# Assignment 2: Policy Gradient

**Andrew ID:** `shriishs`
**Collaborators:** `None`
**NOTE:** Please do **NOT** change the sizes of the answer blocks or plots.

# 5   Small-Scale Experiments

## 5.1   Experiment 1 (Cartpole) – [25 points total]

### 5.1.1   Configurations

---

**Q5.1.1**

```
python rob831/scripts/run_hw2.py --env_name CartPole-v0 -n 100 -b 1000 \
    -dsa --exp_name q1_sb_no_rtg_dsa

python rob831/scripts/run_hw2.py --env_name CartPole-v0 -n 100 -b 1000 \
    -rtg -dsa --exp_name q1_sb_rtg_dsa

python rob831/scripts/run_hw2.py --env_name CartPole-v0 -n 100 -b 1000 \
    -rtg --exp_name q1_sb_rtg_na

python rob831/scripts/run_hw2.py --env_name CartPole-v0 -n 100 -b 5000 \
    -dsa --exp_name q1_lb_no_rtg_dsa

python rob831/scripts/run_hw2.py --env_name CartPole-v0 -n 100 -b 5000 \
    -rtg -dsa --exp_name q1_lb_rtg_dsa

python rob831/scripts/run_hw2.py --env_name CartPole-v0 -n 100 -b 5000 \
    -rtg --exp_name q1_lb_rtg_na
```
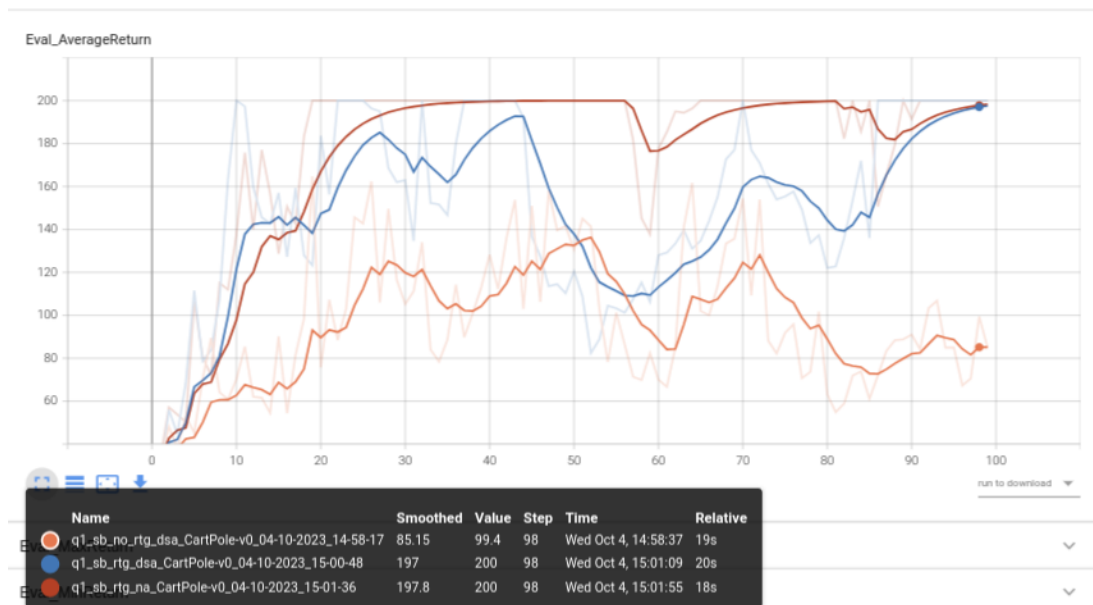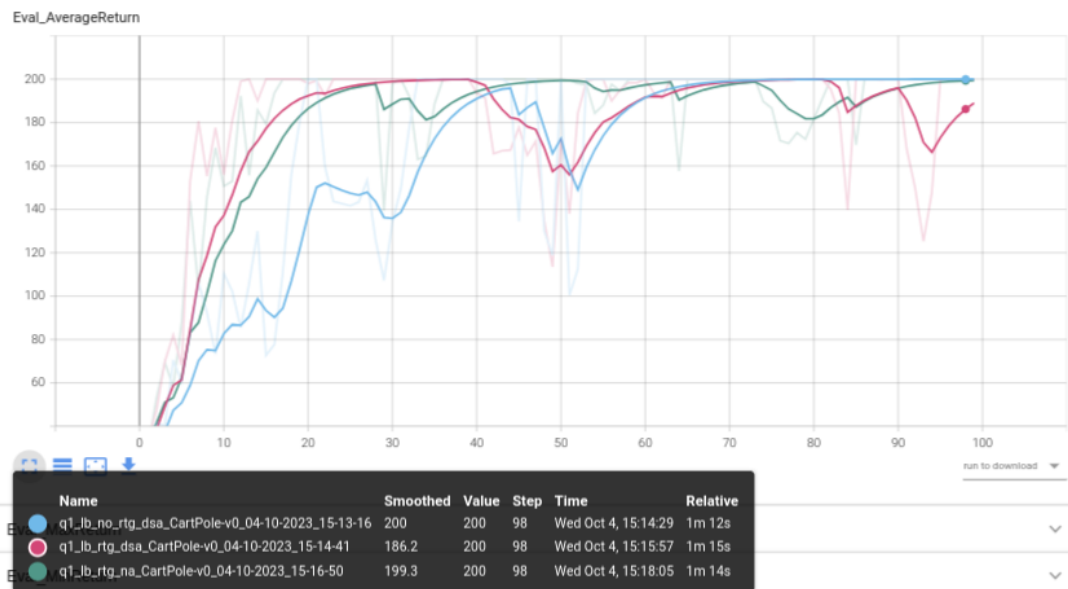
---

