

Flappy Bird AI Reinforcement Learning Project

10

2. Pre-trained Model Integration
 3. Reinforcement Learning Implementation (DQN)
 4. Model Training
 5. Testing and Evaluation
- throughout each step, I have provided explanation

1. Environment Setup

I began by importing all the new data into the database. When the screen is displayed, which makes

```
import numpy as np
import cv2
```

[illegible]

I leveraged MobileNetV2 as a feature extractor

```
from tensorflow.keras.layers import ReLU
```

```
def build_q_model(action_size, learning_rate=0.0005):
    base_model = MobileNetV2(v2={'include_top': 'imagenet', 'include_top': False, 'input_shape': (224, 224, 3)})
    for layer in base_model.layers:
        layer.trainable = False

    x = GlobalAveragePooling2D()(base_model.output)
    # Add a hidden layer for more representational power
    x = Dense(128, activation='relu')(x)
    q_values = Dense(action_size, activation='linear')(x)

    model = Model(inputs=base_model.input, outputs=q_values)
    model.compile(optimizer=Adam(learning_rate=learning_rate), loss='mse')
    return model
```

For my DQN agent, I implemented experience replay and a target network. I set ϵ decay=0.999 so that the agent would remain exploratory for

amount of exploration even late in training.

```

        learning_rate=0.00005,
        gamma=0.99,
        epsilon=1.0,
        epsilon_min=0.01,
        epsilon_decay=0.999,
        memory_size=100000):
    self.action_size = action_size
    self.memory = deque(maxlen=memory_size)
    self.gamma = gamma
    self.epsilon = epsilon
    self.epsilon_min = epsilon_min
    self.epsilon_decay = epsilon_decay
    self.learning_rate = learning_rate

    self.model = build_q_model(self.action_size, learning_rate=self.learning_rate)
    self.target_model = build_q_model(self.action_size, learning_rate=self.learning_rate)
    self.update_target_model()

def update_target_model(self):
    self.target_model.set_weights(self.model.get_weights())

def remember(self, state, action, reward, next_state, done):
    self.memory.append((state, action, reward, next_state, done))

def act(self, state):
    # With slower decay, agent remains exploratory longer, possibly finding better strategies
    if np.random.rand() <= self.epsilon:
        return np.random.randint(self.action_size)
    q_values = self.model.predict(state, verbose=0)
    return np.argmax(q_values[0])

def replay(self, batch_size=64):
    if len(self.memory) < batch_size:
        return
    minibatch = random.sample(self.memory, batch_size)
    states = np.array([mb[0][0] for mb in minibatch])
    next_states = np.array([mb[3][0] for mb in minibatch])

    states_q = self.model.predict(states, verbose=0)
    next_states_q = self.model.predict(next_states, verbose=0)
    next_states_tq = self.target_model.predict(next_states, verbose=0)

    for l, (state, action, reward, next_state, done) in enumerate(minibatch):
        q_values = states_q[l]
        if done:
            q_values[action] = reward
        else:
            next_action = np.argmax(next_states_q[l])
            q_values[action] = reward + self.gamma * next_states_tq[l][next_action]

    self.model.fit(states, states_q, epochs=1, verbose=0)

    if self.epsilon > self.epsilon_min:
        self.epsilon *= self.epsilon_decay

```

```
batch_size = 64
```

```

    range_update_freq = 10
    frame_skip = 1 # act every frame for finer control

    reset = False # if desired, can be set to True

