple heads with 64 and 32 units. The critic had similar architecture except actions which were concatenated in the second hidden layer.

### 6.2.2 Results

We evaluated the algorithm on a single machine. Results in Figure 11 showed that multi-head architecture brought significant performance increase. These models with different multiple heads were pre-trained using previous experience memory in the warm-up phase before training. The single-head model could be regarded as vanilla DDPG. These agents were trained using the same techniques and hyperparameters mentioned above following Algorithm 1. The figure showed that DB-DDPGs outperformed vanilla DDPG on performance ensuring faster training and stability. The model with 10-head scored the highest cumulative reward by more efficient exploration and ensemble training.



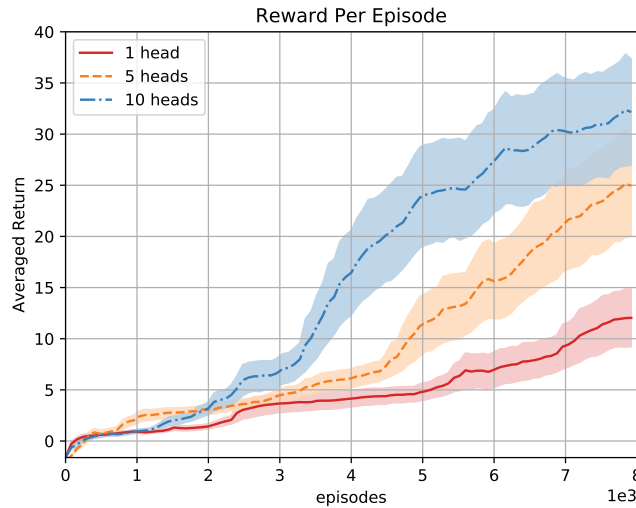


Fig. 11: Episode reward comparison among three modifications of DB-DDPG with different number of bootstrapped heads. Single-head model could be regarded as vanilla DDPG. Multi-head models scored higher rewards than single-head model.

### 6.3 Discussion

We proposed the *Double Bootstrapped DDPG* algorithm for high-dimensional continuous control tasks. It was demonstrated that DB-DDPG with multiple heads led to significantly faster learning than vanilla DDPG while retaining stability in the Learning to Run challenge.

Going forward, DB-DDPG can be not only efficient on a single machine, but parallelizable up to more machines boosting the learning. Furthermore, we believe this architecture is also available for dealing with multi-task reinforcement learning problems, which means each head only concentrates on its own sub-task when training then makes its own professional decision for ensemble collaboratively.

## 7 Plain DDPG with soft target networks

Anton Pechenko

We used DDPG with soft target networks and Ornstein-Uhlenbeck process. Discount factor was 0.99, replay buffer was 1 000 000, batch size was 64. For training we used 7 simulators which was run in parallel.

**Algorithm 1** Double Bootstrapped DDPG

---

**Input:** head number  $K$ , mini-batch size  $N$ , maximum training episode  $E$   
**Initialize:** Randomly initialize critic network  $Q(s, a | \theta^Q)$  with  $K$  outputs  $\{Q_k\}_{k=1}^K$ , and actor network  $\mu(s | \theta^\mu)$  with  $K$  outputs  $\{\mu_k\}_{k=1}^K$ .  
Initialize critic target networks  $\theta_1^{Q'}, \dots, \theta_K^{Q'}$  with weights  $\theta_k^{Q'} \leftarrow \theta_k^Q$  and actor target networks  $\theta_1^{\mu'}, \dots, \theta_K^{\mu'}$  with weights  $\theta_k^{\mu'} \leftarrow \theta_k^\mu$   
Initialize replay buffer  $R$ , masking distribution  $M$   
**for** episode  $e = 1, E$  **do**  
  Initialize a random process  $\mathcal{N}$  for action exploration  
  Receive initial observation state  $s_0$   
  Pick a pair of activated critic and actor networks using  $k \sim \text{Uniform}\{1, \dots, K\}$   
  **for** step  $t = 1, T$  **do**  
    Select action  $a_t$  from actions  $\{a | a^k = \mu(s_t | \theta_k^\mu)\}_{k=1}^K$  according to the policies and exploration noise as following,  

$$a_t = \arg \max_a \left\{ \sum_{i=1}^K Q_i(s_t, a) |_{a=\mu_k(s_t)} \right\}_{k=1}^K + \mathcal{N}_t$$
  
