

Assignment 2: Policy Gradient

Andrew ID: ayanovic

Collaborators: Write the Andrew IDs of your collaborators here (if any).

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

I used the following commands to run the experiments (same as in the assignment instructions):

```
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

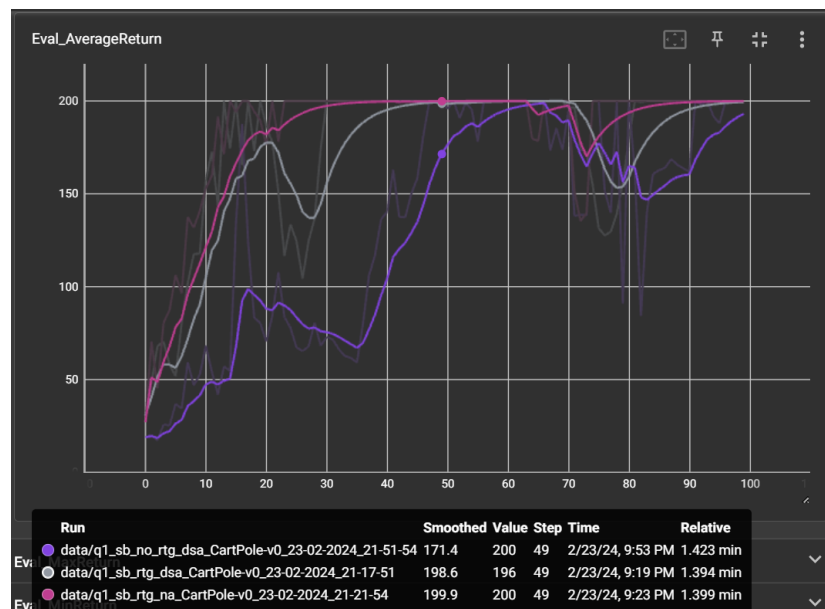
python rob831/scripts/run_hw2.py --env_name CartPole-v0 -n 100 -b 5000 \
    -rtg -dsa --exp_name q1_lb_rtg_dsa

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

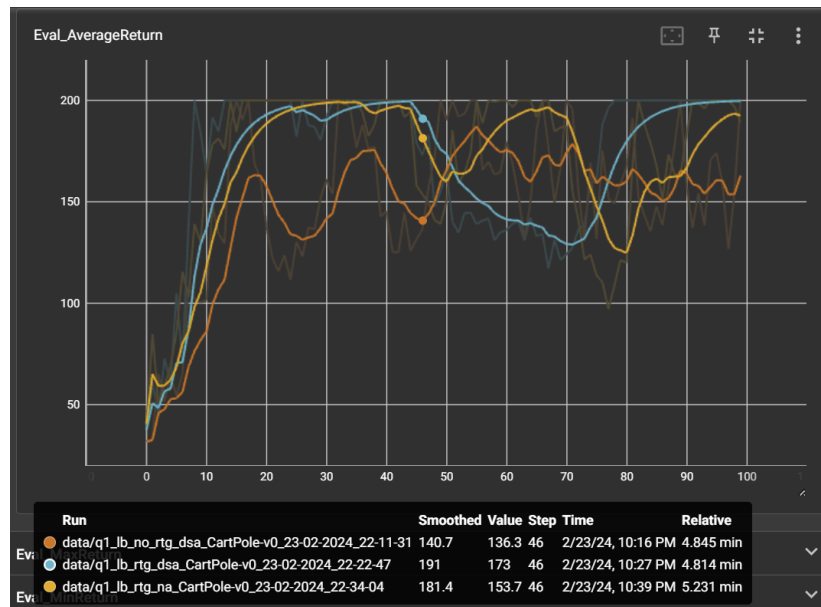
5.1.2.1 Small batch – [5 points]

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

The reward-to-go estimator performs better than the non-reward-to-go estimator.

5.1.3.2 Advantage standardization – [5 points]

Q5.1.3.2

The advantage standardization performs better than the non-standardized advantage.

5.1.3.3 Batch size – [5 points]

Q5.1.3.3

