

## Assignment 2: Policy Gradient

Andrew ID: mukaiy

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

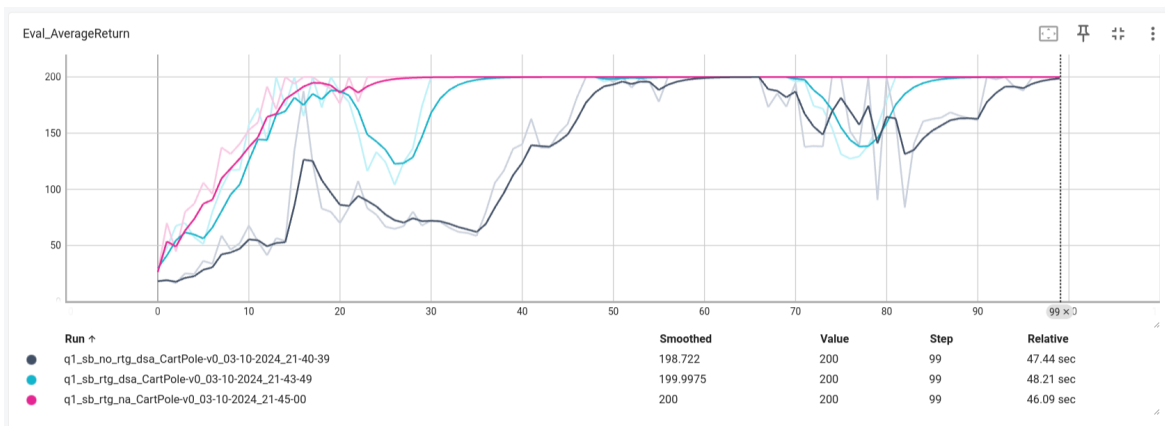
python rob831/scripts/run_hw2.py --env_name CartPole-v0 -n 100 -b 5000 \
-rtg -dsa --exp_name q1_lb_rtg_dsa

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

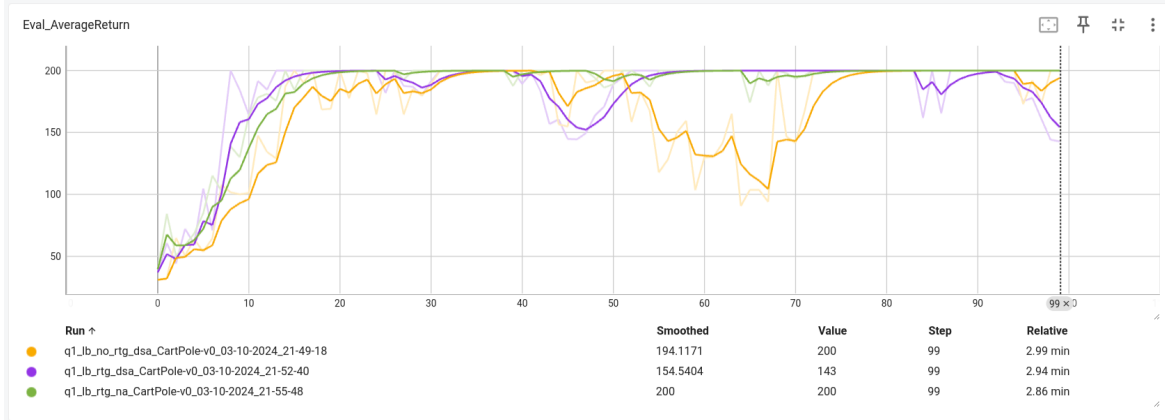
###### 5.1.2.1 Small batch – [5 points]

###### Q5.1.2.1



### 5.1.2.2 Large batch – [5 points]

#### Q5.1.2.2



### 5.1.3 Analysis

#### 5.1.3.1 Value estimator – [5 points]

##### Q5.1.3.1

Obviously reward-to-go because of causality and smaller variance.

#### 5.1.3.2 Advantage standardization – [5 points]

##### Q5.1.3.2

Yes, it does help. It helps to stabilize the learning process by reducing the variance of the advantage estimates.

