
Application of RL to Simulate Performance of an Unguided Hypersonic Glide Vehicle

Justine John A Serdoncillo
Aerospace Engineering and Mechanics
University of Minnesota - Twin Cities
Minneapolis, MN 55414
serdo004@umn.edu

Darryl Williams
Aerospace Engineering and Mechanics
University of Minnesota - Twin Cities
Minneapolis, MN 55414
will7322@umn.edu

Abstract

In this paper, we present a novel approach for simulating and controlling a high-speed aircraft using simplified physics and reinforcement learning (RL). We focus on the case of an unguided hypersonic glide vehicle (HGV), which is a challenging and relevant problem for aerospace engineering. Our goal is to design a control policy that can keep the vehicle on a straight heading, without relying on analytical expressions or numerical approximations of the dynamics. To achieve this, we use neural networks to model the continuous state space and the action space of the system, and apply RL algorithms to learn the optimal policy from data. We compare our approach with conventional methods and demonstrate its effectiveness and robustness in various scenarios. Our results show that RL can offer a viable and efficient alternative for controlling high-speed aircraft.

Keywords: Q-Learning, Neural Networks, Glide Vehicles

1 Problem Introduction

Hypersonic glide vehicles (HGVs) are capable of maneuvering within earth's atmosphere as they glide at high-speeds in excess of Mach 5. As a result of their speed, HGVs are able to travel several thousand kilometers within a short amount of time. Their ability to achieve long distances in short time frames and their use of aerodynamic forces for maneuvering have made them a subject of extensive research, particularly for defense applications. A typical trajectory for HGVs involves a boost phase that brings the vehicle high in the atmosphere, a reentry phase, and a glide phase at speeds of Mach 5-20.

In this paper we use a deep Q-learning approach with experience reply to perform guidance on an HGV during its glide phase. In the implemented framework, an aerodynamic code is integrated into the equation of motion over a non-rotating spherical Earth:

$$\frac{\partial}{\partial t} \begin{bmatrix} \gamma \\ V \\ x \\ h \end{bmatrix} = \begin{bmatrix} \frac{1}{V} \left[\frac{L}{m} + g \cos(\gamma) - \frac{V^2}{R} \cos(\gamma) \right] \\ \frac{D}{m} + g \sin(\gamma) \\ -V \cos(\gamma) \\ -V \sin(\gamma) \end{bmatrix}. \quad (1)$$

We then use an optimal policy π^* , to preform guidance on the HGV initial conditions of 60 km in altitude traveling at 6 km/s. Newtonian Aerodynamic (NA) was used as a lower-fidelity tool to compute the aerodynamic forces acting on the HGV. This approximate method is much cheaper than computational fluid dynamics enabling rapid exploration of the simulated environment but for lower levels of accuracy. Newtonian aerodynamics assumes tangential velocities are unchanged, while the normal velocity results in aerodynamic forces acting on the vehicle as shown in figure 1.

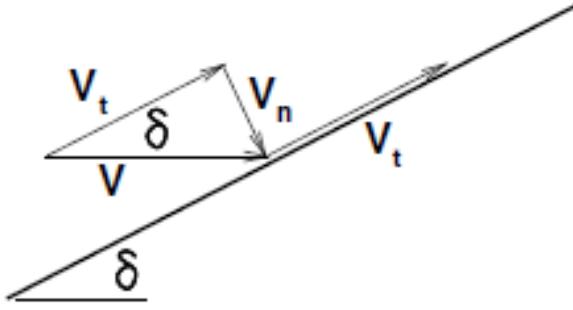


Figure 1: Illustration of tangential and normal velocities on surface.

This approximation is more accurate as speeds are increased into higher mach numbers. The algorithm to compute forces over a tessellated surface is described in the NA Procedure below.

```

 $N = \text{number of triangles elements}$ 
procedure NA( $\alpha$ )
   $i \leftarrow 1$ 
   $\hat{v} = \{\cos(\alpha), \sin(\alpha), 0\}$ 
  Forces =  $\{0, 0, 0\}$ 
  while  $i \leq N$  do
    if  $(\hat{v} \cdot \hat{n}_i) < 0$  then
       $cp_i = 2(\hat{v} \cdot \hat{n}_i)^2$ 
    else
       $cp_i = 0$ 
    end if ▷ Flow is hidden from velocity vector
     $P_i = (cp_i \cdot q) + P_\infty$ 
     $F = F + (P_i \cdot A_i) \cdot (-\hat{n}_i)$  ▷  $A_i$  computed from cross product
     $i \leftarrow i + 1$ 
  end while
   $D = F(2) \sin(\alpha) + F(1) \cos(\alpha)$ 
   $L = F(2) \cos(\alpha) - F(1) \sin(\alpha)$ 
end procedure

```

We then use the aerodynamic forces(L, D) to compute the derivatives in Equation 1. Time was integrated using the 4-th order Runge-Kutta method with a fixed time-step of $\Delta t = 1$ seconds as the time discretization.

