# Reinforcement learning lab report

Adrian-Gabriel Bălănescu

June 18, 2021

# Contents

# Chapter 1

# Game description

*Gravitar* is an arcade style game released by Atari in 1982. The game was known for its high level of difficulty [1].

## 1.1 Gameplay

The game takes place in a fictional solar system. The player controls a spacecraft with the help of five buttons: two to rotate the ship left and right, one to shoot, one to activate the thruster and one for both a tractor beam and a force field. In the side-view levels, the player has to destroy red bunkers that shoot constantly, and can also use the tractor beam to pick up blue fuel tanks. Once all of the bunkers are destroyed, the planet will blow up, and the player will earn a bonus. Once all planets are destroyed, the player will move onto another solar system. The player will lose a life if he crashes into the terrain or gets hit by an enemy's shot, and the game will end immediately if fuel runs out [1].

After completing all 11 planets (or alternatively completing the reactor three times) the player enters the second universe and the gravity will reverse. Instead of dragging the ship towards the planet surface, the gravity pushes it away. In the third universe the landscape becomes invisible and the gravity is positive again. The final, fourth universe, has invisible landscape and reverse gravity. After completing the fourth universe the game starts over [1].



Figure 1.1: Gravitar logo
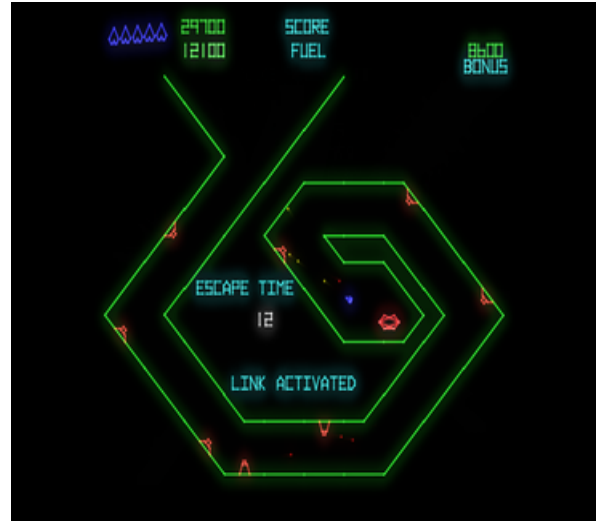
Figure 1.2: A level from Gravitar



Figure 1.3: The reactor

## 1.2 Gravitar in SLM-lab

*SLM-lab* [2] is a software framework for reproducible reinforcement learning (RL) research. It implements most of the canonical RL algorithms in a modular platform that allows for easy experimentation and analysis. It includes the OpenAI Gym [3] environments, which contains a lot of Atari games emulations, one being Gravitar.

In this environment, the observation is an RGB image of the screen, which is an array of shape (210, 160, 3). Each action is repeatedly performed for a duration of $k$ frames, where $k$ is uniformly sampled from $\{2, 3, 4\}$. The action space is discrete and has 18 available commands like fire, up, down, right, left and so on.

The Gravitar environment can be loaded in SLM-lab by specifying in the .json configuration file, at the "env" section the "name" of the environment. In this case, "GravitarNoFrameskip–v4".

# Chapter 2

# Algorithms

In this chapter, the algorithms used will be briefly explained.

## 2.1 Reinforce

Reinforce is a Monte Carlo variant of the policy-gradient algorithm in reinforcement learning. The agent collects samples of an episode using its current policy and uses it to update the policy parameter $\theta$. It is updated as an off-policy algorithm:

$$\nabla_\theta J(\theta) = E_\pi[G_t \nabla_\theta ln \pi_\theta(a_t, s_t)] \tag{2.1}$$

## 2.2 A2C - Advantage Actor-Critic

The Advantage Actor-Critic algorithm relies on two important concepts:

**The Actor-Critic architecture.** This type of algorithms learn two parametrized functions simultaneously. The actor learns the actual policy that the agent uses to act in the world, while the critic learns the value function that evaluates the state-action pairs.

Using the critic to provide a reinforcement signal has the advantage of being able to feed a more continuous signal for training that is useful for example in sparse rewards environments. However, training becomes more difficult because the policy is dependent on the quality of the value function. In the beginning of the training process, the value function is not learned, so the evaluation signal provided doesn't help much to select better actions. Therefore, usually an Advantage function is actually learned.