One thing I don't understand is why changing the bias has no significant impact on training, since standardization changes the mean. Each reward is a function of both state and action, so moving the mean shouldn't be just like adding a state-dependent bias, e.g.,  $V(s)$ , right?

I guess it's because this process only "slightly" biases the advantage estimates, pretty much like the effect of not having accurate  $V(s)$  in the actor-critic method.

### 5.1.3.3 Batch size – [5 points]

#### Q5.1.3.3

I would say that batch size marginally helps to improve the policy faster, and I think the improvement is non-linear, which is why we don't see the time-to-plateau gets  $\frac{1}{5}$ ed.

But more data does help the experiment without reward-to-go to converge faster, since it typically needs more rollouts, i.e., more state/action/reward tuples, because it typically has higher variance.

## 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 10000 -lr 0.001 -rtg \
--exp_name q2_b_10000_r_0.001

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 0.001 -rtg \
--exp_name q2_b_1000_r_0.001

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 0.01 -rtg \
--exp_name q2_b_1000_r_0.01

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 0.05 -rtg \
--exp_name q2_b_1000_r_0.05

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 0.1 -rtg \
--exp_name q2_b_1000_r_0.1

python rob831/scripts/run_hw2.py --env_name InvertedPendulum-v4 \
--ep_len 1000 --discount 0.9 -n 100 -l 2 -s 64 -b 100 -lr 0.1 -rtg \
--exp_name q2_b_100_r_0.1

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 0.01 -rtg \
--exp_name q2_b_2000_r_0.01

python rob831/scripts/run_hw2.py --env_name InvertedPendulum-v4 \
--ep_len 1000 --discount 0.9 -n 100 -l 2 -s 64 -b 500 -lr 0.01 -rtg \
--exp_name q2_b_500_r_0.01

