# Deep Q-Learning for Flappy Bird: Implementation and Report

Călin Ionuț, Lăzurcă Daniel

January 12, 2025

## Abstract

This report presents a Deep Q-Learning (DQN) approach to train an agent to play the Flappy Bird game. The implementation involves designing a convolutional neural network (CNN) for approximating Q-values, using a replay memory for experience replay, and applying techniques such as image preprocessing and epsilon-greedy exploration. This report provides a detailed description of the architecture, methodology, and implementation in Python.

## 1 Introduction

Deep Q-Learning combines reinforcement learning with deep neural networks to solve high-dimensional state-space problems. In this project, the Flappy Bird environment is utilized to test and demonstrate the capabilities of a DQN agent. The key components of the implementation include the agent's architecture, replay memory, and image preprocessing pipeline.

## 2 Code Components

This section describes the core components of the implementation.

### 2.1 DQNAgent Class

The `DQNAgent` class defines the agent that interacts with the environment, learns from experiences, and updates its policy. Key parameters such as `batch_size`, `gamma`, and `epsilon` are adjustable. A snippet of the class implementation is shown below:

```python
def init_weights(l):
    if type(l) == torch.nn.Linear or type(l) == torch.nn.
        Conv2d:
        torch.nn.init.xavier_normal_(l.weight)
        l.bias.data.fill_(0.01)

class DQNAgent:
    def __init__(self, input_shape, num_actions, batch_size
        =32, gamma=0.99, eps=1, eps_min=0.01, eps_decay=0.999,
         replay_buffer=5000):
        self.env = gymnasium.make("FlappyBird-v0",
            render_mode="rgb_array", use_lidar=False)
        self.device = torch.device("cuda" if torch.cuda.
            is_available() else "cpu")
        self.policy_net = DQN(input_shape, num_actions).to(
            self.device)
        self.policy_net.apply(init_weights)
        self.target_net = DQN(input_shape, num_actions).to(
            self.device)
        self.target_net.load_state_dict(self.policy_net.
            state_dict())
        self.target_net.eval()
        self.optimizer = torch.optim.Adam(self.policy_net.
            parameters())
        self.memory = ReplayMemory(replay_buffer)
        self.batch_size = batch_size
        self.gamma = gamma
        self.eps = eps
        self.eps_min = eps_min
        self.eps_decay = eps_decay
```

Listing 1: DQNAgent Implementation

## 2.2 Replay Memory

Replay memory is implemented using a double-ended queue (deque) to store experiences and allow random sampling for training. This improves stability and efficiency. The implementation is as follows:

```python
from collections import deque
from random import sample

class ReplayMemory:
    def __init__(self, capacity):
        self.memory = deque(maxlen=capacity)

    def push(self, state, action, reward, next_state, done):
```

```
 9          self.memory.append((state, action, reward, next_state
               , done))
10
11      def __len__(self):
12          return len(self.memory)
13
14      def sample(self, batch_size):
15          return sample(self.memory, batch_size)
```

<div align="center">Listing 2: Replay Memory</div>

## 2.3 Deep Q-Network (DQN)

The DQN is a convolutional neural network designed to process image inputs and predict Q-values for actions. It includes convolutional layers, fully connected layers, and dropout for regularization. The implementation is as follows:

```
 1  import torch.nn as nn
 2  import torch.nn.functional as F
 3  import torch
 4  import numpy as np
 5
 6  class DQN(nn.Module):
 7      def __init__(self, input_shape, num_actions):
 8          super(DQN, self).__init__()
 9          self.conv1 = nn.Conv2d(in_channels=input_shape[0],
               out_channels=4, kernel_size=8, stride=4)
10          self.conv2 = nn.Conv2d(in_channels=4, out_channels=8,
                kernel_size=4, stride=2)
11          self.fc1 = nn.Linear(self._compute_conv_out(
               input_shape), 64)
12          self.fc2 = nn.Linear(64, num_actions)
13          self.dropout = nn.Dropout(p=0.2)
14
15      def _compute_conv_out(self, shape):
16          x = torch.zeros(1, *shape)
17          x = self.conv1(x)
18          x = self.conv2(x)
19          return int(np.prod(x.size()))
20
21      def forward(self, x):
22          x = F.relu(self.conv1(x))
23          x = F.relu(self.conv2(x))
24          x = x.view(x.size(0), -1)
25          x = F.relu(self.fc1(x))
26          x = self.dropout(x)
27          x = self.fc2(x)
```

```
28            return x
```

Listing 3: DQN Architecture

## 2.4   Image Preprocessing

To improve the agent's performance, we use the last 4 frames as the state, preprocessed as follows:

- Scale down to 72x72, (because 72 is a divisor of the frame's width and height)

- Separate on each color channel

- Apply binarization on red and blue to get the features of the bird, and on green to get features of the pipes

- Combine the binarized masks for red and blue to get the bird mask, which we expand using a convolution

- Apply the bird mask over the pipe, such that the pipes are white and the bird is gray. The rest is black.