2 Methodology

2.1 Simulation Environment

Our policy $\pi(a = \alpha|s)$, when given a state, selects an angle of attack (α) such that forces on the vehicle, lift (L) and drag (D), are modulated. With these forces, we can then integrate the equation of motion for a fixed δt to arrive at our new state. We note that the simulated environment is not representative of actual flight. Rather, this model for system dynamics is sought as a proof of concept and to have tractable training times of the algorithm used. The relative disparity between system dynamics comparing a computational fluid dynamics approach to NA is shown in Figure 3.

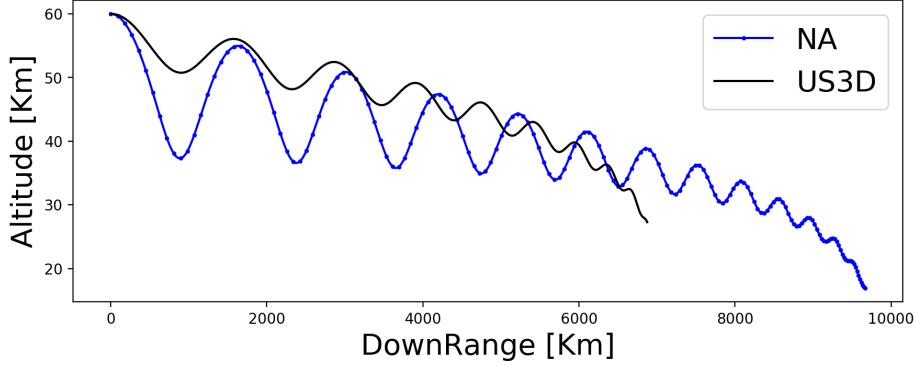


Figure 2: Comparison between full trajectory using US3D (CFD) and NA.

2.2 Markov Decision Process model

We implemented an epsilon greedy approach for exploration. Modeled as an episodic problem that terminated once the altitude went below 1 km or velocity decreased below 1 km/s. A summary of major hyperparameters and state space below:

- State Space: γ, h, V
- Action: α
- $\epsilon_{i+1} = 1 + \log(\epsilon_i / (0.6 * n))$
- Problem: Episodic with gamma = 0.9
- Rewards: $-\gamma^2 + d$

2.3 Neural Networks for continuous State Space RL

Artificial neural networks have been used for extensively for nonlinear function approximation. It characterized by applying a linear model to a feature vector then applying a non-linear mapping. This neural networks are often convoluted within each other its to have several instances of linear and non-linear mappings applied to them. This type of neural network architecture is known as a deep neural network. For the implemented method we decided on a neural network composed of two 24 neuron hidden layers, The non-linear activation decided was the use of the rectified linear unit.

In our case, we are using a deep Q-learning methodology by using an input tuple of the state, action, new state, and the reward. Our compute time was broken into the following:

- Exploration: first 25%
 - choosing random actions to explore the state transitions
 - store [state, action, state prime, reward] for each combination
- Training: second 65%
 - reduction of epsilon to reduce exploration
 - trains using stored memory
 - input is the state and action, output is the next state and reward
- Testing: last 10%
 - test neural net by running on new intial conditions and trajectories

The figure below shows the algorithm that we based our neural network on. The code for the environment and neural network as seen in the appendices below.

Algorithm 4 Deep Q-Learning with Experience Replay

```
1: Initialize replay memory  $D$  with capacity  $N$ 
2: Initialize  $Q$ -values with random weights
3: for  $t = 1 \dots T$  do
4:   With probability  $\epsilon$  select a random action  $a_t$ ,
5:   otherwise select  $a_t = \max_a Q^*(\phi(s_t), a, w_{t-1})$ 
6:   Execute action  $a_t$ 
7:   Receive reward  $r_t$  and new state  $s_{t+1}$ 
8:   Store transition  $(\phi(s_t), a_t, r_t, \phi(s_{t+1}))$  into  $D$ 
9:   Sample experience  $(\phi(s_j), a_j, r_j, \phi(s_{j+1}))$  from  $D$ 
10:  Set  $y_j = \begin{cases} r_j & \text{if } s_{j+1} \text{ is terminal state} \\ r_j + \gamma \max_{a_{j+1}} Q(\phi(s_{j+1}), a_{j+1}, w_{t-1}) & \text{otherwise} \end{cases}$ 
11:  Perform gradient descent update using  $[y_j - Q(\phi(s_j), a_j, w_t)]^2$  as loss
12: end for
```

Figure 3: Deep Q-Learning algorithm

Figure 4

3 Results

To tackle this project, there were lots of iterations that was needed in order to create a proper model environment and for it to achieve our desired objectives. At first, the initial parameters were used to simulate and run the environment but adding a negative reward at the end of the episode. At the end, our test cases showed the HGV to glide in a really flat trajectory and then descends down relatively fast. This seemed like what we initially wanted but from the data, it can be seen that it ended the episode early at around 1/10 of the original downrange using a constant angle of attack.