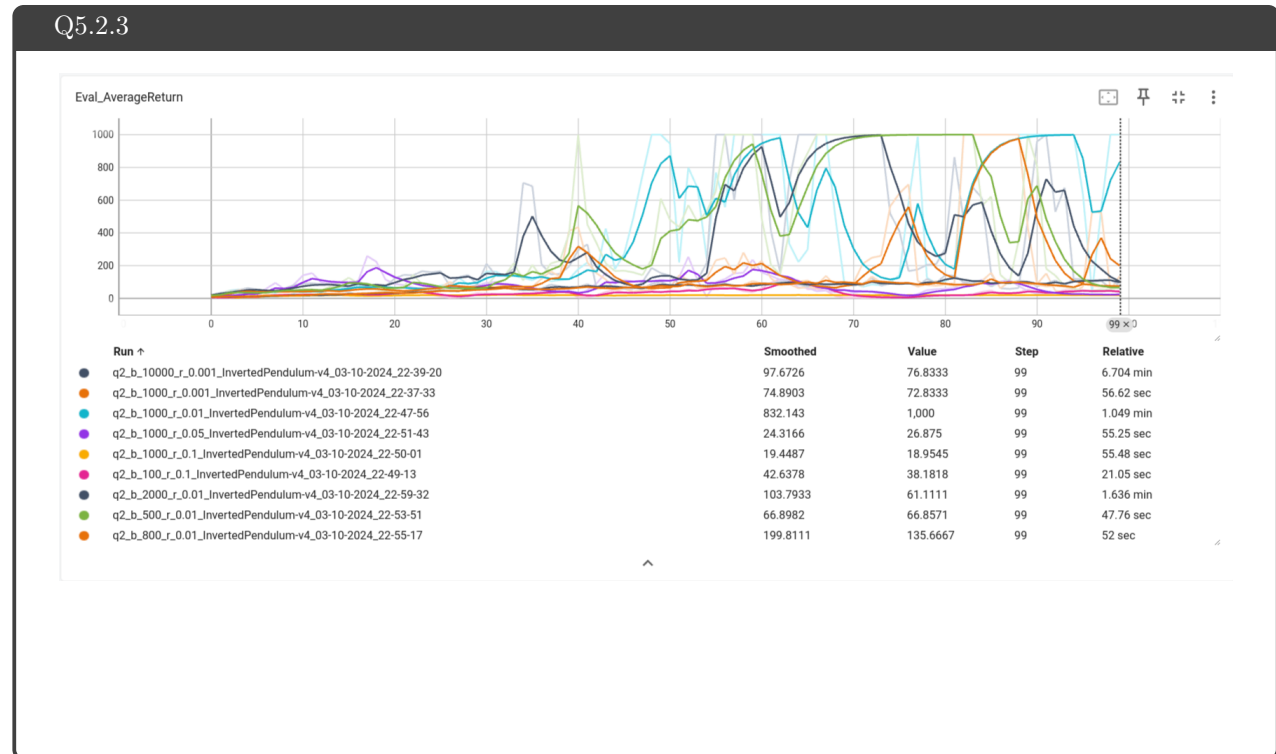
python rob831/scripts/run_hw2.py --env_name InvertedPendulum-v4 \
--ep_len 1000 --discount 0.9 -n 100 -l 2 -s 64 -b 800 -lr 0.01 -rtg \
--exp_name q2_b_800_r_0.01
```

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

#### Q5.2.2

500  
0.01

### 5.2.3 Plot – [5 points]



## 7 More Complex Experiments

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

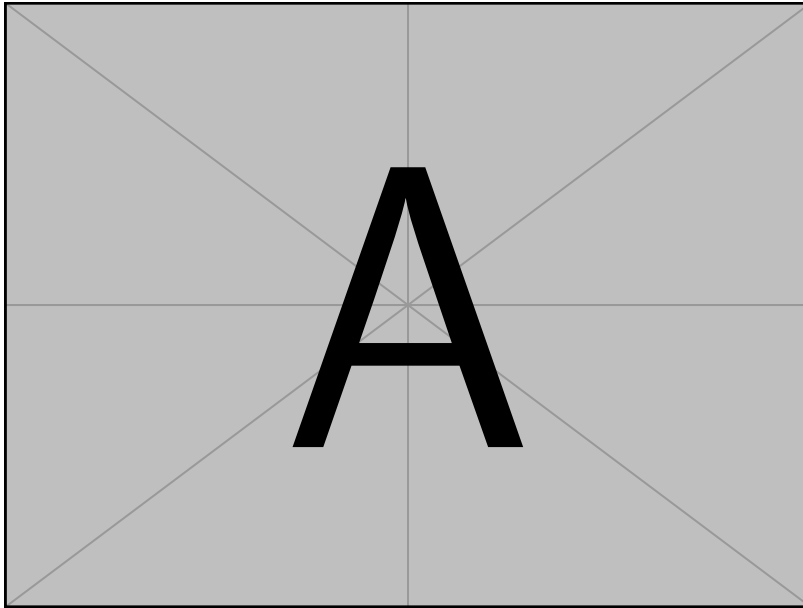
#### 7.1.1 Configurations

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]

Q7.1.2

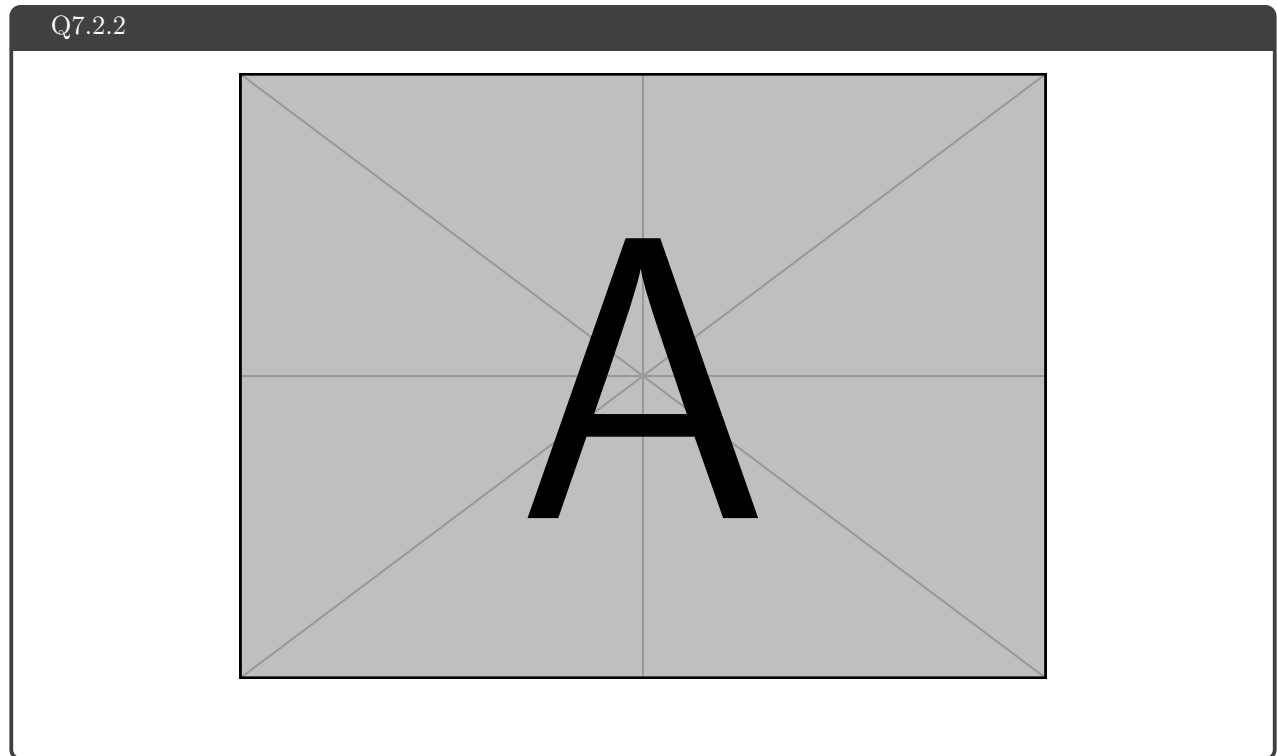


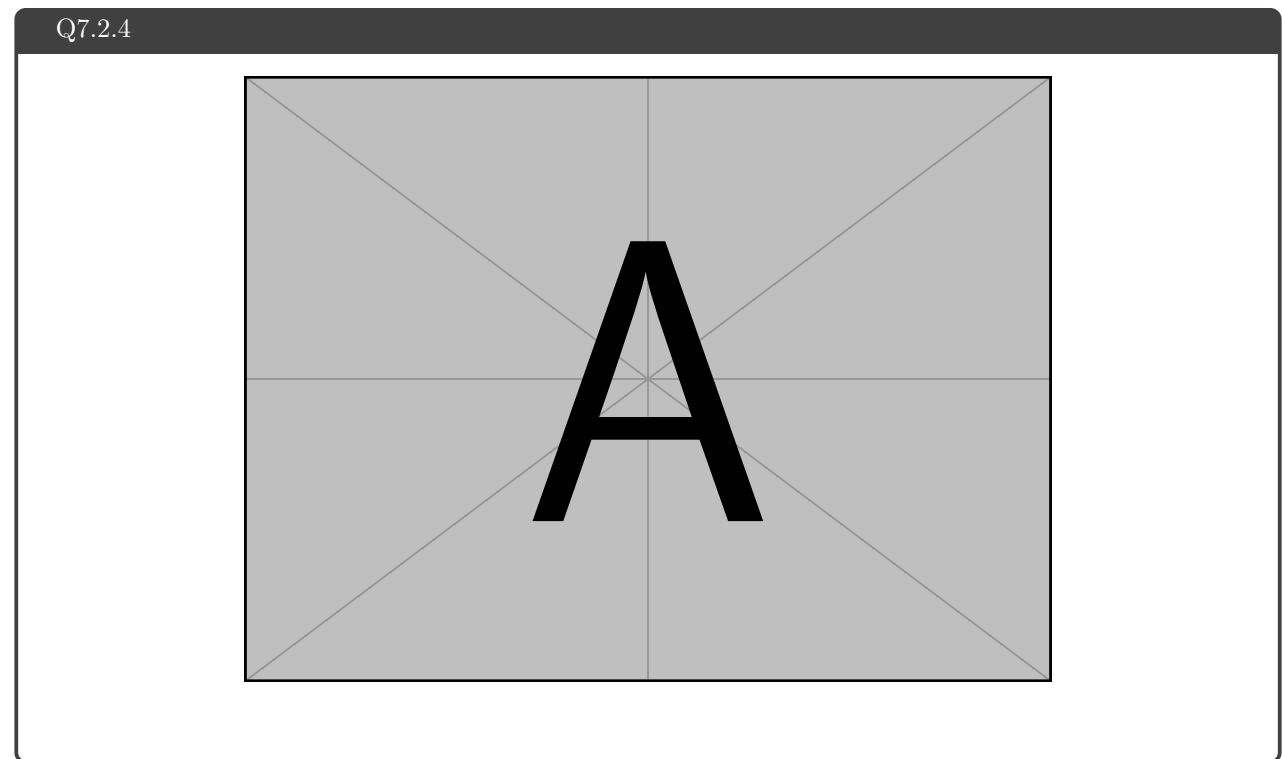
## 7.2 Experiment 4 (HalfCheetah) – [30 points]

### 7.2.1 Configurations

Q7.2.1