The larger batch size appears to perform better by reducing average variance of the large-batch runs.

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

5.2.1 Configurations – [5 points]

Q5.2.1

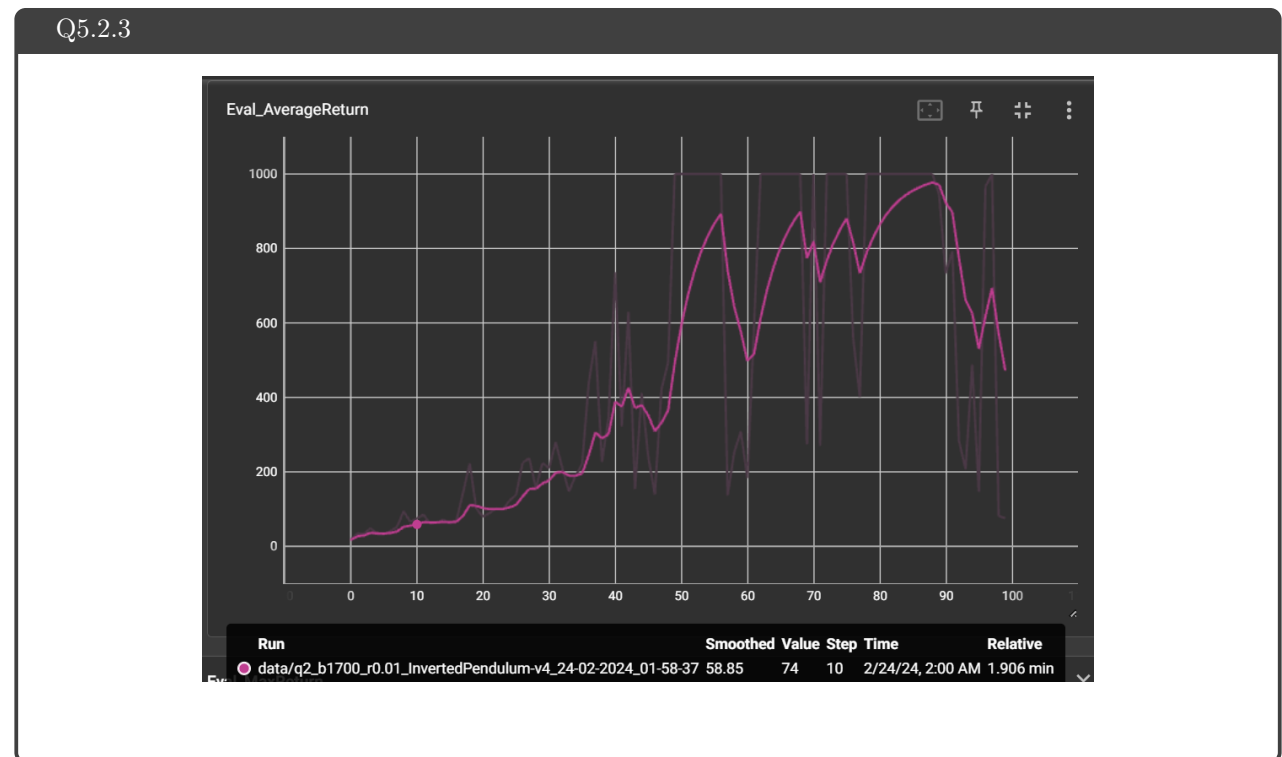
```
python rob831/scripts/run_hw2.py --env_name InvertedPendulum-v4 \
--ep_len 1000 --discount 0.9 -n 100 -l 2 -s 64 -b <b*> -lr <r*> -rtg \
--exp_name q2_b<b*>_r<r*>
python rob831/scripts/run_hw2.py --env_name InvertedPendulum-v4 \
--ep_len 1000 --discount 0.9 -n 100 -l 2 -s 64 -b 1000 -lr 0.001 \
--exp_name q2_b1000_r0.001
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 0.005 \
--exp_name q2_b3000_r0.005
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 0.01 \
--exp_name q2_b3000_r0.01
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 0.05 \
--exp_name q2_b3000_r0.05
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 0.007 \
--exp_name q2_b2000_r0.007
python rob831/scripts/run_hw2.py --env_name InvertedPendulum-v4 \