Due to the performance of the trajectory for the heading angle, thoughts on if the issue was more on the neural net was done. The episodes was increased to 1000 from 100, the gamma changed from 0.7 to 0.85 to highlight the importance of future actions so it doesn't try to crash earlier. Lastly due to the increase in episode number, the minibatch was increased to 64. After doing all of these changes, the effect on the trajectories was still the same so the focus changed not on the neural network but on the rewards structure.

With the same effect but much higher number of episodes, the number of episodes was greatly decreased to 100 again. At the same time, the gamma was increased to 0.95. To reward going higher distances, an extra reward of the downrange was added while also penalizing if the terminating downrange was below 1000km. This had better results, and incremental improvements was desired.

In order to address this, splits in the exploration, training and testing were done. To have a bigger source of memory for training, the exploration was increased to 35% of the episodes, while training account for 65% of the episodes. At the same time, some limitations were removed which was now just the altitude because a low velocity eventually gets a low altitude soon anyway. The number of episodes was increased to 250 with a minibatch of 32 as well. There were some further explorations with our rewards structure by being more lenient with the heading angle so an absolute value was used instead.

Lastly, due to a recent collaboration and ideation, adding generality for the controller would be great so there were now different initial conditions added so the state space may be greatly increased. At the same time, extra limitations were added due to the physical representation of the model. This should have been added earlier but was disregarded for the sake of exploration.

We now summarize the runs of the DNN Q-learning. Two of our runs are listed Figure 5. We see consistency that between the two figures, as the network progresses in training, the algorithm learns slowly to keep the heading angle near zero while maximizing range. Its evident by subfigure 5b that it has developed a policy for incentivizing max range. This is supported by the early failed regions in subfigure 5a where the trajectories did not exceed 1,000 km in down range distance.