**Reinforcing signal based on the Advantage function.** The advantage function measures how much an action is better or worse than a policy's average action in a certain state. It is defined as follows:

$$A^\pi(s_t, a_t) = Q^\pi(s_t, a_t) - V^\pi(s_t) \tag{2.2}$$

It has the following properties:

- The advantage is a relative measure. It quantifies whether a particular action $a$ in a state $s$ will improve the policy relative to the value $V^\pi(s)$.

- If an action is worse than the average action, but the expected return is still positive - the action taken should become less probable. This is the behavior of the advantage function, in contrast with other methods that might actually encourage that action.

- The advantage function is able to capture long-term effects of an action, because it considers all future time steps, while ignoring the effects of all the actions to date.

### 2.2.1 Estimating the advantage with n-Step returns

The n-Step estimate is formulated as follows:

$$A^\pi_{Nstep}(s_t, a_t) = Q^\pi(s_t, a_t) - V^\pi(s_t) \approx r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \cdots + \gamma^n r_{t+n} + \gamma^{n+1}\hat{V}^\pi(s_{t+n+1}) - \hat{V}^\pi(s_t) \tag{2.3}$$

The hyperparameter n - the number of actual rewards, controls the bias-variance trade-off.

### 2.2.2 Estimating the advantage with Generalized Advantage Estimation

GAE is defined as an exponentially weighted average of all of the n-step forward return advantages:

$$A^\pi_{GAE}(s_t, a_t) = \sum_{l=0}^{\infty} (\gamma\lambda)^l \delta_{t+l} \tag{2.4}$$

where, $\delta_t = r_t + \gamma V^\pi(s_{t+1}) - V^\pi(s_t)$

The variance is controlled through the decay coefficient $\lambda$

## 2.3 PPO - Proximal Policy Optimization

PPO [4] aims at avoiding the problem of performance collapse through a surrogate objective that is guaranteed to have monotonic improvement. It works by extending either RE-INFORCE or an actor-critic algorithm with the modified objective. It is data efficient and reliable as TRPO [5] but uses first-order optimization. Let $r_t(\theta)$ denote the probability ratio $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}$ so r($\theta_{old}$) = 1. TRPO maximizes a "surrogate" objective:

$$L^{CPI}(\theta) = \hat{E}_t \left[ \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)} \hat{A}_t \right] = \hat{E}_t \left[ r_t(\theta)\hat{A}_t \right] \tag{2.5}$$

where $CPI$ refers to a conservative policy iteration.[4] Without a constraint, maximization of $L^{CPI}$ would lead to an excessively large policy update, therefore, PPO modifies the objective to penalize changes to the policy that moves $r_t(\theta)$ away from 1:

$$J^{CLIP}(\theta) = \hat{E} \left[ min \left( r_t(\theta)\hat{A}, clip \left( r_t(\theta), 1 - \epsilon, 1 + \epsilon \right) \hat{A} \right) \right] \tag{2.6}$$

where $\epsilon$ is a hyperparameter. The intuition is as follows: The *clip* term modifies the surrogate objective by clipping the probability ratio, which removes the incentive to move $r_t$ outside the interval $[1 - \epsilon, 1 + \epsilon]$ [4].

# Chapter 3

# Experimental Results

## 3.1 Reinforce

For training Reinforce, the following configuration has been used:

```
1    {
2      "reinforce_gravitar": {
3        "agent": [{
4          "name": "Reinforce",
5          "algorithm": {
6            "name": "Reinforce",
7            "action_pdtype": "default",
8            "action_policy": "default",
9            "center_return": true,
10           "explore_var_spec": null,
11           "gamma": 0.99,
12           "entropy_coef_spec": {
13             "name": "linear_decay",
14             "start_val": 0.01,
15             "end_val": 0.001,
16             "start_step": 1000,
17             "end_step": 300000
18           },
19           "training_frequency": 1
20         },
21         "memory": {
22           "name": "OnPolicyReplay",
23           "batch_size": 64
24         },
25         "net": {
26           "type": "ConvNet",
27           "shared": true,
28           "conv_hid_layers": [
29           [32, 8, 4, 0, 1],
30           [64, 4, 2, 0, 1],
31           [32, 3, 1, 0, 1]
32           ],
33           "fc_hid_layers": [512],
34           "hid_layers_activation": "relu",
```

```
35        "init_fn": "orthogonal_",
36        "normalize": true,
37        "batch_norm": false,
38        "clip_grad_val": 0.5,
39        "loss_spec": {
40          "name": "MSELoss"
41        },
42        "optim_spec": {
43          "name": "Adam",
44          "lr": 0.002
45        },
46        "lr_scheduler_spec": null,
47        "gpu": true
48      }
49    }],
50    "env": [{
51      "name": "GravitarNoFrameskip-v4",
52      "frame_op": "concat",
53      "frame_op_len": 4,
54      "reward_scale": "sign",
55      "num_envs": 16,
56      "max_t": null,
57      "max_frame": 300000
58    }],
59    "body": {
60      "product": "outer",
61      "num": 1
62    },
63    "meta": {
64      "distributed": false,
65      "eval_frequency": 10000,
66      "max_session": 1,
67      "max_trial": 1
68    }
69  }
70 }
```