```
python rob831/scripts/run_hw2.py --env_name HalfCheetah-v4 --ep_len 150 \  
  --discount 0.95 -n 100 -l 2 -s 32 -b 10000 -lr 0.02 \  
  --exp_name q4_search_b10000_lr0.02  
python rob831/scripts/run_hw2.py --env_name HalfCheetah-v4 --ep_len 150 \  
  --discount 0.95 -n 100 -l 2 -s 32 -b 10000 -lr 0.02 -rtg \  
  --exp_name q4_search_b10000_lr0.02_rtg  
python rob831/scripts/run_hw2.py --env_name HalfCheetah-v4 --ep_len 150 \  
  --discount 0.95 -n 100 -l 2 -s 32 -b 10000 -lr 0.02 --nn_baseline \  
  --exp_name q4_search_b10000_lr0.02_nnbaseline  
python rob831/scripts/run_hw2.py --env_name HalfCheetah-v4 --ep_len 150 \  
  --discount 0.95 -n 100 -l 2 -s 32 -b 10000 -lr 0.02 -rtg --nn_baseline \  
  --exp_name q4_search_b10000_lr0.02_rtg_nnbaseline
```

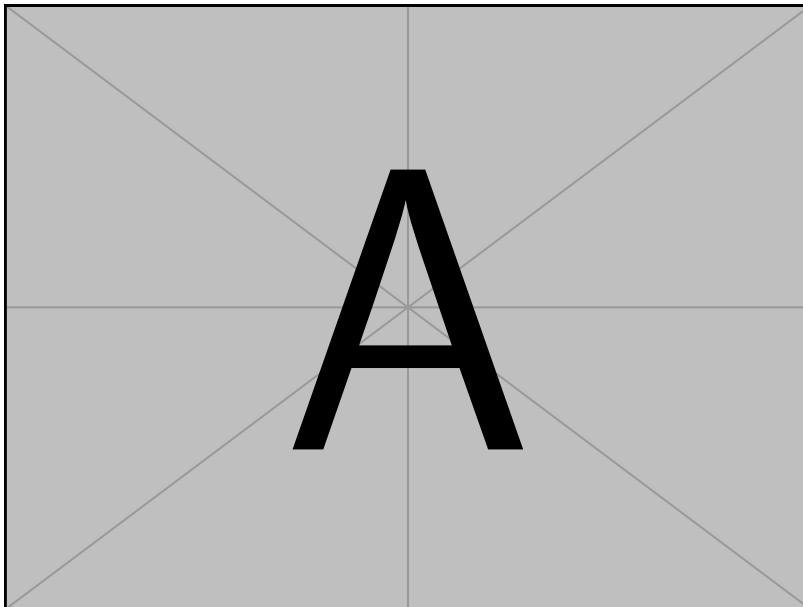
**7.2.2 Plot – [10 points]****7.2.3 (Optional) Optimal  $b^*$  and  $r^*$  – [3 points]**

**7.2.4 (Optional) Plot – [10 points]****7.2.5 (Optional) Describe how  $b^*$  and  $r^*$  affect task performance – [7 points]**

Q7.2.5

**7.2.6 (Optional) Configurations with optimal  $b^*$  and  $r^*$  – [3 points]****Q7.2.6**