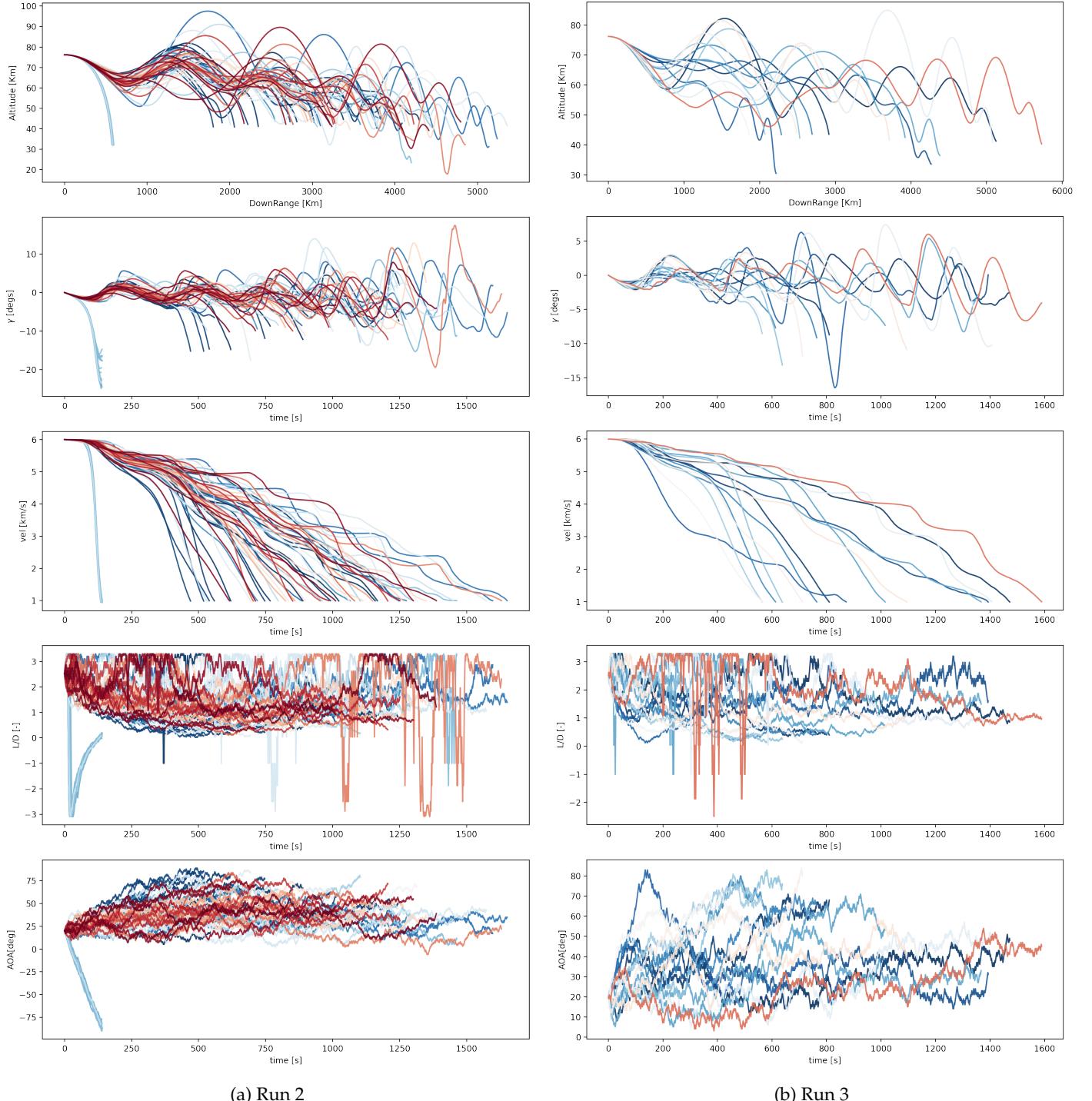
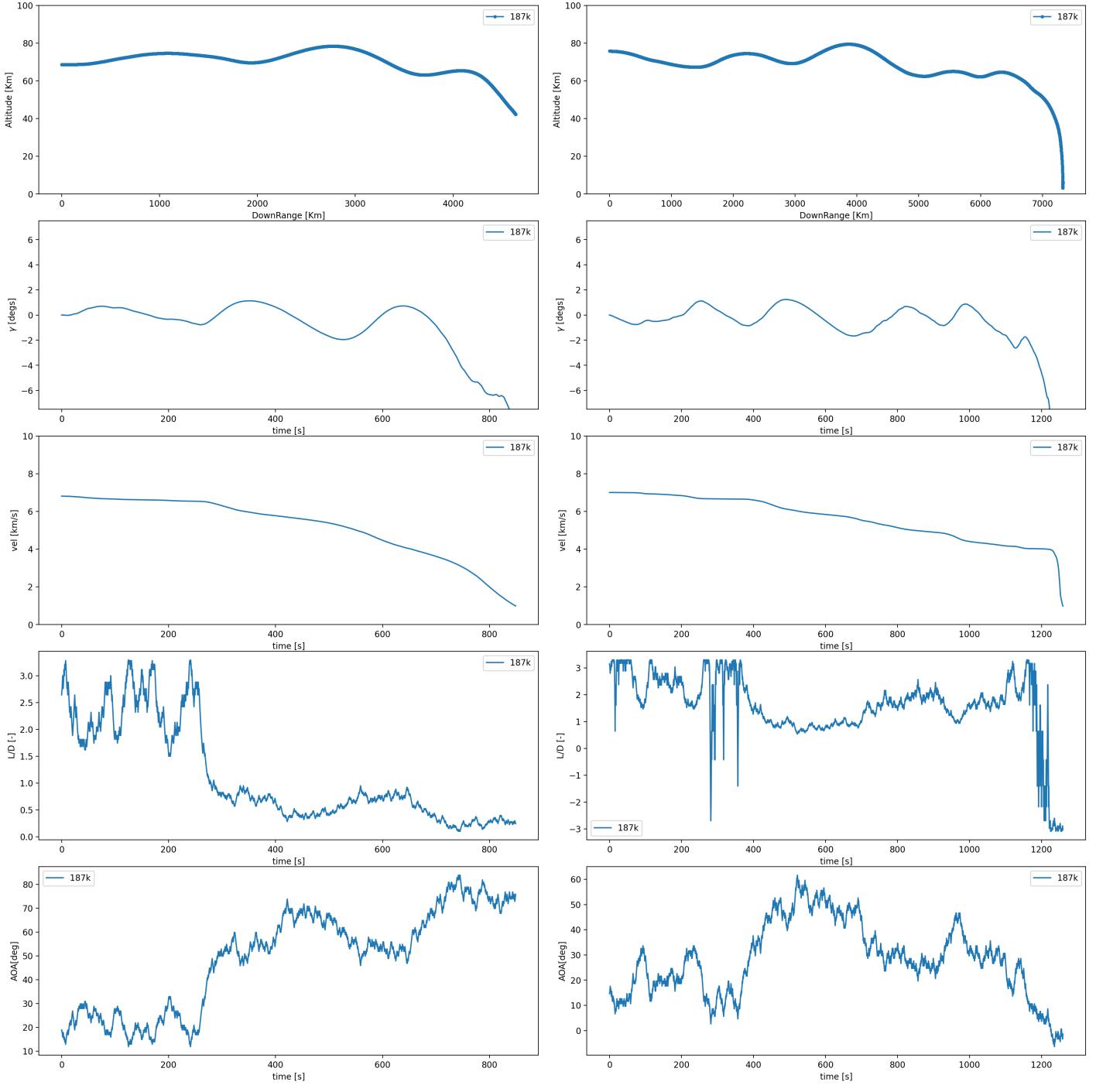