--ep_len 1000 --discount 0.9 -n 100 -l 2 -s 64 -b 1700 -lr 0.007 \
--exp_name q2_b1700_r0.007
python rob831/scripts/run_hw2.py --env_name InvertedPendulum-v4 \
--ep_len 1000 --discount 0.9 -n 100 -l 2 -s 64 -b 1700 -lr 0.009 \
--exp_name q2_b1700_r0.009
python rob831/scripts/run_hw2.py --env_name InvertedPendulum-v4 \
--ep_len 1000 --discount 0.9 -n 100 -l 2 -s 64 -b 1700 -lr 0.01 \
--exp_name q2_b1700_r0.01
```

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

Q5.2.2

batch size=1700, learning rate=0.01 Note: Although the graph shows sudden drops in performance, this training was able to reach the maximum reward of 1000.

5.2.3 Plot – [5 points]



7 More Complex Experiments

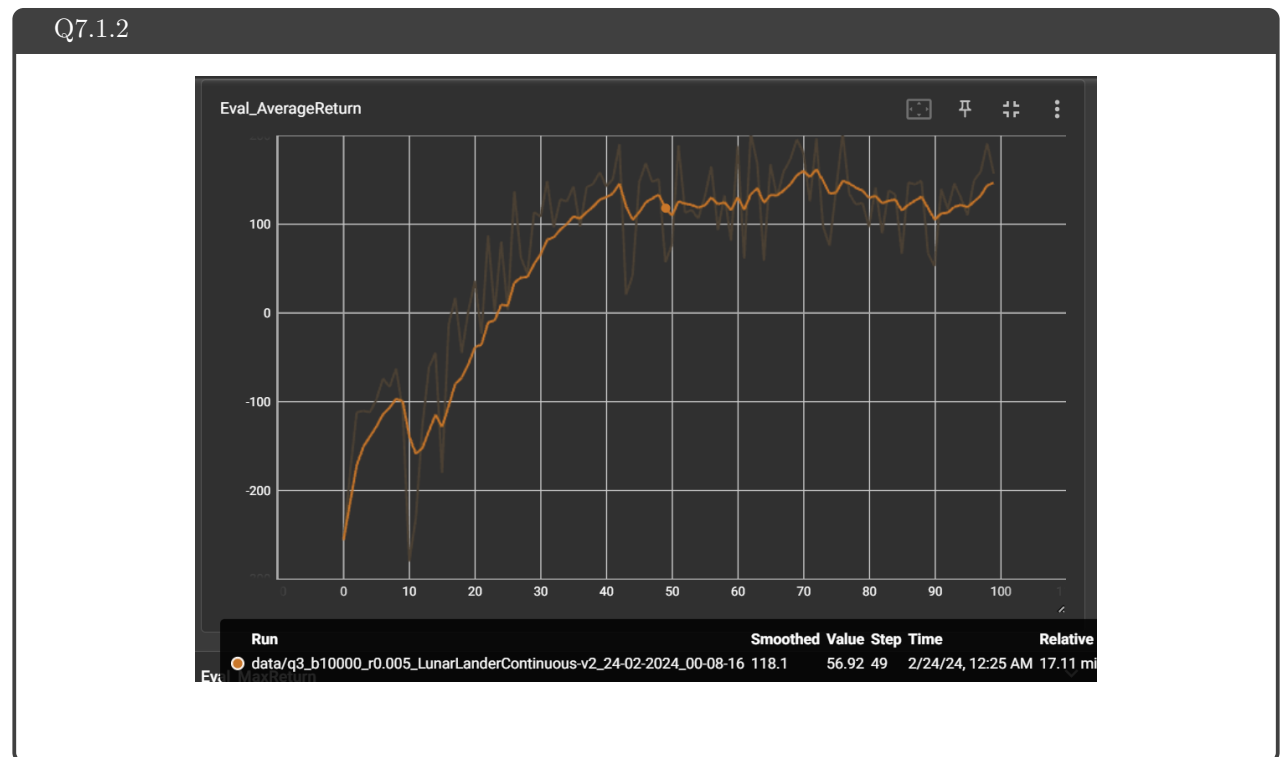
7.1 Experiment 3 (LunarLander) – [10 points total]

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]



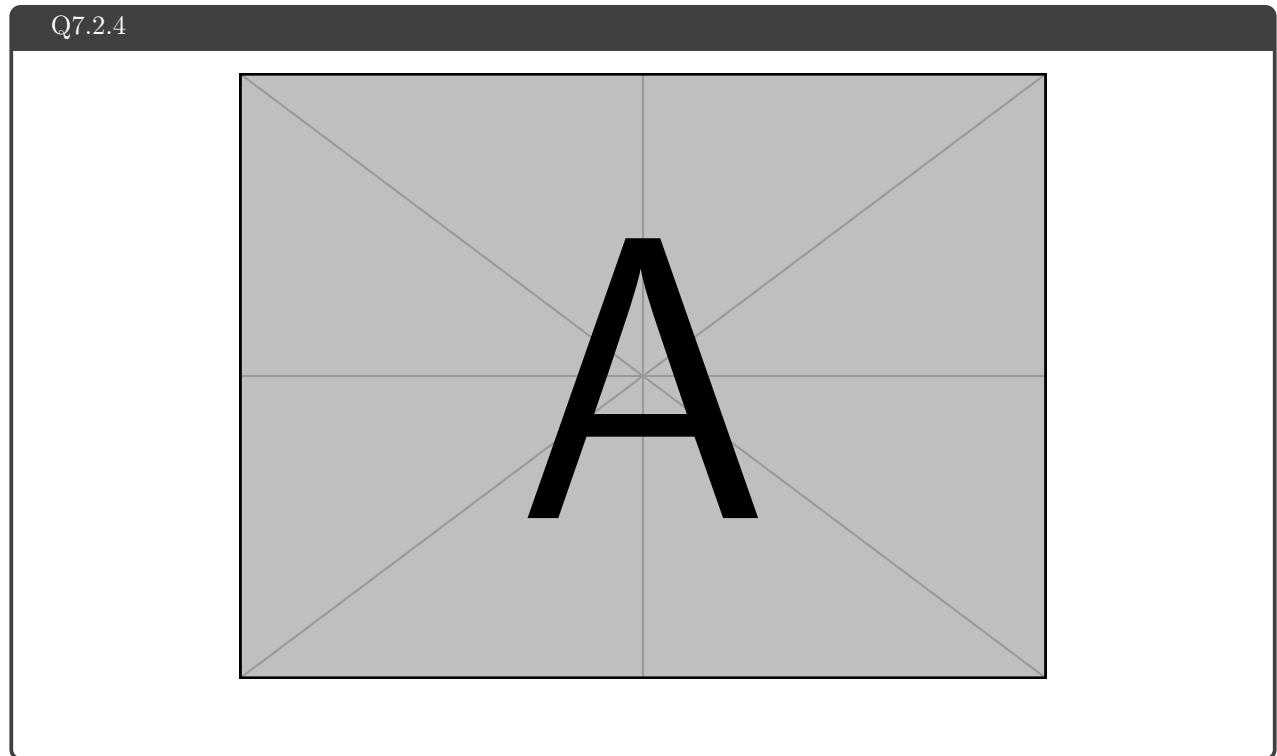
7.2 Experiment 4 (HalfCheetah) – [30 points]

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
```