    Execute action  $a_t$  then observe reward  $r_t$  and new state  $s_{t+1}$   
    Sample bootstrapped mask  $m_t \sim M$   
    Store transition  $(s_t, a_t, r_t, s_{t+1}, m_t)$  in  $R$   
    Sample a random minibatch of  $m$  transitions  $(s_i, a_i, r_i, s_{i+1}, m_i)$   
    Set  $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1} | \theta_k^{\mu'})) | \theta_k^{Q'}$   
    Update critic network by minimizing the loss:  $L = \frac{1}{N} \sum_i m_i^k (y_i - Q(s_i, a_i | \theta_k^Q))^2$   
    Update actor network using the sampled policy gradient:  

$$\nabla_{\theta_k^\mu} \approx \frac{1}{N} \sum_i \nabla_a Q(s, a | \theta_k^Q) |_{s=s_i, a=\mu(s_i)} \nabla_{\theta_k^\mu} \mu(s | \theta_k^\mu) |_{s=s_i}$$
  
    Update the target networks:  

$$\theta_k^{Q'} \leftarrow \tau \theta_k^{Q'} + (1 - \tau) \theta_k^Q$$
  

$$\theta_k^{\mu'} \leftarrow \tau \theta_k^{\mu'} + (1 - \tau) \theta_k^\mu$$
  
  **end for**  
**end for**

---

## 7.1 Method

### 7.1.1 State description

To remove redundancy from the state space and make training more efficient it is often reasonable to exploit state space symmetry. That is why we used *first person view* transformation of the observation vector. That means we subtracted coordinates and angle of pelvis from others bones coordinates and angles. Also we assigned zero to  $X$  coordinate of the pelvis to collapse observation space along  $X$  coordinate to make run performance independent of  $X$  coordinate. We had no information about ground reaction forces in observation vector. Since it seemed very important for running it was needed to estimate it. To that end, we constructed a state vector from a sequence of the three last observation vectors.

### 7.1.2 Training process

We used two separated multilayer perceptrons [26] for actor and critic with 5 hidden layers, 512 neurons each. In three out of seven agents we random 20-40 actions at start and then started collecting  $(S, A, R, S')$  experience items, i.e. a previous state state, an action, a reward and a result state. This increases variety of starting points and increases quality of replay buffer experience. For 30% of simulations we turned off obstacles in order to increase the variety of replay buffer. Training consisted of two steps. The first step was optimization with Adam [17] with  $1e-4$  learning rate till agents increase its score at maximum, which, in our case, turned out to be approximately 37 meters. In the second step, we trained the network with the stochastic gradient descent (SGD) using the learning rate of  $5e-5$ . During the SGD step agents decreased running velocity but increase robustness to falls and obstacles overcoming resulting in 26 meters score.

## 7.2 Experiments and results

The implementation of this solution can be found as a part of the RL-Server software available at github<sup>4</sup>. The RL-Server is a python application using Tensorflow [1], with DQN [20] and DDPG algorithms included. The entire training process, including replay buffer, training and inference steps is performed within the RL-Server application. A lightweight client code can be included in parallel running environments. It enables multiple languages thanks to an open communication protocol.

## 8 PPO with reward shaping

Adam Stelmaszczyk and Piotr Jarosik

We trained our final model with PPO on 80 cores in 5 days using reward shaping, extended and normalized observation vector. We recompiled OpenSim with lower accuracy to have about 3x faster simulations.

### 8.1 Methods

Our general approach consisted of two phases: exploration of popular algorithms and exploitation of the most promising one. In the first phase we tried 3 algorithms (in the order we tried them): Deep Deterministic Policy Gradient (DDPG) (keras-rl implementation [24]), Proximal Policy Optimization (PPO) [33] (OpenAI baselines implementation [6]), Evolution Strategies (ES) [28] (OpenAI implementation [29]). We also tried 3 improvements: changing the reward function during training (reward shaping), improving observations (feature engineering) and normalizing observations.

We started our experiments with DDPG (without improvements) and we could not achieve good results. That was probably because of the bad normalization or not enough episodes. We had problems parallelizing keras-rl and we were using only one process. Therefore, we switched to PPO and ES, for which learning plots and parallelization looked better. We were incrementally adding improvements to these two techniques. In the

---

<sup>4</sup> <https://github.com/parilo/rl-server>

end we used the default hyperparameters for all the algorithms, their values can be found in the full source code [35].

### 8.1.1 Reward shaping

We guided learning to promising areas by shaping the reward function. We employed: a penalty for straight legs, a penalty for `pelvis.x` greater than `head.x` (causing a skeleton to lean forward), adding 0.01 reward for every time step (to help take the first step and get out of local maximum) and using velocity instead of distance passed, found in [9]. Using velocity rewards passing the same distance in less time steps.