Figure 5: Output runs of DNN Q-learning. Blue signifies earlier portions of the training while red signifies later in training.

When we implement the policy after updating the results of two runs are shown in Figure 6. We see that the policy performance is acceptable for the earlier portions of the trajectory. However, as expected, once velocities are lower more lift is needed to maintain level flight. Hence, we see notable trends upward in subfigure 6a. We see in subfigure 6b just how sensitive initial chooses effect guidance for the policy. We note this by the subtle decrease in subfigure 6a as opposed to the increase in 6b leading to very different angle of attack progression but similar trajectory profiles in altitudes.



(a) Testing optimal policy run 1

(b) Testing optimal policy run 2

Figure 6: Online run of π^* after training

4 Discussion

Compared with other MDP problems, this is more complex due to the usage of the ODE and there is no analytical solution for the controls problem. However, we have showed a new approach for simulating and controlling a hypersonic glide vehicle using simplified physics and reinforcement learning. Some of the key insights that we have discovered while doing this project is that firstly, Choosing a suitable neural network architecture for the policy and value functions was not trivial, as different architectures had different trade-offs in terms of complexity, performance, and generalization. We experimented with various architectures and compared their results to find the best one for our problem. Additionally,

we have experienced that tuning the hyperparameters of the neural network and the reinforcement learning algorithm was also a difficult task, as they affected the convergence, stability, and robustness of the learning process. We used a trial-and-error approach and followed some best practices to find the optimal values for the hyperparameters.

We would report that our main issue was designing a proper reward function which was crucial for achieving our objective of maintaining a flat heading angle. We found that different reward functions had different impacts on the final output, and that some reward functions did not encourage the desired behavior. We also had to consider the trade-off between accuracy and computation time, as higher accuracy required more computation time and vice versa.

Some key limitations include the number of episodes that the model was trained on, even if it was then modeled to be generalizable for any initial condition. At the same time, we also had difficulty in balancing multiple goals, such as minimizing the heading angle but also not having it to terminate too early. Afterwards, due to the scope of the simplified physics it terminates early after the altitude reaches a high height. Qualitatively, this trajectory is good because there was not that much reduction in velocity while also having a straight heading angle. Unfortunately, this better controls can't be used or it won't be physically intuitive. In the end, we had to settle for something mediocre but was still able to perform better than the constant angle-of-attack.