```
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
```

**7.2.7 (Optional) Plot for four runs with optimal  $b^*$  and  $r^*$  – [7 points]****Q7.2.7****8 Implementing Generalized Advantage Estimation**



## 8.1 Experiment 5 (Hopper) – [20 points]

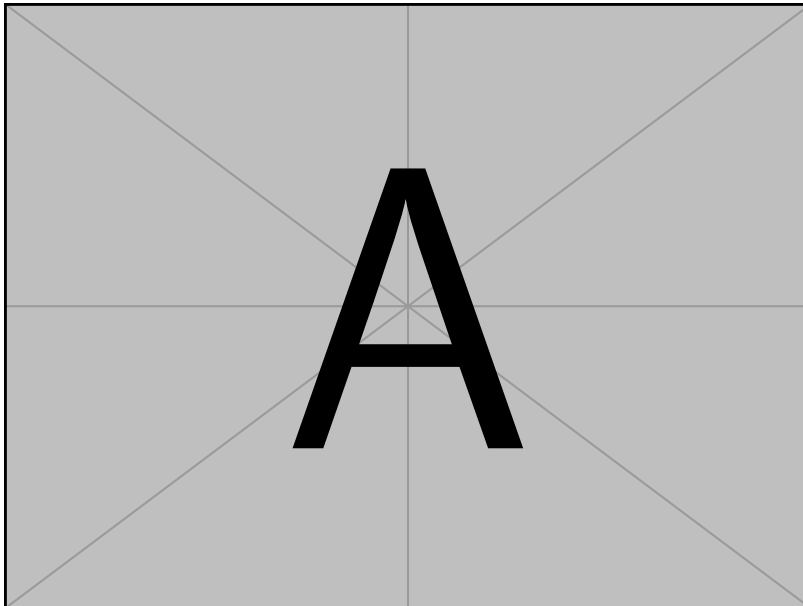
### 8.1.1 Configurations

#### Q8.1.1