The process_image function is implemented as follows:

```
1  def process_image(image):
2      image = image[20:380, :, :]
3      image = cv2.resize(image, (72, 72))
4
5      bird_1 = image[:, :, 0]
6      bird_1[bird_1 > 227] = 255
7      bird_1[bird_1 <= 227] = 0
8
9      bird_2 = image[:, :, 2]
10     bird_2[bird_2 > 138] = 255
11     bird_2[bird_2 <= 138] = 0
12
13     bird = np.zeros((72, 72))
14     bird[(bird_1 == 255) | (bird_2 == 255)] = 255
15     bird_mask = apply_convolution(bird)
16
17     pipes = image[:, :, 1]
18     pipes[pipes > 136] = 255
19     pipes[pipes <=136] = 0
20
21     final_state = pipes
```

```
22    final_state = np.where(bird_mask == 255, 127, final_state
          )
23
24    final_state = np.where(final_state == 255, 1, final_state
          )
25    final_state = np.where(final_state == 127, 0.5,
          final_state)
26
27    return final_state
```

Listing 4: Image Preprocessing

## 2.5 Testing Player

This agent would take a model that was pretrained by our implementations and uses it to play Flappy Bird. It tries to play 5 games and outputs the final rewards received, whilst showcasing the gameplay itself. It's implementation is shown below:

```
1  import torch
2  import gymnasium
3  import flappy_bird_gymnasium
4  from utils.dqn import DQN
5  from utils.img_processor import process_image
6  import numpy as np
7
8  class DQNPlayer:
9      def __init__(self, model_path, input_shape, num_actions,
           seed=42):
10         self.env_rgb = gymnasium.make("FlappyBird-v0",
               render_mode="rgb_array",use_lidar = False)
11         self.env_human = gymnasium.make("FlappyBird-v0",
               render_mode="human",use_lidar = False)
12         self.device = torch.device("cuda" if torch.cuda.
               is_available() else "cpu")
13
14         self.env_rgb.reset(seed=seed)
15         self.env_human.reset(seed=seed)
16
17         self.policy_net = DQN(input_shape, num_actions).to(
               self.device)
18         self.policy_net.load_state_dict(torch.load(model_path
               , map_location=self.device))
19         self.policy_net.eval()
20
21      def play(self, episodes=5):
22         for episode in range(episodes):
23             seed = np.random.randint(0, 10000)
```

```python
24              state_rgb, _ = self.env_rgb.reset(seed=seed)
25              self.env_human.reset(seed=seed)
26
27              frame = self.env_rgb.render()
28              if frame is None:
29                  raise ValueError("Rendered frame is None.
                        Check environment's render_mode.")
30              state = process_image(frame)
31              state = np.stack([state] * 4, axis=0)
32
33              done = False
34              total_reward = 0
35
36              while not done:
37                  self.env_human.render()
38
39                  state_tensor = torch.FloatTensor(state).
                        unsqueeze(0).to(self.device)
40                  with torch.no_grad():
41                      action = self.policy_net(state_tensor).
                            argmax(dim=1).item()
42
43                  _, reward, done, _, _ = self.env_rgb.step(
                        action)
44                  self.env_human.step(action)
45                  total_reward += reward
46
47
48                  frame = self.env_rgb.render()
49                  if frame is None:
50                      raise ValueError("Rendered frame is None
                            during gameplay.")
51                  next_state = process_image(frame)
52                  state = np.append(state[1:], [next_state],
                        axis=0)
53              print(f"Episode {episode + 1}: Total Reward = {
                    total_reward:.1f}")
54
55          self.env_rgb.close()
56          self.env_human.close()
57
58  if __name__ == "__main__":
59      model_path = "models/best_model_62.6.pth"
60      player = DQNPlayer(model_path, input_shape=(4, 72, 72),
            num_actions=2, seed=42)
61      player.play(episodes=5)
```

Listing 5: Testing player

# 3   Training

The training loop runs for a specified number of epochs. During each episode, the agent interacts with the environment, collects experiences, and updates the policy network using the Bellman equation. Target networks are periodically updated to stabilize learning.

## 3.1   Hyperparameters and Results

The performance of the DQN agent depends significantly on the choice of hyperparameters. The following table summarizes the hyperparameters used in this implementation: We tested two sets of hyperparameters , modifying the replay buffer and the number of epochs after witch the model trains on mini-batches.

| Hyperparameter | Value |
|---|---|
| Batch Size | 32 |
| Replay Memory Capacity | 50,000 |
| Discount Factor ($\gamma$) | 0.95 |
| Initial Epsilon ($\epsilon$) | 1.0 |
| Minimum Epsilon ($\epsilon_{min}$) | 0.01 |
| Epsilon Decay Rate | 0.999 |
| Learning Rate | 0.001 |
| Target Network Update Frequency | 5 epochs |
| Training Episodes | 10,000 |
| Input Image Dimensions | (72, 72) |

Table 1: Second set of hyperparameters used in training the DQN agent.