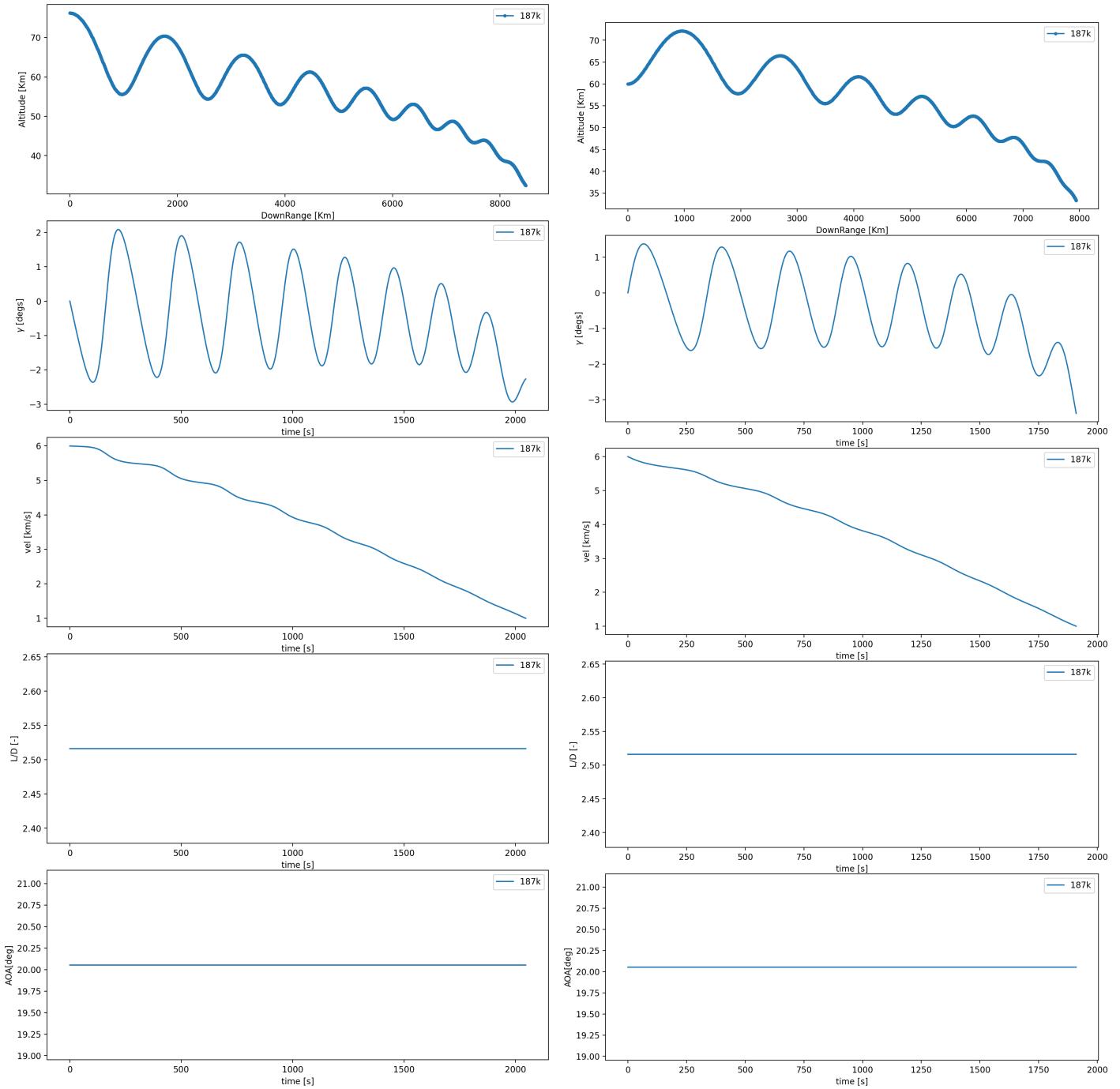
5 Conclusions

In conclusion, we have showed the feasibility of using a neural net for reinforcement learning of a hypersonic glide vehicle. This was achieved by running multiple episodes and was split into an exploration, training, and testing state. In the exploration stage, the action taken was completely random, while in the training stage, the randomness of the action choosing was slowly decreasing. With our goal of having a flat heading angle, the new trajectory of the HGV is much different than the constant angle of attack case. There were lots of issues that were faced during the simulation for each run because of the limitations on some of the state space values so it is still physically representative and accurate. Overall, in addition with training hyperparameters, this was addressed by changing the rewards, termination and neural net hyperparameters. And was able to get a proper control system using the first few timesteps of the gliding.

For future work, extended training time and a more complex neural network architecture may train the controller better to fit our objectives more accurately. A neural network generally gets more complex with various combinations of sequential blocks and activation functions as the model to be trained gets more complicated. In this case, with a 3 layer fully connected layer, this is definitely not enough as the number of nodes cannot fully represent a complex 4 state ordinary differential equation. Different work in the realm of machine learning has been done with the use of LSTMs or Residual Blocks which add complexity. Additionally, it can be noted that neural networks have a tendency to overfit the data so the use of dropout layers may benefit the model.