### 5.1.2   Plots

### 5.1.2.1   Small batch – [5 points]

---

**Q5.1.2.1**



---

### 5.1.2.2    Large batch – [5 points]

Q5.1.2.2



| Name | Smoothed | Value | Step | Time | Relative |
|------|----------|-------|------|------|----------|
| q1_lb_no_rtg_dsa_CartPole-v0_04-10-2023_15-13-16 | 200 | 200 | 98 | Wed Oct 4, 15:14:29 | 1m 12s |
| q1_lb_rtg_dsa_CartPole-v0_04-10-2023_15-14-41 | 186.2 | 200 | 98 | Wed Oct 4, 15:15:57 | 1m 15s |
| q1_lb_rtg_na_CartPole-v0_04-10-2023_15-16-50 | 199.3 | 200 | 98 | Wed Oct 4, 15:18:05 | 1m 14s |

### 5.1.3    Analysis

### 5.1.3.1    Value estimator – [5 points]

Q5.1.3.1

It can be seen from both the small and large batch experiments that the reward-to-go estimator has better performance without advantage standardization. The difference is seen more when the batch size is smaller.

### 5.1.3.2    Advantage standardization – [5 points]

Q5.1.3.2

Yes, advantage standardization helped the policy reach a high return much faster and remain more stable than the case without it.

### 5.1.3.3   Batch size – [5 points]

---

**Q5.1.3.1**

Using a larger batch size helps the policy reach a high return much faster. The policy then fluctuates about the converged value. The performance is much better.