### 8.1.2 Feature engineering

We changed all `x` positions from absolute to relative to `pelvis.x`. Thanks to that similar poses were represented by similar observations. We also extended the observation vector from 41 values to 82 by adding: the remaining velocities and accelerations of the body parts, ground touch indicator for toes and talus, a queue of two obstacles: the next one (preserved from the original observation) and the previous one. Without this, when passing over an obstacle, agent would lose sight of the obstacle underneath it as it would immediately switch to the next one.

### 8.1.3 Normalizing observations

We logged the mean and standard deviation of all the observations to see if the normalization was done correctly. By default, baselines PPO implementation used a filter which automatically normalized every observation. For every observation, it was keeping its running mean and standard deviation. It did normalization by subtracting the mean and dividing by std. This worked well for most of the observations, however for some it was problematic. For constants, e.g. the strength of `psoas`, the standard deviation would be 0. The code in that case was just passing this observation as it is. The magnitude of an observation was treated as an importance when passed to a network. Too big values would saturate all the other smaller inputs (which may be more important). Also, the first strength of `psoas` had a different value than the following ones (due to a bug in the challenge environment, later fixed). So, the filter would calculate some arbitrary mean with standard deviation and later use them. Another problem was that some observations were most of the time were close to 0, but were shooting up in some moments to greater values, e.g. velocity. This resulted in huge values (because initial standard deviation was close to 0), saturating the network.

Because all of these problems, we skipped the auto normalizing and manually normalized every observation. Iteratively, we were running our model and visualizing mean with standard deviation for all the observations. Then we were correcting the normalization of observations which mean was far from 0 or standard deviation far from 1.