Appendix

Trajectories with constant angle-of-attack



(a) $\text{vel} = 6\text{km/s}$, $\text{alt}=76.2 \text{ km}$, $\alpha=20.05$

(b) $\text{vel} = 6\text{km/s}$, $\text{alt}=60.0 \text{ km}$, $\alpha=20.05$

Hypersonic Glide Vehicle Environment

```
1 import numpy as np
2 from na_python_version import *
3 neqn = 5
4 import os
5
6 class plane():
7     def __init__(self):
8         os.chdir(r'C:\Users\GreenFluids_VR\Documents\GitHub\CUDA__Project\code\python')
9
10    def render(self, eps=999, state=None):
11        fname = "../geometry/glider_187k.stl"
12        self.fresults = "../data/Run 9\trajNA187_PY_RK4_NN_" + str(eps) + ".txt"
13        rho, P_inf, L_D = 0.0, 0.0, 0.0
14        self.t = 0
15        self.dt = 0.5
16        self.t1 = np.copy(self.dt)
17
18        y = np.zeros(5)
19        yp = np.zeros(5)
20        ntri=187140
21        na = NA(ntri)
22        read_STL(fname,na)
23        self.cumreward = 0
24        if state is None:
25            y[0] = 0.0      # gamma (0 rads)
26            y[1] = np.random.randint(5000, 8000) #6000.0   # vel (6 km/s)
27            y[2] = np.random.randint(50000, 76200) #76200   # altitude (60 km)
28            y[3] = 0.0      # downrange (0 km)
29            y[4] = np.random.uniform(0.10, 0.35) #0.35   # alpha
30        else:
31            y = np.copy(state)
32
33        self.na = na
34        self.y = y
35        self.yp = yp
36        EOM(self.na, self.t, self.y, self.yp)
37        self.L_D = na.L / na.D
38        self.startwrite = False
39        self.write()
40
41        return self.y
42
43    def step(self, action):
44        del_alpha = np.pi/180 if action == 1 else -np.pi/180
45        state = np.copy(self.y)
46        self.y[4] += del_alpha
47        done = False
48        self.y = runge_kutta_4(self.na, self.t, self.y, self.dt, self.yp)
49        reward = -(180/np.pi * self.y[0]) ** 2
50        if self.y[1] < 1e3 or self.y[2] > 85e3 or self.y[4] < -np.pi/2 or self.y[4] > np.pi/2:
51            done = True
52            reward += -1e5
53        if self.y[3] < 1e3 or self.y[2] > 85e3:
54            reward += -1e6
55        self.t = self.t1
56        self.t1 += self.dt
57        self.L_D = self.na.L / self.na.D
58        self.cumreward += reward
59        self.write()
60        return state, self.y, reward, done
61
62    def write(self):
63        if self.startwrite == False:
64            with open(self.fresults, 'w') as outputFile:
65                outputFile.write(f"{self.t} {self.y[0]} {self.y[1]} {self.y[2]} {self.y[3]} {self.y[4]} {self.L_D} {self.cumreward}\n")
```

```
66     self.startwrite = True
67 else:
68     with open(self.fresults, 'a') as outputFile:
69         outputFile.write(f'{self.t} {self.y[0]} {self.y[1]} {self.y[2]} {self.y[3]} {self
.y[4]} {self.L_D} {self.cumreward}\n')
```