---

## 5.2   Experiment 2 (InvertedPendulum) – [15 points total]

### 5.2.1   Configurations – [5 points]

---

**Q5.2.1**

```
python rob831/scripts/run_hw2.py --env_name InvertedPendulum-v4 \
    --ep_len 1000 --discount 0.9 -n 100 -l 2 -s 64 -b 1000 -lr 1e-2 -rtg \
    --exp_name q2_b1000_r1e-2

python rob831/scripts/run_hw2.py --env_name InvertedPendulum-v4 \
    --ep_len 1000 --discount 0.9 -n 100 -l 2 -s 64 -b 1000 -lr 2e-2 -rtg \
    --exp_name q2_b1000_r2e-2

python rob831/scripts/run_hw2.py --env_name InvertedPendulum-v4 \
    --ep_len 1000 --discount 0.9 -n 100 -l 2 -s 64 -b 2000 -lr 1e-2 -rtg \
    --exp_name q2_b2000_r1e-2

python rob831/scripts/run_hw2.py --env_name InvertedPendulum-v4 \
    --ep_len 1000 --discount 0.9 -n 100 -l 2 -s 64 -b 2000 -lr 2e-2 -rtg \
    --exp_name q2_b2000_r2e-2

python rob831/scripts/run_hw2.py --env_name InvertedPendulum-v4 \
    --ep_len 1000 --discount 0.9 -n 100 -l 2 -s 64 -b 3000 -lr 1e-2 -rtg \
    --exp_name q2_b3000_r1e-2

python rob831/scripts/run_hw2.py --env_name InvertedPendulum-v4 \
    --ep_len 1000 --discount 0.9 -n 100 -l 2 -s 64 -b 3000 -lr 2e-2 -rtg \
    --exp_name q2_b3000_r2e-2

python rob831/scripts/run_hw2.py --env_name InvertedPendulum-v4 \
    --ep_len 1000 --discount 0.9 -n 100 -l 2 -s 64 -b 4000 -lr 1e-2 -rtg \
    --exp_name q2_b4000_r1e-2

python rob831/scripts/run_hw2.py --env_name InvertedPendulum-v4 \
    --ep_len 1000 --discount 0.9 -n 100 -l 2 -s 64 -b 4000 -lr 2e-2 -rtg \
    --exp_name q2_b4000_r2e-2
```
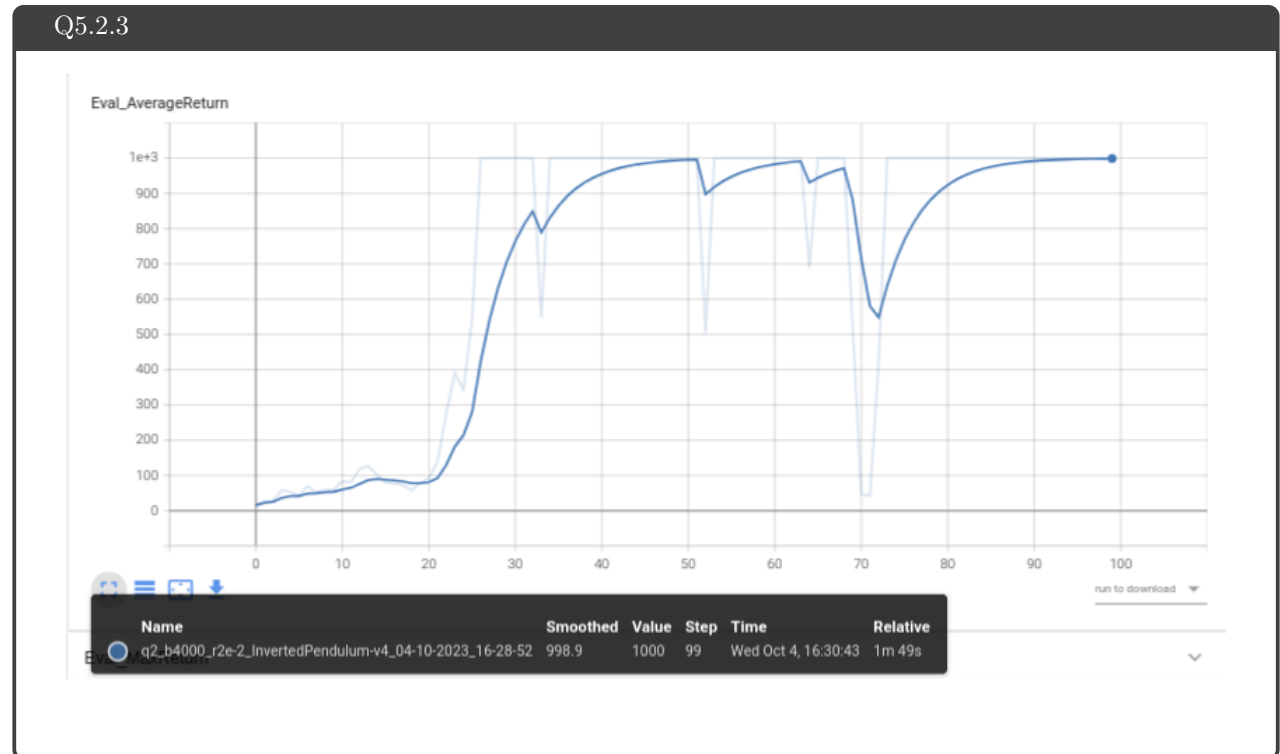
---