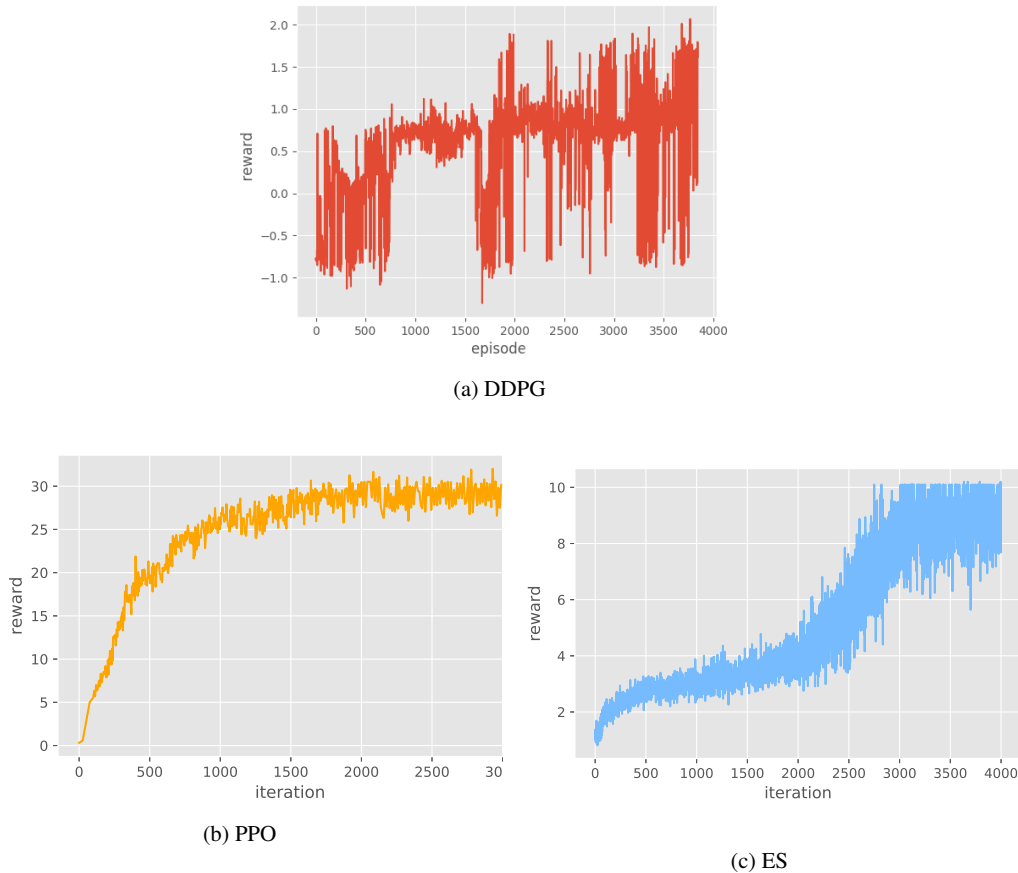


Fig. 12: Mean rewards (after reward shaping, so higher than during grading) achieved by DDPG, PPO and ES in one training run. The x axis in the DDPG plot stands for episodes, in PPO and ES - for iterations (multiple episodes), because of the different implementations used (keras-rl for DDPG and OpenAI baselines for PPO & ES). The mean reward of DDPG was very variable. We experienced the most stable learning with PPO. ES gave us worse results, after 3000 iterations the skeleton stayed still for 1000 time steps, scoring around 10, because we rewarded 0.01 for every survived time step.

## 8.2 Experiments and results

We found that OpenSim simulator was the bottleneck, accounting for about 99% of the CPU time. To speed it up about 3x, we changed the simulator's accuracy from default 0.1% to 3%, following <https://github.com/ctmakro/stanford-osrl#the-simulation-is-too-slow>. Any change made in our environment could have introduced a bias when grading on the server, however we didn't notice any significant score changes.

We conducted our experiments on: *Chuck*, an instance with 80 CPU cores (Intel Xeon), provided by Faculty of Mathematics, Informatics, and Mechanics, University of Warsaw. *Devbox*, an instance with 16 CPU and 4 CUDA GPU cores, provided by Institute of Fundamental Technological Research, Polish Academy of Sciences. AWS c5.9xlarge and c4.8xlarge instances, in the last 4 days, sponsored by Amazon.

We checked OpenSim simulator performance on Nvidia GPU cores (with NVBLAS support), however it didn't reduce the computation time. Final training was done on *Chuck* due to its large number of CPU cores. Different training runs were often resulting in very different models, some of which were stuck in local maxima early. Because of that, we started 5 to 10 separate training runs. We monitored their plots and visualized the current policy from time to time. If we judged that the model was not promising - e.g. it was stuck in local maxima or looked inferior to some other one - we stopped it and gave the resources to the more promising trainings.

In Figure 12 we present mean episode reward obtained by DDPG, PPO and ES during one training run. We achieved the best result using PPO, our score in the final stage was 18.63 meters. The final score was taken as the best score out of 5 submissions. The remaining scores were: 18.53, 16.14, 15.19, 14.5. The average of our 5 final submissions was 16.6.

### 8.3 Discussion

The average score during DDPG training was fluctuating a lot, sometimes it could also drop and never regain. We tried to tune the hyperparameters, without any gain though. Our problems were most probably due to bad data normalization or not enough episodes. The average score with PPO was not fluctuating as much and did not suddenly drop. That is why we switched to PPO and stayed with it until the end of the competition. ES usage in the Learning to Run environment should be more thoroughly examined.

There are a few things we would do differently now. We would try DDPG OpenAI baselines implementation. We would use simpler and well-known environment in the beginning, e.g. Walker2d and reproduce the results. We would make sure normalization is done correctly. We would try Layer Normalization instead of tedious manual normalization. We would tune hyperparameters in the end and in a rigorous way [11]. We would repeat an action  $n$  times (so-called *frame skip*). We would learn also on mirror flips of observations as shown in [23]. Finally, we would use TensorBoard or similar for visualizations.

## 9 Leveraging emergent muscles-activation patterns: from DDPG in a continuous action space to DQN with a set of actions

Andrew Melnik, Malte Schilling and Helge Ritter

A continuous action space with a large number of dimensions allows for complex manipulations. However, often, only a limited number of points in the action space is used. Furthermore, approaches like Deep Deterministic Policy Gradient (DDPG) may stick to a local optimum, where different optima have different sets of points in use. Therefore, to generalize over several local optima, we collected such points in the action space from different local optima and leveraged them in Deep Q-Network (DQN) [21] learning.

### 9.1 Methods

Our approach consisted of two parts. In the first part, we applied the DDPG model to the Learning to Run environment. Our model consisted of actor and critic sub-networks with parameters used as recommended by the Getting Started guide<sup>5</sup>. For the initial exploration of the continuous 18-dimensional action space, we used Ornstein-Uhlenbeck (OU) process [37] to generate temporally correlated noise which we added to the actor Neural Network (NN) output values. The model learned reliable muscles-activation patterns to perform successful locomotion (Fig. 13). After training, when the agent reached a performance plateau, the outputs of the actor NN became either equal to 0 or 1 (vectors of 18 binary values). We let the agents run further and collected a set of actor NN outputs (Table. 5, Fig. 15). To generalize over many successful locomotion strategies, we collected patterns from different independently trained agents.