The following parameters are set:

- The algorithm is set to Reinforce

- The action policy is the default one

- The agent uses a convolutional network with three layers connected to a single linear layer with ReLU activations.

- The optimizer is Adam.

- Algorithm is trained for 300k steps.

- Evaluation occurs every 10k steps.

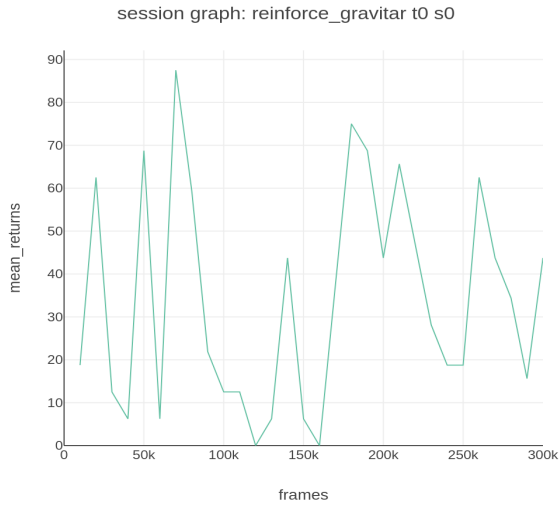  The following returns are obtained:
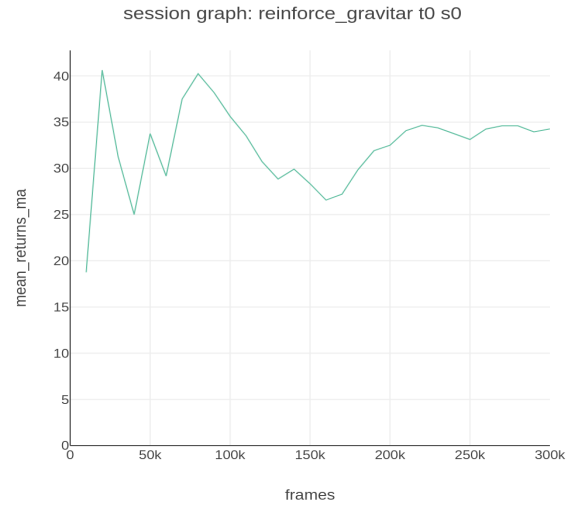
Figure 3.1: Trial graph



Figure 3.2: Trial graph with moving average

Figure 3.3: Reinforce trial graphs from SLM lab. On the vertical axis the total rewards are shown and on the horizontal axis the number of frames.

## 3.2 A2C

For training A2C, two experiments have been considered:

1. Training with generalized advantage function (GAE) The config file is as follows:

```
1    {
2      "a2c_gae_gravitar": {
3        "agent": [{
4          "name": "A2C",
5          "algorithm": {
6            "name": "ActorCritic",
7            "action_pdtype": "default",
8            "action_policy": "default",
9            "explore_var_spec": null,
10           "gamma": 0.99,
11           "lam": 0.95,
12           "num_step_returns": null,
13           "entropy_coef_spec": {
14             "name": "no_decay",
15             "start_val": 0.01,
16             "end_val": 0.01,
17             "start_step": 0,
18             "end_step": 0
19           },
20           "val_loss_coef": 0.5,
21           "training_frequency": 32
22         },
23         "memory": {
24           "name": "OnPolicyBatchReplay"
25         },
```