### 5.2.2   smallest b* and largest r* (same run) – [5 points]

---

**Q5.2.2**

Smallest **b\*** = 4000
Largest **r\*** = 2e-2

---

### 5.2.3    Plot – [5 points]

---

Q5.2.3

Eval_AverageReturn



| Name | Smoothed | Value | Step | Time | Relative |
|------|----------|-------|------|------|----------|
| q2_b4000_r2e-2_InvertedPendulum-v4_04-10-2023_16-28-52 | 998.9 | 1000 | 99 | Wed Oct 4, 16:30:43 | 1m 49s |

---

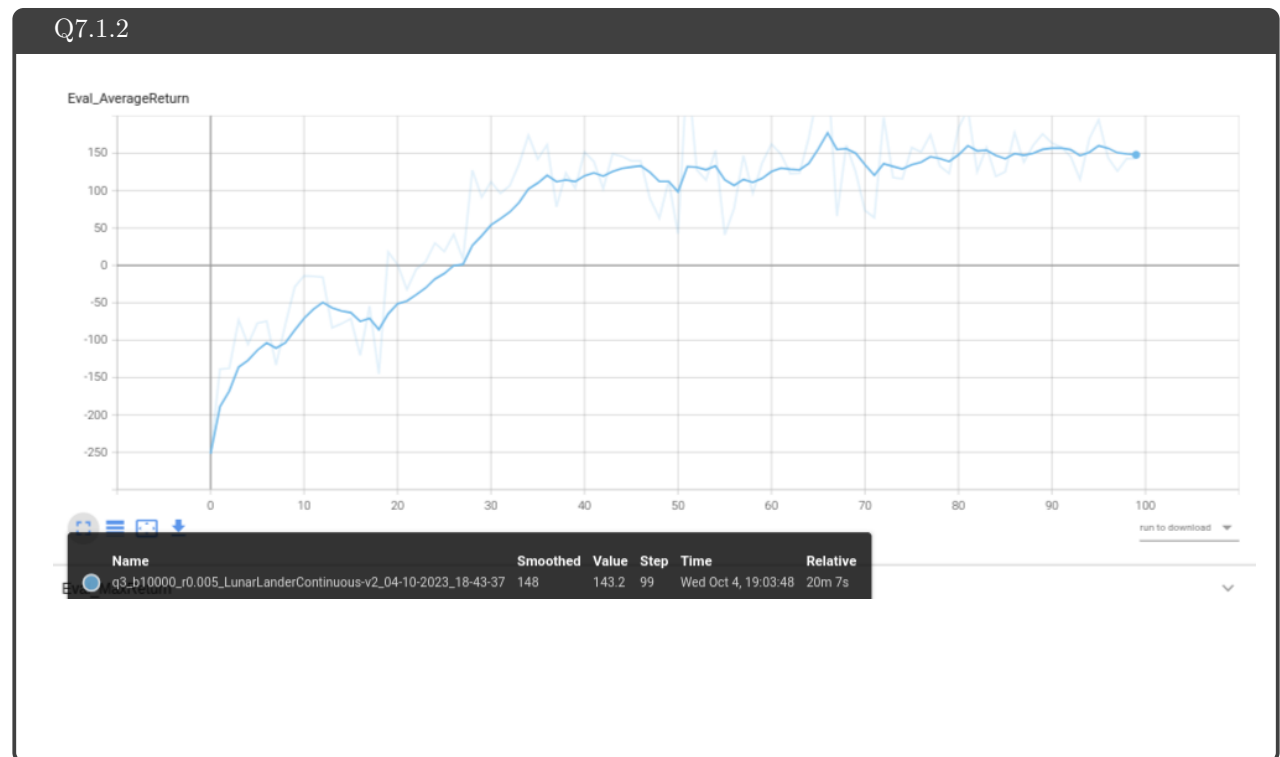# 7    More Complex Experiments

## 7.1    Experiment 3 (LunarLander) – [10 points total]

### 7.1.1    Configurations

---

Q7.1.1

```
python rob831/scripts/run_hw2.py \
    --env_name LunarLanderContinuous-v4 --ep_len 1000
    --discount 0.99 -n 100 -l 2 -s 64 -b 10000 -lr 0.005 \
    --reward_to_go --nn_baseline --exp_name q3_b10000_r0.005
```