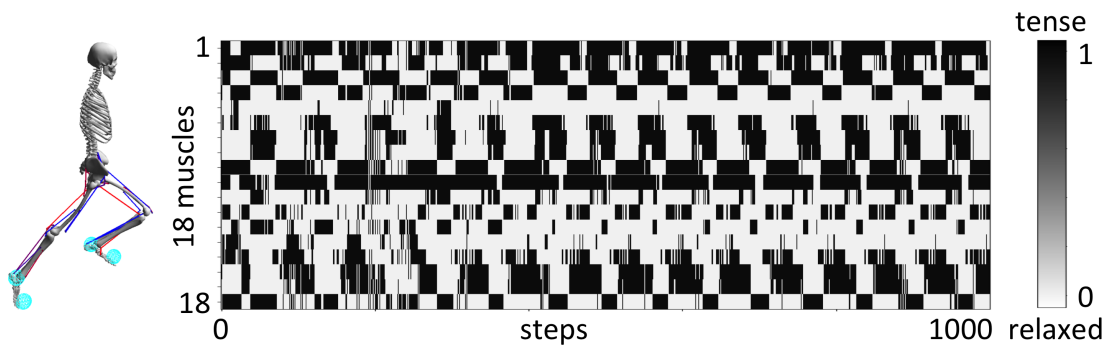


Fig. 13: Muscle activations in an episode. Columns are actor NN outputs after training.

In the second part, we used the set of collected patterns in two different scenarios. In the first scenario, we continued with training of DDPG models and used Q-values of the set of collected patterns for exploration, as, according to the critic NN, actor’s outputs were close, but not equal to the highest Q-values of the set of collected patterns (Fig. 16). The collected set allowed us to vastly reduce the action-search-space to a moderate set of useful muscles-activation patterns. Exploration by selection of patterns with highest Q-values allowed, in many cases, to increase the score further by 10-20%, after the DDPG agent reached a plateau of performance. In the second scenario, following the DQN approach, we used solely the critic NN (Fig. 14) to train new agents. To get Q-values for the set of collected muscles-activation patterns we concatenated them with state values and fed the batch to the critic NN (Fig. 14).

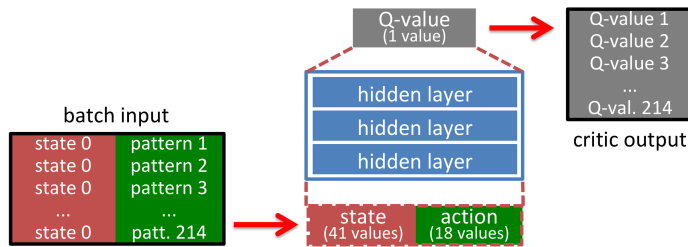


Fig. 14: Q-value estimation for a set of collected patterns.

<sup>5</sup> <https://github.com/stanfordnml/osim-rl>

Patt. ID	Muscles-Activation Pattern	Occur. freq.(%)	Sum (%)
1	110100001101001110	8.1	8.1
2	110100001100000110	4.5	12.6
3	000100001101001110	4.0	16.6
4	111001110110100001	3.5	20.1
5	101001110110100001	3.5	23.6
6	111001101100000110	2.7	26.3
7	110100001100011110	2.4	28.7
8	010100001101001110	2.3	31.0
...	...	...	...
214	000000001101001001	0.1	100.0

Table 5: Patterns of muscles activation and occurrence frequencies.

Authors Suppressed Due to Excessive Length



Fig. 15: Pareto distribution of occurrence frequencies of the muscles-activation patterns.

While a binary vector of 18 values would have  $2^{18} = 262144$  potential combinations, we found only a set of 214 muscles-activation patterns for trained models. That reduced exploration space to a moderate set of meaningful patterns. Certain patterns occurred more frequently than other (Table 5) with the 8 most frequent patterns representing already more than 30 % of all executed actions. About a half of the collected patterns (108 patterns) occurred only once (per episode of 1000 steps). Occurrence frequencies of the collected patterns (Table 5) demonstrated a Pareto distribution (Fig. 15).