8

```
26          "net": {
27           "type": "ConvNet",
28           "shared": true,
29           "conv_hid_layers": [
30           [32, 8, 4, 0, 1],
31           [64, 4, 2, 0, 1],
32           [32, 3, 1, 0, 1]
33           ],
34           "fc_hid_layers": [512],
35           "hid_layers_activation": "relu",
36           "init_fn": "orthogonal_",
37           "normalize": true,
38           "batch_norm": false,
39           "clip_grad_val": 0.5,
40           "use_same_optim": false,
41           "loss_spec": {
42            "name": "MSELoss"
43           },
44           "actor_optim_spec": {
45            "name": "RMSprop",
46            "lr": 7e-4,
47            "alpha": 0.99,
48            "eps": 1e-5
49           },x
50           "critic_optim_spec": {
51            "name": "RMSprop",
52            "lr": 7e-4,
53            "alpha": 0.99,
54            "eps": 1e-5
55           },
56           "lr_scheduler_spec": null,
57           "gpu": true
58          }
59         }],
60         "env": [{
61          "name": "GravitarNoFrameskip-v4",
62          "frame_op": "concat",
63          "frame_op_len": 4,
64          "reward_scale": "sign",
65          "num_envs": 16,
66          "max_t": null,
67          "max_frame": 1e6
68         }],
69         "body": {
70          "product": "outer",
71          "num": 1
72         },
73         "meta": {
74          "distributed": false,
75          "eval_frequency": 10000,
76          "log_frequency": 10000,
77          "rigorous_eval": 0,
78          "max_session": 1,
79          "max_trial": 1
```

9

```
80          }
81        }
82      }
```

The following parameters are set:

- The algorithm is A2C
- The action policy is the default one
- GAE is used as an Advantage function estimator
- The agent uses a convolutional network with three layers connected to a single linear layer with ReLU activations.
- The optimizer is RMSProp.
- Algorithm is trained for 1M steps.
- Evaluation occurs every 10k steps.

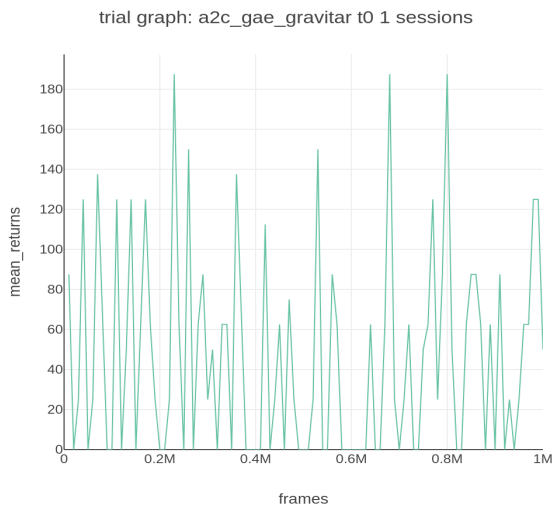The following returns are obtained:
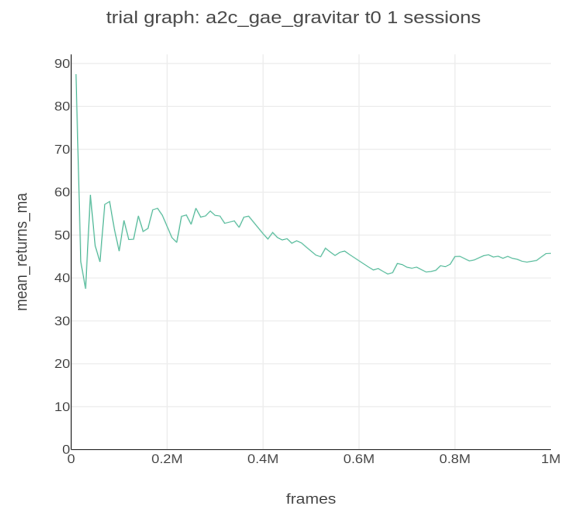


Figure 3.4: Trial graph

Figure 3.5: Trial graph with moving average

Figure 3.6: A2C trial graphs from SLM lab. On the vertical axis the total rewards are shown and on the horizontal axis the number of frames.

2. Training with n-Step Returns

   The config file is almost the same as the GAE version, it differs at the *num_step_returns* parameter.