### 7.1.2    Plot – [10 points]
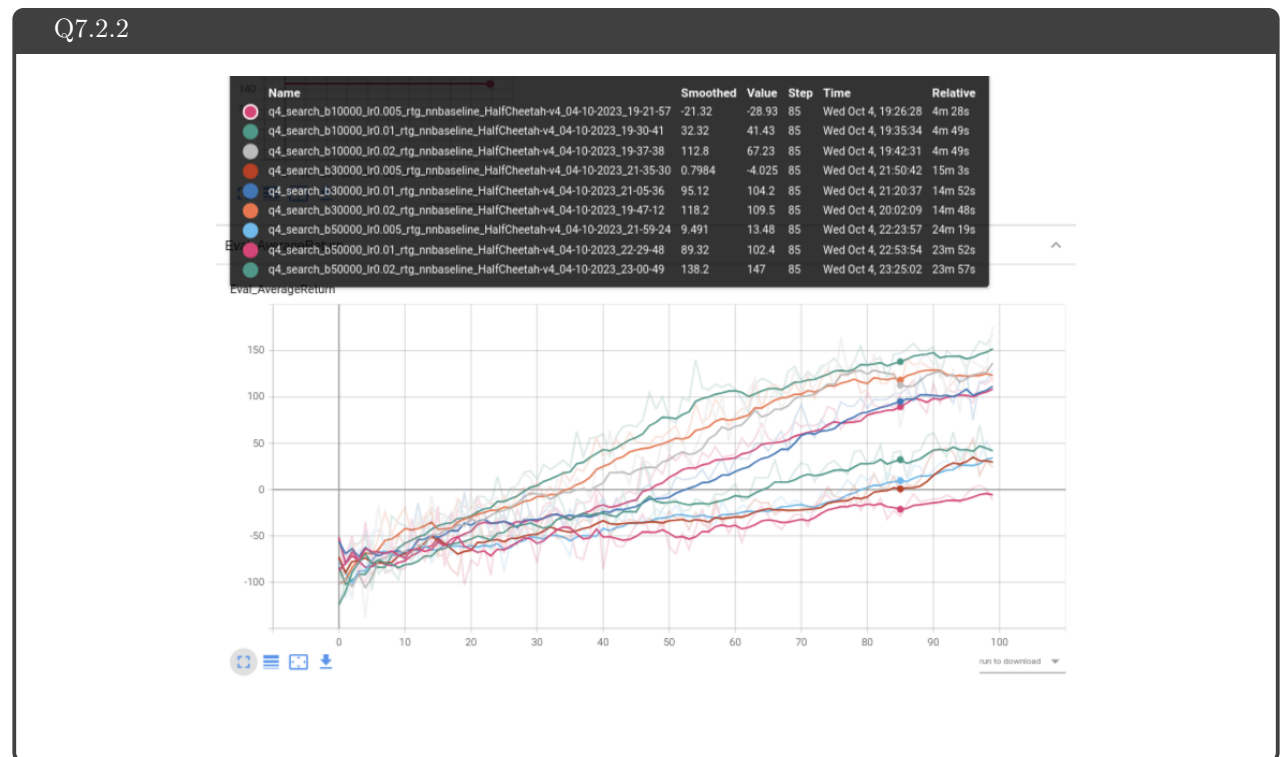


## 7.2    Experiment 4 (HalfCheetah) – [30 points]

### 7.2.1    Configurations

Q7.2.1

```
# b ∈ [10000, 30000, 50000], r ∈ [0.005, 0.01, 0.02]
python rob831/scripts/run_hw2.py --env_name HalfCheetah-v4 --ep_len 150 \
    --discount 0.95 -n 100 -l 2 -s 32 -b <b> -lr <r> -rtg --nn_baseline \
    --exp_name q4_search_b<b>_lr<r>_rtg_nnbaseline
```