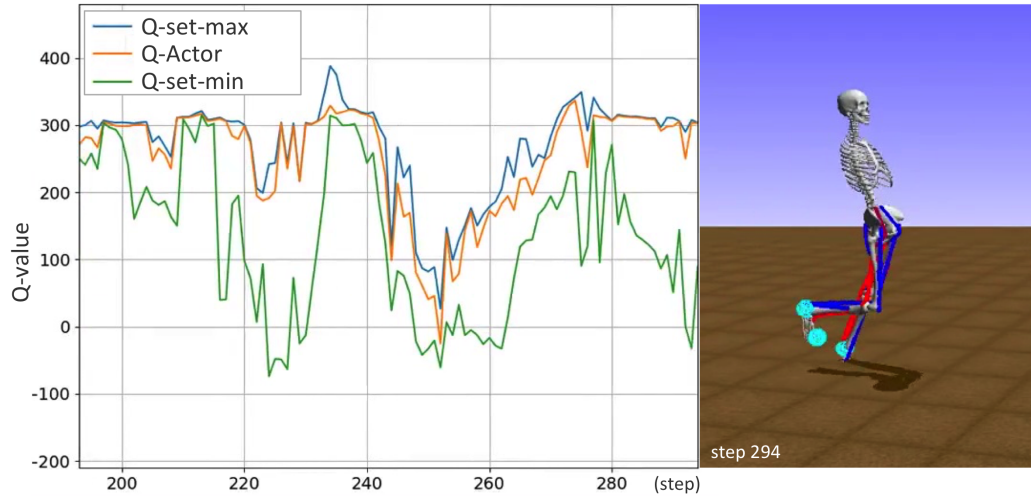


Fig. 16: Q-value range. Upper and lower curves represent maximum and minimum Q-values for the set of collected muscles-activation patterns in given states. The middle curve represents Q-values of the muscles-activation patterns proposed by the actor NN.



## 9.2 Discussion

We presented alternative (DQN based) approaches to explore and select actions in the continuous action space. However, the prerequisite is to have a set of meaningful actions for the task. Exploiting the set of collected muscles-activation patterns for further training has shown to lead to better performance and overall we want to consider in the future how to bootstrap the information further from emerging activation patterns and updating this collection. As a possible extension, we could take into account the reflection symmetry of muscles-activation patterns for the left and right legs. The set of 214 unique muscles-activation patterns for 18 muscles (two legs) contains only a set of 56 unique muscles-activation patterns for 9 muscles (per leg). In a way, this work is related to the ideas of hierarchical reinforcement learning [10] and the work on learning Dynamic Movement Primitives [30, 13] which are attractor systems of a lower dimensionality on the lower levels of such hierarchical systems.

## 10 Affiliations and acknowledgments

**Organizers:** Łukasz Kidziński, Carmichael Ong, Jennifer Hicks and Scott Delp are affiliated with Department of Bioengineering, Stanford University. Sharada Prasanna Mohanty, Sean Francis and Marcel Salathé are affiliated with Ecole Polytechnique Federale de Lausanne. Sergey Levine is affiliated with University of California, Berkeley.

**Team PKU (place 2nd, Section 2):** Zhewei Huang and Shuchang Zhou are affiliated with Beijing University. **Team reason8.ai (place 3rd, Section 3):** Mikhail Pavlov, Sergey Kolesnikov and Sergey Plis are affiliated with reason8.ai. **Team IMCL (place 4th, Section 4):** Zhibo Chen, Zhizheng Zhang, Jiale Chen and Jun Shi are affiliated with Immersive Media Computing Lab, University of Science and Technology of China. **Team deepsense.ai (place 6th, Section 5):** Henryk Michalewski is affiliated with Institute of Mathematics, Polish Academy of Sciences and deepsense.ai. Piotr Miłoś and Błażej Osiński are affiliated with Faculty of Mathematics, Informatics, and Mechanics, University of Warsaw and deepsense.ai. **Team THU-JGeek (place 8th, Section 6):** Zhuobin Zheng, Chun Yuan and Zhihui Lin are affiliated with Tunghai University. **Team Anton Pechenko (place 16th, Section 7):** Anton Pechenko is affiliated with Yandex. **Team Adam Stelmaszczyk (place 22nd, Section 8):** Adam Stelmaszczyk is affiliated with Faculty of Mathematics, Informatics, and Mechanics, University of Warsaw. Piotr Jarosik is affiliated with Institute of Fundamental Technological Research, Polish Academy of Sciences. **Team Andrew Melnik (place 28th, Section 9):** Andrew Melnik, Malte Schilling and Helge Ritter are affiliated with CITEC, Bielefeld University.

Team deepsense.ai was supported by the PL-Grid Infrastructure: Prometheus and Eagle supercomputers, located respectively in the Academic Computer Center Cyfronet at the AGH University of Science and Technology in Kraków and the Supercomputing and Networking Center in Poznań. The deepsense.ai team also expresses gratitude to NVIDIA and Goodeep for providing additional computational resources used during the experiment.