As we can see from the results, using a larger replay buffer tends to lead to a more stable learning curve, yet the Neural Network learns at a slower pace. Meanwhile, using a smaller value for the replay buffer, the network starts learning somewhat faster, but reaches a pleatou much quicker and doesn't show that much improvement.
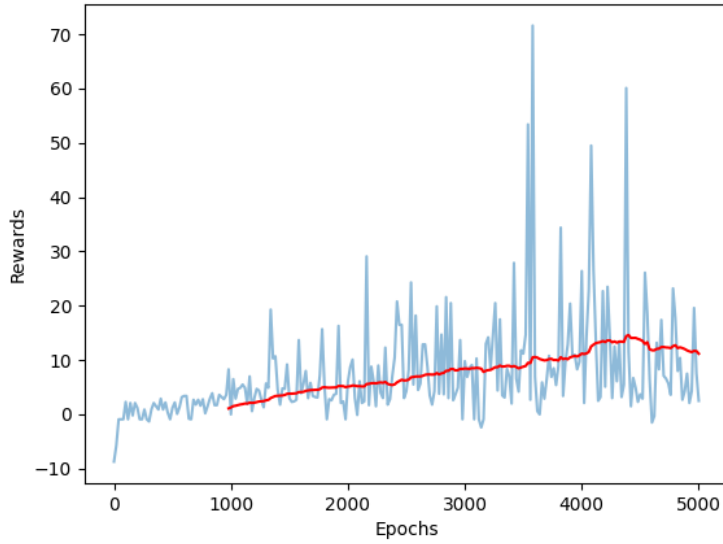
Figure 1: Learning curve: Cumulative return achieved per episode during training using the second set of hyperparameters.
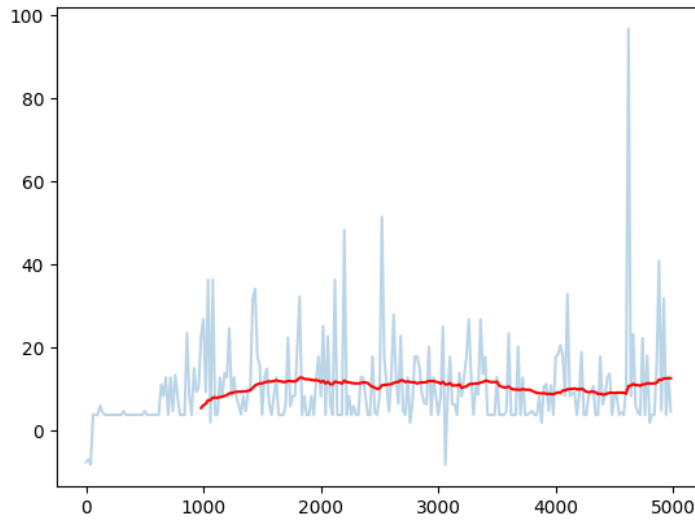


Figure 2: Learning curve: Cumulative return achieved per episode during training using the first set of hyperparameters.

| Hyperparameter | Value |
| --- | --- |
| Batch Size | 32 |
| Replay Memory Capacity | 25,000 |
| Discount Factor ($\gamma$) | 0.95 |
| Initial Epsilon ($\epsilon$) | 1.0 |
| Minimum Epsilon ($\epsilon_{min}$) | 0.01 |
| Epsilon Decay Rate | 0.999 |
| Learning Rate | 0.001 |
| Target Network Update Frequency | 10 epochs |
| Training Episodes | 10,000 |
| Input Image Dimensions | (72, 72) |

Table 2: First set of hyperparameters used in training the DQN agent.

# 4 Other solutions

At first, we tried implementing a Actor Critic model for the Neural Network, but we saw some problems that we found hard to fix with that approach. Mainly, we saw the model tending to prioritize getting quick, easy and guaranteed rewards in the short term, rather than taking risks in order to achieve some long term, higher rewards. Because of that , the scores received in training reached pleatous really quickly and were not very high, meaning the model was pretty weak at completing it's task. This was the reason why we switched to DQN, where a single neural network is doing the job of both the Actor and the Critic, and it shows massive improvements.

# 5 Conclusion

This project demonstrates the application of Deep Q-Learning to a challenging environment like Flappy Bird. The implementation highlights key aspects such as efficient memory management, CNN-based Q-function approximation, and effective image preprocessing. Future work could involve experimenting with hyperparameters and exploring alternative architectures.