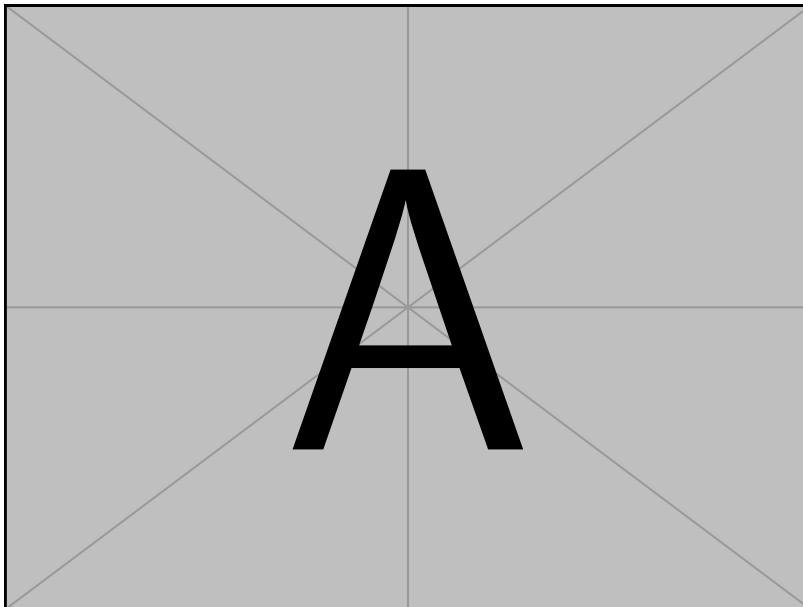
7.2.2 Plot – [10 points]**7.2.3 (Optional) Optimal b^* and r^* – [3 points]**