The challenge was co-organized by the Mobilize Center, a National Institutes of Health Big Data to Knowledge (BD2K) Center of Excellence supported through Grant U54EB020405. The challenge was partially sponsored by Nvidia who provided DGX Station<sup>TM</sup> for the first prize in the challenge, and GPUs Titan V for the second and the third prize, by Amazon Web Services who provided 30000 USD in cloud credits for participants, and by Toyota Research Institute who funded one travel grant.

## References

1. Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G.S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., Zheng, X.: TensorFlow: Large-scale machine learning on heterogeneous systems (2015). URL <https://www.tensorflow.org/>. Software available from tensorflow.org
2. Anonymous: Distributional policy gradients. International Conference on Learning Representations (2018). URL <https://openreview.net/forum?id=SyZipzbCb>
3. Ba, J.L., Kiros, J.R., Hinton, G.E.: Layer normalization. arXiv preprint arXiv:1607.06450 (2016)
4. Bratko, I., Urbančič, T., Sammut, C.: Behavioural cloning: Phenomena, results and problems. IFAC Proceedings Volumes **28**(21), 143 – 149 (1995). DOI [https://doi.org/10.1016/S1474-6670\(17\)46716-4](https://doi.org/10.1016/S1474-6670(17)46716-4). URL <http://www.sciencedirect.com/science/article/pii/S1474667017467164>
5. Clevert, D.A., Unterthiner, T., Hochreiter, S.: Fast and accurate deep network learning by exponential linear units (elus). arXiv preprint arXiv:1511.07289 (2015)
6. Dhariwal, P., Hesse, C., Plappert, M., Radford, A., Schulman, J., Sidor, S., Wu, Y.: OpenAI Baselines. <https://github.com/openai/baselines> (2017)
7. Dietterich, T.G., et al.: Ensemble methods in machine learning. Multiple classifier systems **1857**, 1–15 (2000)
8. Dorigo, M., Colombetti, M.: Robot Shaping: An Experiment in Behavior Engineering. MIT Press, Cambridge, MA, USA (1997)
9. Heess, N., Sriram, S., Lemmon, J., Merel, J., Wayne, G., Tassa, Y., Erez, T., Wang, Z., Eslami, A., Riedmiller, M., et al.: Emergence of locomotion behaviours in rich environments. arXiv preprint arXiv:1707.02286 (2017)
10. Heess, N., Wayne, G., Tassa, Y., Lillicrap, T.P., Riedmiller, M.A., Silver, D.: Learning and transfer of modulated locomotor controllers. CoRR **abs/1610.05182** (2016). URL <http://arxiv.org/abs/1610.05182>
11. Henderson, P., Islam, R., Bachman, P., Pineau, J., Precup, D., Meger, D.: Deep Reinforcement Learning that Matters. ArXiv e-prints (2017)
12. Hessel, M., Modayil, J., Van Hasselt, H., Schaul, T., Ostrovski, G., Dabney, W., Horgan, D., Piot, B., Azar, M., Silver, D.: Rainbow: Combining improvements in deep reinforcement learning. arXiv preprint arXiv:1710.02298 (2017)
13. Ijspeert, A., Nakanishi, J., Pastor, P., Hoffmann, H., Schaal, S.: Dynamical movement primitives: Learning attractor models for motor behaviors. Neural Computation **25**, 328–373 (2013). URL <http://www-clmc.usc.edu/publications/I/ijspeert-NC2013.pdf>. Clmc
14. Jaśkowski, W., Lykkebø, O.R., Toklu, N.E., Trifterer, F., Buk, Z., Koutník, J., Gomez, F.: Reinforcement Learning to Run... Fast. In: S. Escalera, M. Weimer (eds.) NIPS 2017 Competition Book. Springer, Springer (2018)
15. Ji, S., Xu, W., Yang, M., Yu, K.: 3d convolutional neural networks for human action recognition. IEEE transactions on pattern analysis and machine intelligence **35**(1), 221–231 (2013)
16. Kidziński, Ł., Sharada, M.P., Ong, C., Hicks, J., Francis, S., Levine, S., Salathé, M., Delp, S.: Learning to run challenge: Synthesizing physiologically accurate motion using deep reinforcement learning. In: S. Escalera, M. Weimer (eds.) NIPS 2017 Competition Book. Springer, Springer (2018)
17. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. CoRR **abs/1412.6980** (2014). URL <http://arxiv.org/abs/1412.6980>
18. Klambauer, G., Unterthiner, T., Mayr, A., Hochreiter, S.: Self-normalizing neural networks. arXiv preprint arXiv:1706.02515 (2017)
19. Lillicrap, T.P., Hunt, J.J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., Wierstra, D.: Continuous control with deep reinforcement learning. arXiv preprint arXiv:1509.02971 (2015)
20. Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., Riedmiller, M.A.: Playing atari with deep reinforcement learning. CoRR **abs/1312.5602** (2013). URL <http://arxiv.org/abs/1312.5602>
21. Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A.A., Veness, J., Bellemare, M.G., Graves, A., Riedmiller, M., Fidjeland, A.K., Ostrovski, G., et al.: Human-level control through deep reinforcement learning. Nature **518**(7540), 529–533 (2015)
22. Osband, I., Blundell, C., Pritzel, A., Van Roy, B.: Deep exploration via bootstrapped dqn. In: Advances in Neural Information Processing Systems, pp. 4026–4034 (2016)
23. Pavlov, M., Kolesnikov, S., Plis, S.M.: Run, skeleton, run: skeletal model in a physics-based simulation. ArXiv e-prints (2017)
24. Plappert, M.: keras-rl. <https://github.com/matthiasplappert/keras-rl> (2016)
25. Plappert, M., Houthoofd, R., Dhariwal, P., Sidor, S., Chen, R.Y., Chen, X., Asfour, T., Abbeel, P., Andrychowicz, M.: Parameter space noise for exploration. arXiv preprint arXiv:1706.01905 (2) (2017)