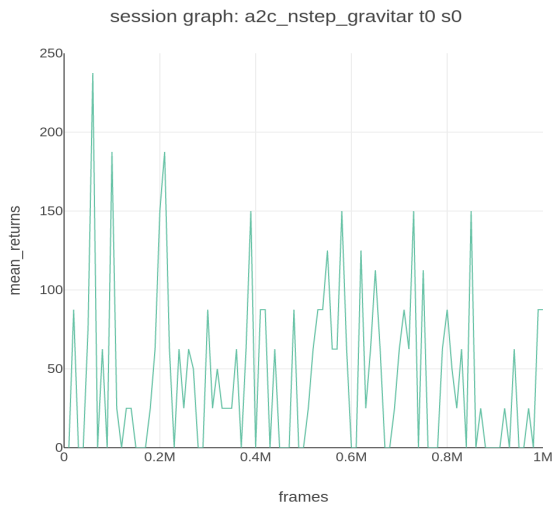   The following returns are obtained:
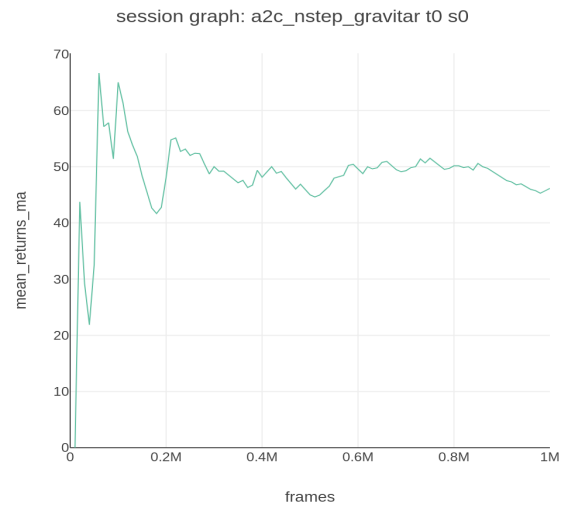


Figure 3.7: Trial graph

Figure 3.8: Trial graph with moving average

Figure 3.9: A2C trial graphs from SLM lab. On the vertical axis the total rewards are shown and on the horizontal axis the number of frames.

## 3.3  PPO

For training the Proximal Policy Optimization algorithm, the following configuration has been used:

```
1    {
2     "ppo_gravitar": {
3      "agent": [{
4       "name": "PPO",
5       "algorithm": {
6        "name": "PPO",
7        "action_pdtype": "default",
8        "action_policy": "default",
9        "explore_var_spec": null,
10       "gamma": 0.99,
11       "lam": 0.70,
12       "clip_eps_spec": {
13        "name": "no_decay",
14        "start_val": 0.10,
15        "end_val": 0.10,
16        "start_step": 0,
17        "end_step": 0
18       },
```

```
19          "entropy_coef_spec": {
20           "name": "no_decay",
21           "start_val": 0.01,
22           "end_val": 0.01,
23           "start_step": 0,
24           "end_step": 0
25          },
26          "val_loss_coef": 0.5,
27          "time_horizon": 128,
28          "minibatch_size": 256,
29          "training_epoch": 4
30         },
31         "memory": {
32          "name": "OnPolicyBatchReplay",
33         },
34         "net": {
35          "type": "ConvNet",
36          "shared": true,
37          "conv_hid_layers": [
38          [32, 8, 4, 0, 1],
39          [64, 4, 2, 0, 1],
40          [32, 3, 1, 0, 1]
41          ],
42          "fc_hid_layers": [512],
43          "hid_layers_activation": "relu",
44          "init_fn": "orthogonal_",
45          "normalize": true,
46          "batch_norm": false,
47          "clip_grad_val": 0.5,
48          "use_same_optim": false,
49          "loss_spec": {
50           "name": "MSELoss"
51          },
52          "actor_optim_spec": {
53           "name": "Adam",
54           "lr": 2.5e-4,
55          },
56          "critic_optim_spec": {
57           "name": "Adam",
58           "lr": 2.5e-4,
59          },
60          "lr_scheduler_spec": {
61           "name": "LinearToZero",
62           "frame": 1e7
63          },
64          "gpu": true
65         }
66        }],
67        "env": [{
68         "name": "GravitarNoFrameskip-v4",
69         "frame_op": "concat",
70         "frame_op_len": 4,
71         "reward_scale": "sign",
72         "num_envs": 4,
```

```
73        "max_t": null ,
74        "max_frame": 1e6
75      }] ,
76      "body": {
77       "product": "outer",
78       "num": 1
79      },
80      "meta": {
81       "distributed": false ,
82       "eval_frequency": 10000,
83       "log_frequency": 10000,
84       "rigorous_eval": 0,
85       "max_session": 1,
86       "max_trial": 1
87      }
88    }
89  }
```

The following parameters are set:

- The action policy is the default policy for a discrete action space.