```
#  $\lambda \in [0, 0.95, 0.99, 1]$ 
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 < $\lambda$ > \
  --exp_name q5_b2000_r0.001_lambda< $\lambda$ >
```

### 8.1.2 Plot – [13 points]

#### Q8.1.2



### 8.1.3 Describe how $\lambda$ affects task performance – [7 points]

#### Q8.1.3

## 9 Bonus! (optional)

### 9.1 Parallelization – [15 points]

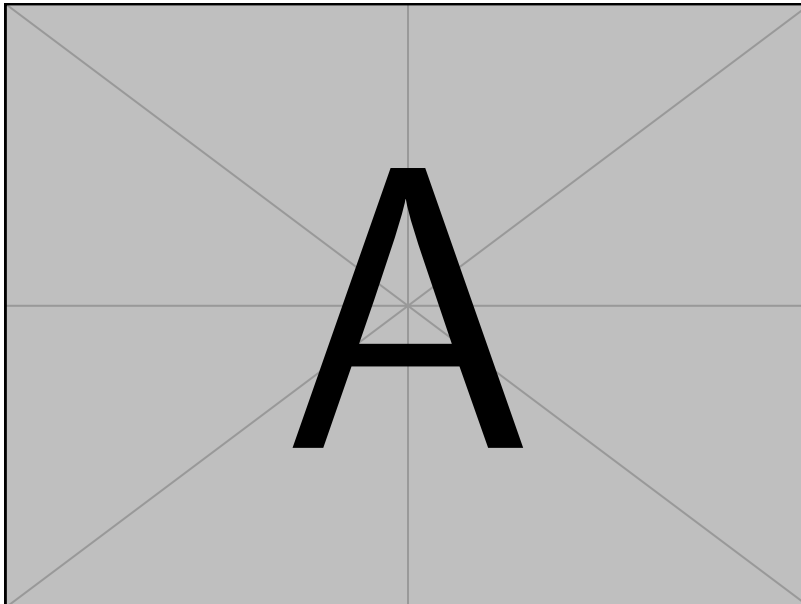
Q9.1

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 \
```