26. Rumelhart, D.E., Hinton, G.E., Williams, R.J.: Parallel distributed processing: Explorations in the microstructure of cognition, vol. 1. chap. Learning Internal Representations by Error Propagation, pp. 318–362. MIT Press, Cambridge, MA, USA (1986). URL <http://dl.acm.org/citation.cfm?id=104279.104293>
27. Rusu, A.A., Rabinowitz, N.C., Desjardins, G., Soyer, H., Kirkpatrick, J., Kavukcuoglu, K., Pascanu, R., Hadsell, R.: Progressive neural networks. CoRR [abs/1606.04671](https://arxiv.org/abs/1606.04671) (2016). URL <http://arxiv.org/abs/1606.04671>
28. Salimans, T., Ho, J., Chen, X., Sidor, S., Sutskever, I.: Evolution Strategies as a Scalable Alternative to Reinforcement Learning. ArXiv e-prints (2017)
29. Salimans, T., Ho, J., Chen, X., Sidor, S., Sutskever, I.: Starter code for evolution strategies. <https://github.com/openai/evolution-strategies-starter> (2017)
30. Schaal, S.: Dynamic movement primitives -a framework for motor control in humans and humanoid robotics. In: H. Kimura, K. Tsuchiya, A. Ishiguro, H. Witte (eds.) Adaptive Motion of Animals and Machines, pp. 261–280. Springer Tokyo, Tokyo (2006). DOI 10.1007/4-431-31381-8\_23. URL [https://doi.org/10.1007/4-431-31381-8\\_23](https://doi.org/10.1007/4-431-31381-8_23)
31. Schaul, T., Quan, J., Antonoglou, I., Silver, D.: Prioritized experience replay. arXiv preprint arXiv:1511.05952 (2015)
32. Schulman, J., Wolski, F., Dhariwal, P., Radford, A., Klimov, O.: Proximal policy optimization algorithms. CoRR [abs/1707.06347](https://arxiv.org/abs/1707.06347) (2017). URL <http://arxiv.org/abs/1707.06347>
33. Schulman, J., Wolski, F., Dhariwal, P., Radford, A., Klimov, O.: Proximal Policy Optimization Algorithms. ArXiv e-prints (2017)
34. Silver, D., Lever, G., Heess, N., Degris, T., Wierstra, D., Riedmiller, M.: Deterministic policy gradient algorithms. In: Proceedings of the 31st International Conference on Machine Learning (ICML-14), pp. 387–395 (2014)
35. Stelmaszczyk, A., Jarosik, P.: Our NIPS 2017: Learning to Run source code. <https://github.com/AdamStelmaszczyk/learning2run> (2017)
36. Sutton, R.S., Precup, D., Singh, S.: Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning. Artificial Intelligence **112** (1999)
37. Uhlenbeck, G.E., Ornstein, L.S.: On the theory of the brownian motion. Physical review **36**(5), 823 (1930)
38. Wiering, M., Schmidhuber, J.: HQ-learning. Adaptive Behaviour **6** (1997)