Neural Network with Deep Q-Learning

```
1 # %% Importing
2 import random
3 import torch
4 import torch.nn as nn
5 import torch.optim as optim
6 import torch.nn.functional as F
7 import numpy as np
8 import os
9
10
11 from sim_NeuralNet import plane
12
13 EPISODES = 100
14 TRAIN_START = 0.5
15 TRAIN_DUR = 0.5
16 EXPLORE_EPI_END = int(TRAIN_START*EPISODES)
17 TEST_EPI_START = int((TRAIN_START + TRAIN_DUR)*EPISODES )
18 EPS_START = 1.0
19 EPS_END = 0.05
20 EPS_DECAY = 1+np.log(EPS_END) / (TRAIN_DUR*EPISODES)
21 GAMMA = 0.9
22 LR = 0.001
23 MINIBATCH_SIZE = 25
24 ITERATIONS = 40
25 REP_MEM_SIZE = 10000
26
27 neqn = 5
28
29 use_cuda = torch.cuda.is_available()
30 FloatTensor = torch.cuda.FloatTensor if use_cuda else torch.FloatTensor
31 LongTensor = torch.cuda.LongTensor if use_cuda else torch.LongTensor
32 ByteTensor = torch.cuda.ByteTensor if use_cuda else torch.ByteTensor
33
34 class QNet(nn.Module):
35     def __init__(self, state_space_dim, action_space_dim):
36         nn.Module.__init__(self)
37         self.l1 = nn.Linear(state_space_dim, 24)
38         self.l2 = nn.Linear(24, 24)
39         self.l3 = nn.Linear(24, action_space_dim)
40     def forward(self, x):
41         x = F.relu(self.l1(x))
42         x = F.relu(self.l2(x))
43         x = self.l3(x)
44         return x
45 class ReplayMemory: ## (D)
46     def __init__(self, capacity):
47         self.capacity = capacity
48         self.memory = []
49     def push(self, s, a, r, ns):
50         self.memory.append((FloatTensor([s]),
51                            a,
52                            FloatTensor([ns]),
53                            FloatTensor([r])))
54         if len(self.memory) > self.capacity:
55             del self.memory[0]
56     def sample(self, MINIBATCH_SIZE):
57         return random.sample(self.memory, MINIBATCH_SIZE)
58     def __len__(self):
59         return len(self.memory)
60
61 class QNetAgent:
62     def __init__(self, stateDim, actionDim):
63         self.sDim = stateDim
64         self.aDim = actionDim
65         self.model = QNet(self.sDim, self.aDim)
66         if use_cuda:
```

```

67     self.model.cuda()
68     self.optimizer = optim.Adam(self.model.parameters(), LR)
69     self.lossCriterion = torch.nn.MSELoss()
70     self.memory = ReplayMemory(REP_MEM_SIZE)
71
72     self.steps_done = 0
73     self.episode_durations = []
74     self.avg_episode_duration = []
75     self.epsilon = EPS_START
76     self.epsilon_history = []
77     self.mode = ""
78     self.eps_done = 0
79
80 def select_action(self, state):
81     self.steps_done += 1
82     if random.random() > self.epsilon:
83         with torch.no_grad():
84             return self.model(FloatTensor(state)).data.max(1)[1].view(1, 1)
85     else:
86         return LongTensor([[random.randrange(2)]])
87
88
89 def run_episode(self, e, environment):
90     state = environment.render(self.eps_done)
91     self.eps_done += 1
92     done = False
93     steps = 0
94
95     if e < EXPLORE_EPI_END:
96         self.epsilon = EPS_START
97         self.mode = "Exploring"
98     elif EXPLORE_EPI_END <= e <= TEST_EPI_START:
99         self.epsilon = self.epsilon*EPS_DECAY
100        self.mode = "Training"
101    elif e > TEST_EPI_START:
102        self.epsilon = 0.025
103        self.mode = "Testing"
104    self.epsilon_history.append(self.epsilon)
105    while True:
106        action = self.select_action(FloatTensor([state]))
107        c_action = action.data.cpu().numpy()[0,0]
108        state, next_state, reward, done = environment.step(c_action)
109        self.memory.push(state, action, reward, next_state)
110        if EXPLORE_EPI_END <= e <= TEST_EPI_START:
111            self.learn()
112        state = next_state
113        steps += 1
114        if done:
115            print("{} Mode: {} | Episode {} Duration {} steps | epsilon {}"
116                  .format(e, steps, "\u033[92m" if steps >= 195 else "\u033[99m", self.epsilon, self.
117 mode))
118            self.episode_durations.append(steps)
119            break
120    return environment.cumreward
121
122 def finaltest(self, eps, state):
123     state = environment.render(eps=eps, state=state)
124     self.epsilon = 0.025
125     while True:
126         action = self.select_action(FloatTensor([state]))
127         c_action = action.data.cpu().numpy()[0,0]
128         state, next_state, reward, done = environment.step(c_action)
129         state = next_state
130         steps += 1
131         if done:
132             print(" | Duration {} steps | epsilon {}"
133                   .format(steps, self.epsilon))

```

```

133         break
134
135
136     def learn(self):
137         if len(self.memory) < MINIBATCH_SIZE:
138             return
139         for i in range(ITERATIONS):
140
141             experiences = self.memory.sample(MINIBATCH_SIZE)
142             batch_state, batch_action, batch_next_state, batch_reward = zip(*experiences)
143
144             batch_state = torch.cat(batch_state)
145             batch_action = torch.cat(batch_action)
146             batch_reward = torch.cat(batch_reward)
147             batch_next_state = torch.cat(batch_next_state)
148
149             current_q_values = self.model(batch_state).gather(1, batch_action)
150
151             max_next_q_values = self.model(batch_next_state).detach().max(1)[0]
152             expected_q_values = batch_reward + (GAMMA * max_next_q_values)
153
154             loss = self.lossCriterion(current_q_values, expected_q_values.unsqueeze(1))
155
156             self.optimizer.zero_grad()
157             loss.backward()
158             self.optimizer.step()
159
160 if __name__ == "__main__":
161     environment = plane()
162     state_size = 5
163     action_size = 2
164     rewards = np.zeros(EPIISODES)
165
166     agent = QNetAgent(state_size, action_size)
167     for e in range(EPIISODES):
168         reward = agent.run_episode(e, environment)
169         rewards[e] = reward
170     print('Complete')
171     test_epi_duration = agent.episode_durations[TEST_EPI_START-EPIISODES+1:]
172     print("Average Test Episode Duration", np.mean(test_epi_duration))
173
174     os.chdir(r'C:\Users\GreenFluids_VR\Documents\GitHub\CUDA__Project\code\data')
175     torch.save(agent.model.state_dict(), "weights1.h5")
176     with open('rewards.txt', 'w') as outputFile:
177         for e in range(EPIISODES):
178             outputFile.write(f"{rewards[e]}\n")
179     agent.finaltest(state=[0, 6e3, 60e3, 0, 0.35])
180     agent.finaltest(state=[0, 6e3, 75e3, 0, 0.35])

```