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]

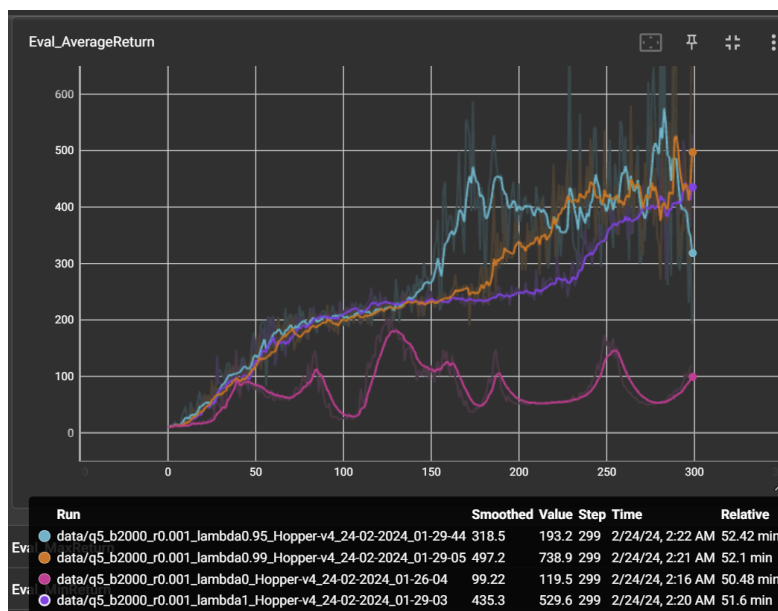
8.1.1 Configurations

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 λ affects task performance – [7 points]

Q8.1.3

The λ value affects the task performance and variance of the return. At $\lambda = 0$, the return is the lowest, and does not reach good performance. For $\lambda = 0.95$, the return is higher, but the variance is also higher. For $\lambda = 0.99$, the return is highest, and the variance is lower on average compared to $\lambda = 0.95$. For $\lambda = 1$, the return is slightly lower than $\lambda = 0.99$, however the variance is the lowest, which is unexpected as it is supposed to be equivalent to the vanilla neural network baseline estimator without any variance reduction. Moreover, when disabling the GAE estimator, the return is equivalent on average to $\lambda = 1$.

9 Bonus! (optional)

9.1 Parallelization – [15 points]

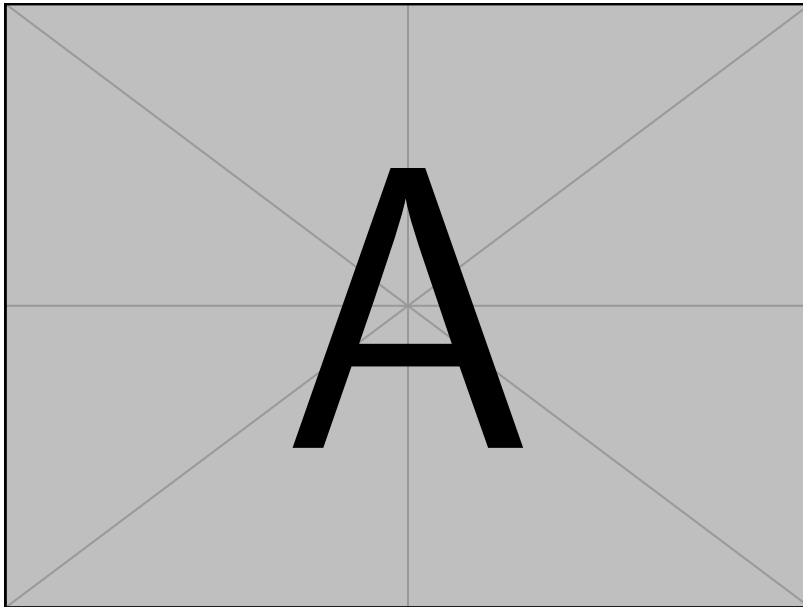
Q9.1

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 \
```