- Generalized advantage estimation is used.

- The clipping parameters are set in lines 12-17

- In the lines 34-43 the network architecture is set as a 3 layer convolutional network with a hidden fully connected layer and ReLU activations.

- The opimizer is set to Adam both for the actor and the critic networks.

- The environment is Gravitar.

- Evaluation is done every 10k steps.

- The algorithm is trained for 1M steps.

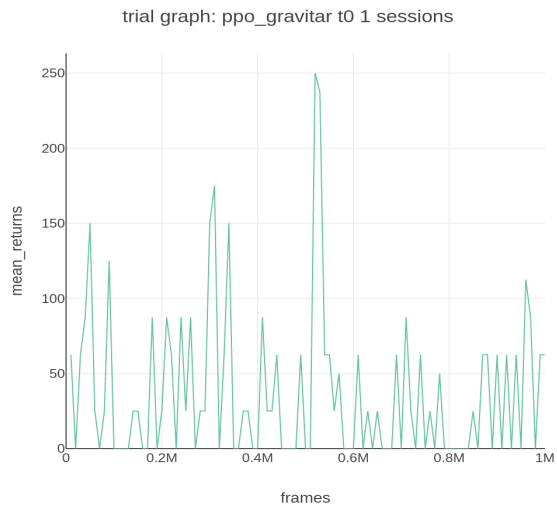The following returns are obtained:
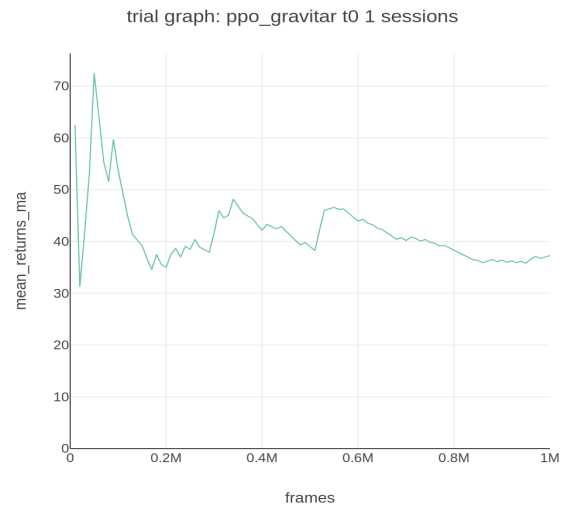


Figure 3.10: Trial graph



Figure 3.11: Trial graph with moving average

Figure 3.12: PPO trial graphs from SLM lab. On the vertical axis the total rewards are shown and on the horizontal axis the number of frames.

# Chapter 4

# Conclusions

To conclude, Gravitar is a very difficult game to solve. Not only due to it's more complex physics and controls, but also due to the fact that the levels are increasingly difficult and quite different from one to another. No algorithms from the ones tested was able to decently solve this environment. As a comparison, apart from Reinforce, which (with no surprise) performed the poorest, the other algorithms performed similarly, with the n-step A2C having the best mean returns average of approx. 50. Also, all the algorithms suffered performance degradation as more steps were taken.

# Bibliography

[1] *Gravitar*.          `https://en.wikipedia.org/w/index.php?title=Gravitar&oldid=` `1018480060`, April 2021. Page Version ID: 1018480060. 1, 1.1

[2] Wah Loon Keng and Laura Graesser. Slm lab. `https://github.com/kengz/SLM-Lab`, 2017. 1.2

[3] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016. 1.2

[4] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal Policy Optimization Algorithms. *arXiv:1707.06347 [cs]*, August 2017. arXiv: 1707.06347. 2.3, 2.3, 2.3

[5] John Schulman, Sergey Levine, Philipp Moritz, Michael I. Jordan, and Pieter Abbeel. Trust Region Policy Optimization. *arXiv:1502.05477 [cs]*, April 2017. arXiv: 1502.05477. 2.3