### 7.2.2 Plot – [10 points]

**Q7.2.2**



| Name | Smoothed | Value | Step | Time | Relative |
|---|---|---|---|---|---|
| q4_search_b10000_lr0.005_rtg_nnbaseline_HalfCheetah-v4_04-10-2023_19-21-57 | -21.32 | -28.93 | 85 | Wed Oct 4, 19:26:28 | 4m 28s |
| q4_search_b10000_lr0.01_rtg_nnbaseline_HalfCheetah-v4_04-10-2023_19-30-41 | 32.32 | 41.43 | 85 | Wed Oct 4, 19:35:34 | 4m 49s |
| q4_search_b10000_lr0.02_rtg_nnbaseline_HalfCheetah-v4_04-10-2023_19-37-38 | 112.8 | 67.23 | 85 | Wed Oct 4, 19:42:31 | 4m 49s |
| q4_search_b30000_lr0.005_rtg_nnbaseline_HalfCheetah-v4_04-10-2023_21-35-30 | 0.7984 | -4.025 | 85 | Wed Oct 4, 21:50:42 | 15m 3s |
| q4_search_b30000_lr0.01_rtg_nnbaseline_HalfCheetah-v4_04-10-2023_21-05-36 | 95.12 | 104.2 | 85 | Wed Oct 4, 21:20:37 | 14m 52s |
| q4_search_b30000_lr0.02_rtg_nnbaseline_HalfCheetah-v4_04-10-2023_19-47-12 | 118.2 | 109.5 | 85 | Wed Oct 4, 20:02:09 | 14m 48s |
| q4_search_b50000_lr0.005_rtg_nnbaseline_HalfCheetah-v4_04-10-2023_21-59-24 | 9.491 | 13.48 | 85 | Wed Oct 4, 22:23:57 | 24m 19s |
| q4_search_b50000_lr0.01_rtg_nnbaseline_HalfCheetah-v4_04-10-2023_22-29-48 | 89.32 | 102.4 | 85 | Wed Oct 4, 22:53:54 | 23m 52s |
| q4_search_b50000_lr0.02_rtg_nnbaseline_HalfCheetah-v4_04-10-2023_23-00-49 | 138.2 | 147 | 85 | Wed Oct 4, 23:25:02 | 23m 57s |

### 7.2.3 Optimal b* and r* – [3 points]

**Q7.2.3**

Optimal b* = 50000, r* = 0.02

### 7.2.4 Describe how b* and r* affect task performance – [7 points]

**Q7.2.4**

It can be observed that fixing a batch size and increasing learning rate improves the performance. Similarly, increasing the batch size with a constant learning rate also improves the performance.

### 7.2.5   Configurations with optimal b* and r* – [3 points]

**Q7.2.5**

```
python rob831/scripts/run_hw2.py --env_name HalfCheetah-v4 --ep_len 150 \
    --discount 0.95 -n 100 -l 2 -s 32 -b <b*> -lr <r*> \
    --exp_name q4_b<b*>_r<r*>

python rob831/scripts/run_hw2.py --env_name HalfCheetah-v4 --ep_len 150 \
    --discount 0.95 -n 100 -l 2 -s 32 -b <b*> -lr <r*> -rtg \
    --exp_name q4_b<b*>_r<r*>_rtg

python rob831/scripts/run_hw2.py --env_name HalfCheetah-v4 --ep_len 150 \
    --discount 0.95 -n 100 -l 2 -s 32 -b <b*> -lr <r*> --nn_baseline \
    --exp_name q4_b<b*>_r<r*>_nnbaseline

python rob831/scripts/run_hw2.py --env_name HalfCheetah-v4 --ep_len 150 \
    --discount 0.95 -n 100 -l 2 -s 32 -b <b*> -lr <r*> -rtg --nn_baseline \
    --exp_name q4_b<b*>_r<r*>_rtg_nnbaseline
```
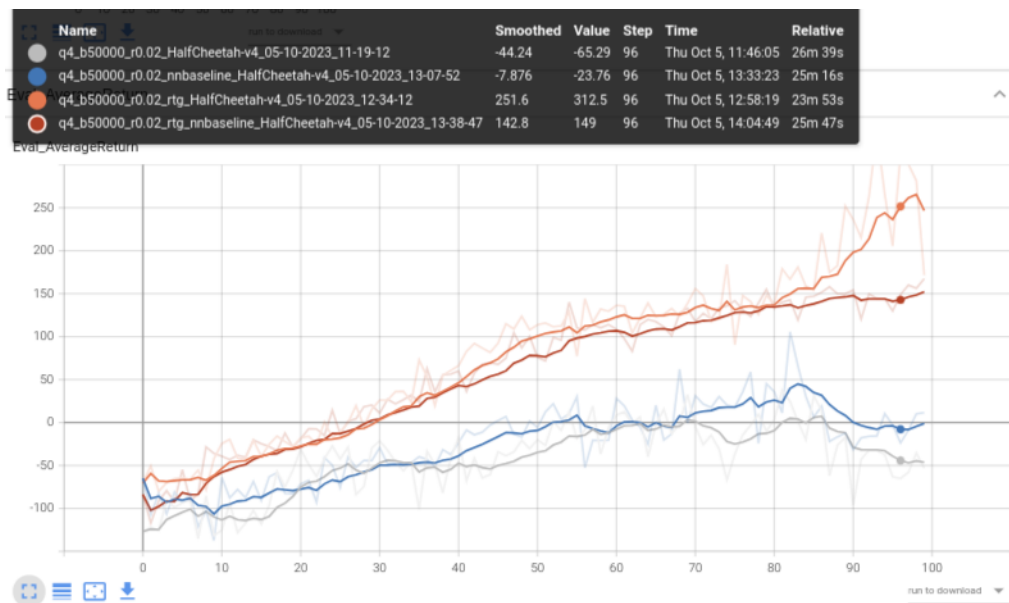
### 7.2.6   Plot for four runs with optimal b* and r* – [7 points]

**Q7.2.6**



# 8   Implementing Generalized Advantage Estimation
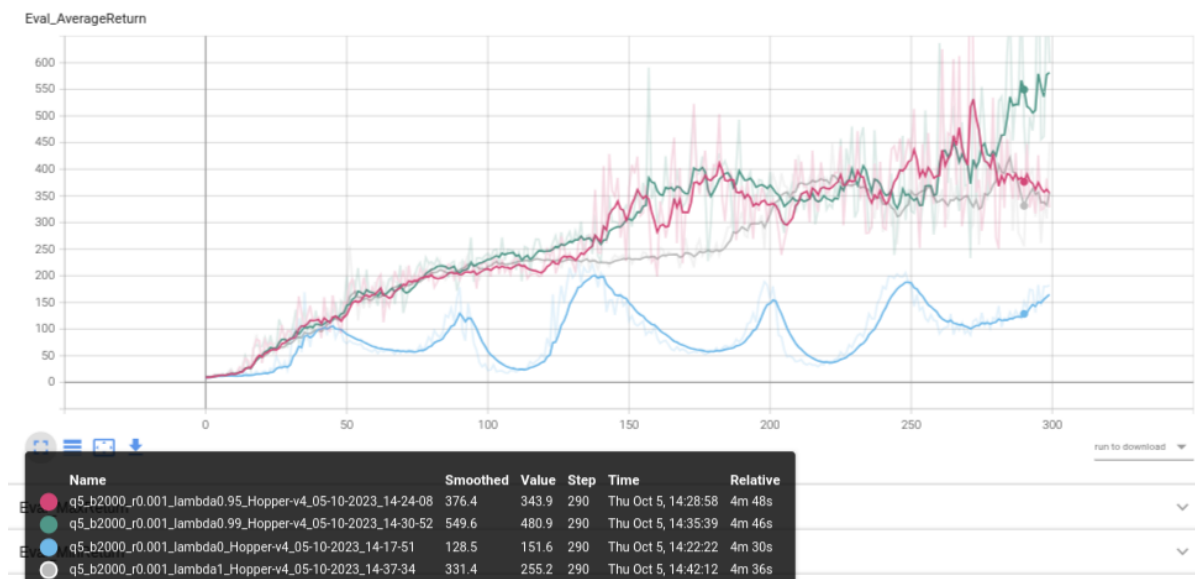
## 8.1 Experiment 5 (Hopper) – [20 points]

### 8.1.1 Configurations

**Q8.1.1**

```
# λ ∈ [0, 0.95, 0.99, 1]
python rob831/scripts/run_hw2.py \
    --env_name Hopper-v4 --ep_len 1000
    --discount 0.99 -n 300 -l 2 -s 32 -b 2000 -lr 0.001 \
    --reward_to_go --nn_baseline --action_noise_std 0.5 --gae_lambda <λ> \
    --exp_name q5_b2000_r0.001_lambda<λ>
```

### 8.1.2 Plot – [13 points]

**Q8.1.2**



Eval_AverageReturn

| Name | Smoothed | Value | Step | Time | Relative |
|------|----------|-------|------|------|----------|
| q5_b2000_r0.001_lambda0.95_Hopper-v4_05-10-2023_14-24-08 | 376.4 | 343.9 | 290 | Thu Oct 5, 14:28:58 | 4m 48s |
| q5_b2000_r0.001_lambda0.99_Hopper-v4_05-10-2023_14-30-52 | 549.6 | 480.9 | 290 | Thu Oct 5, 14:35:39 | 4m 46s |
| q5_b2000_r0.001_lambda0_Hopper-v4_05-10-2023_14-17-51 | 128.5 | 151.6 | 290 | Thu Oct 5, 14:22:22 | 4m 30s |
| q5_b2000_r0.001_lambda1_Hopper-v4_05-10-2023_14-37-34 | 331.4 | 255.2 | 290 | Thu Oct 5, 14:42:12 | 4m 36s |

### 8.1.3 Describe how λ affects task performance – [7 points]

**Q8.1.3**

As taught in class, $\lambda$ serves as a control for the bias-variance tradeoff, where increasing $\lambda$ decreases bias and increases variance. It can be seen that $\lambda = 0$ does not learn well. Setting $\lambda$ to 0.95 and 0.99 gives good results with 0.99 being the best in practice.

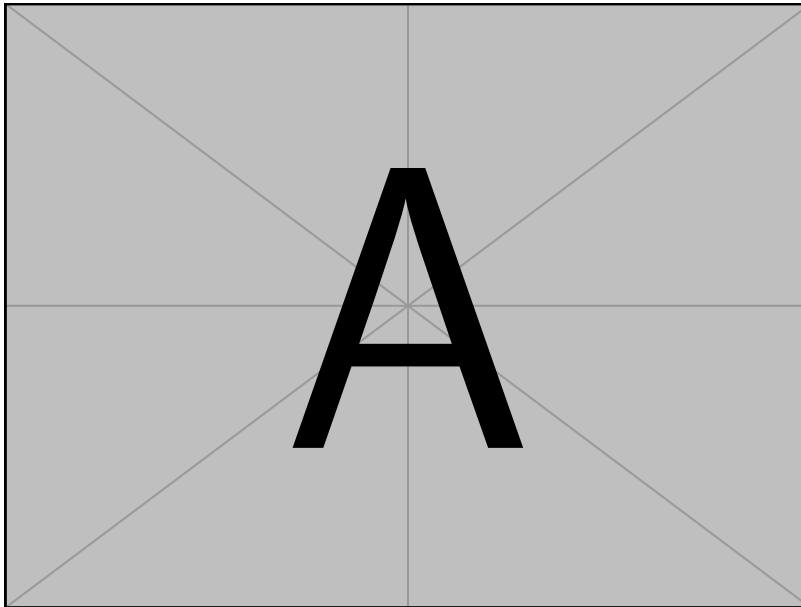# 9   Bonus! (optional)

## 9.1   Parallelization – [15 points]

> **Q9.1**
>
> Difference in training time:
>
> ```
> python rob831/scripts/run_hw2.py \
> ```

## 9.2   Multiple gradient steps – [5 points]

> **Q9.1**
>
> 
>
> ```
> python rob831/scripts/run_hw2.py \
> ```