```

```

gamma=0.99,
epsilon=1.0,
epsilon_min=0.01,
epsilon_decay=0.999,
memory_size=100000)

scores = []
for episode in tqdm(range(n_episodes)):
    env.reset()
    done = False
    total_reward = 0
    state = preprocess(env.getScreenRGB())

    while not done:
        # With no frame_skip, the agent has full control every frame.
        action_idx = agent.act(state)
        action = action_set[action_idx]
        original_reward = env.act(action)
        done = env.game_over()
        reward = get_reward(env, original_reward)

        next_state = preprocess(env.getScreenRGB())
        agent.remember(state, action_idx, reward, next_state, done)

        state = next_state
        total_reward += reward

    if len(agent.memory) > batch_size:
        agent.replay(batch_size)

scores.append(total_reward)
if episode % target_update_freq == 0:
    agent.update_target_model()

if (episode+1) % 50 == 0:
    avg_last_50 = np.mean(scores[-50:])
    print(f"Episode: {episode+1}, Avg Score (last 50 eps): {avg_last_50:.2f}, Epsilon: {agent.epsilon:.2f}")

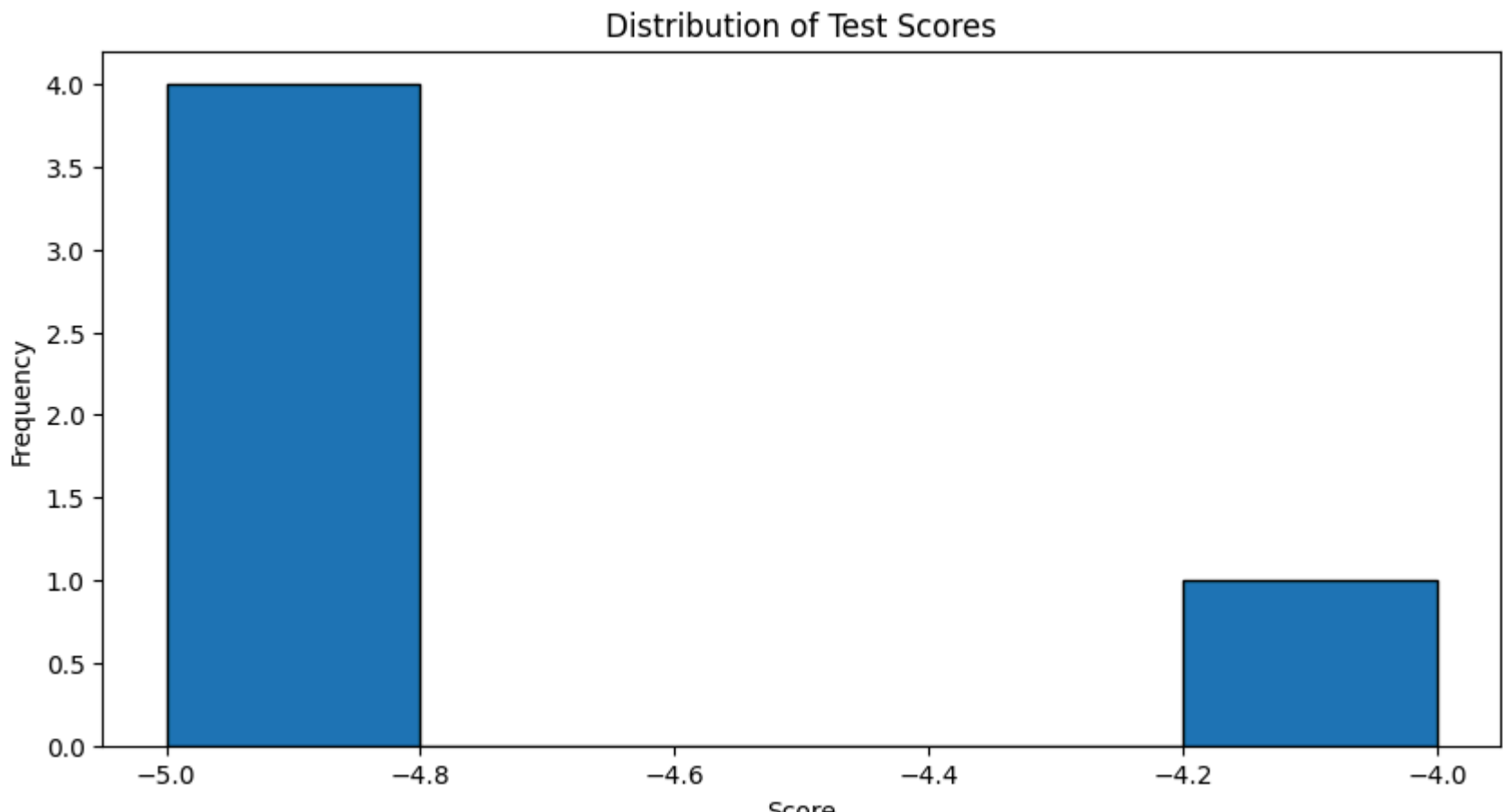
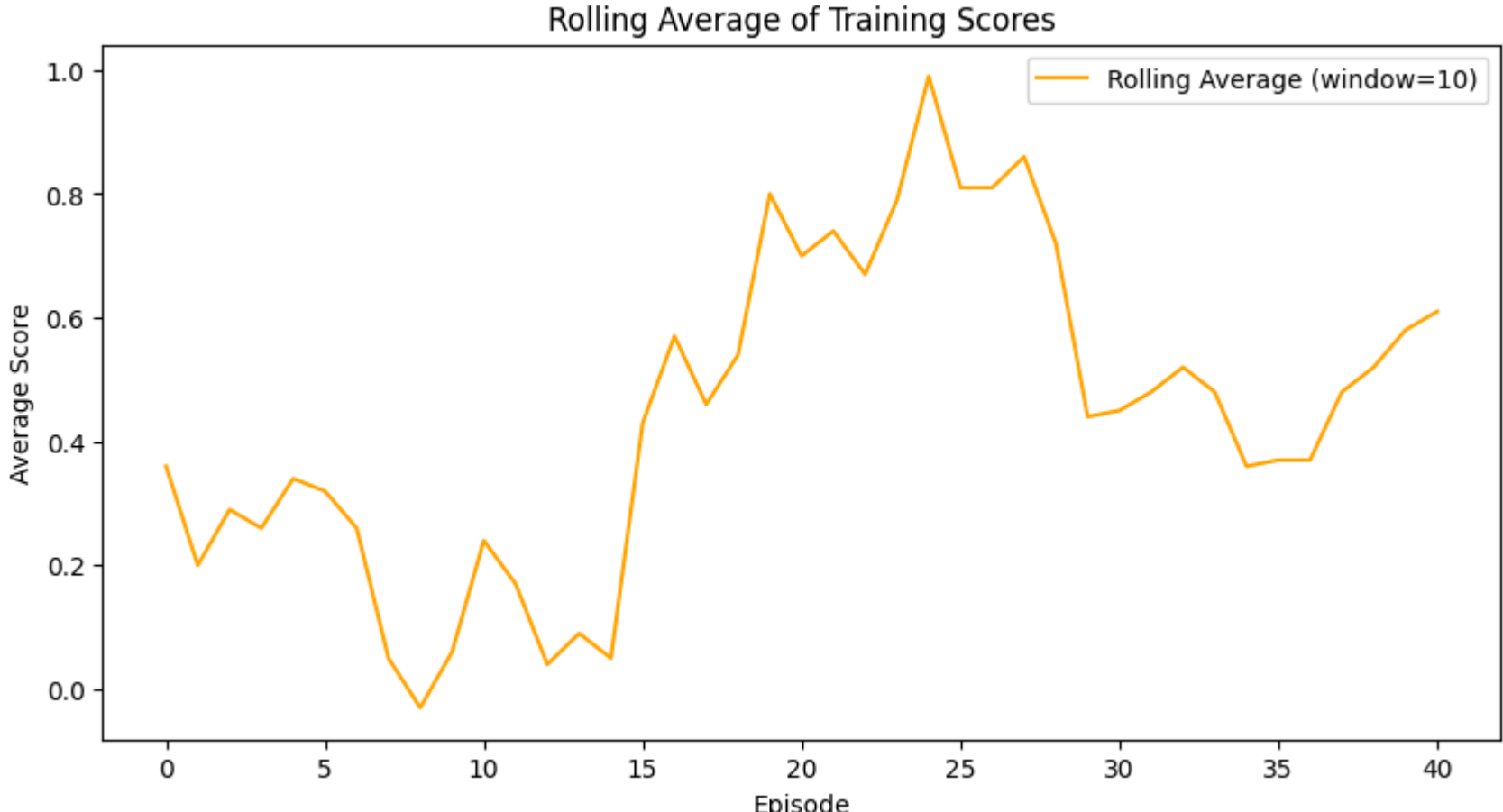
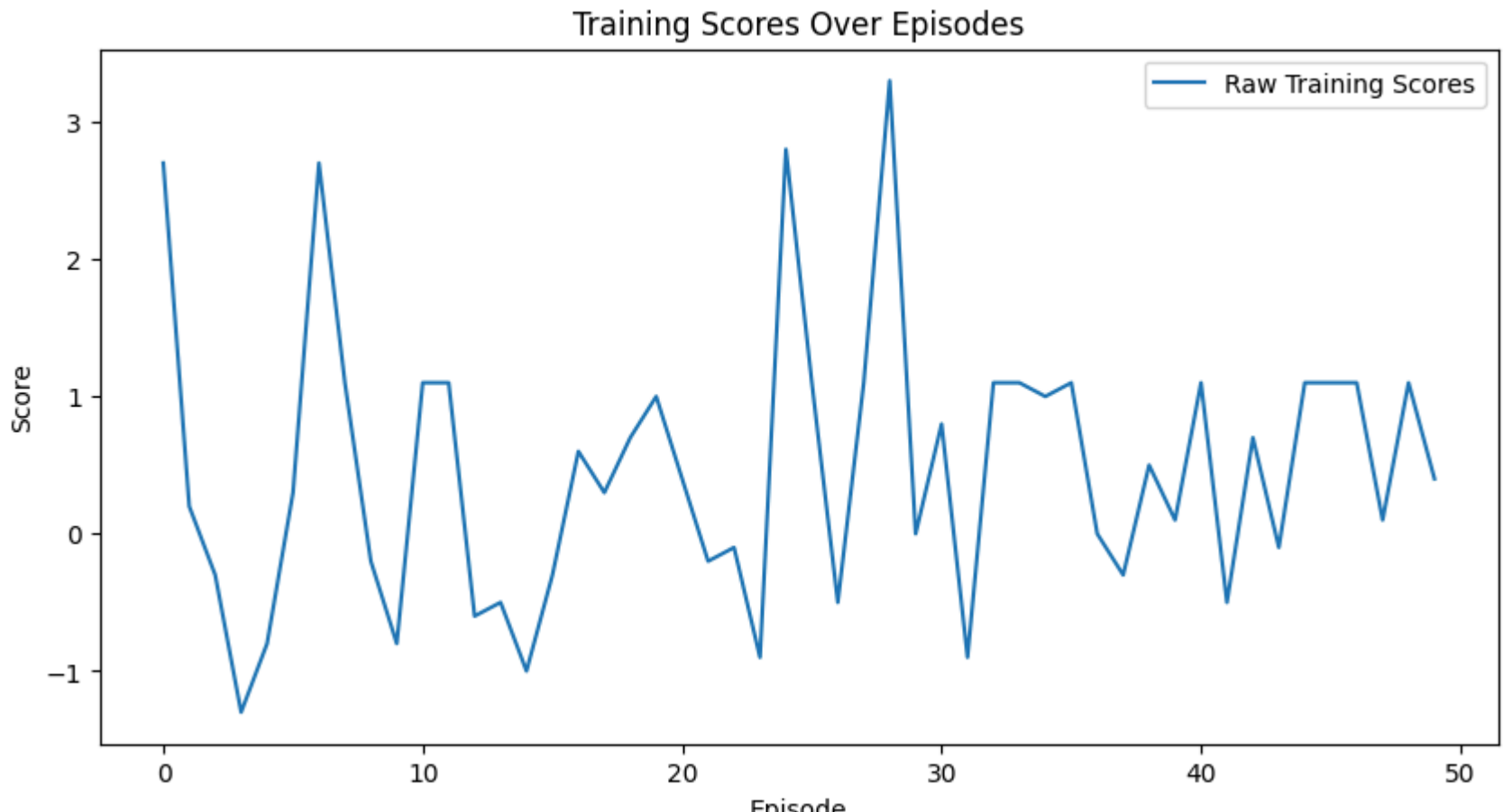
2024-12-10 11:25:15.863014: I metal_plugin/src/device/metal_device.cc:1154] Metal device set to: Apple M3 Max
2024-12-10 11:25:15.863036: I metal_plugin/src/device/metal_device.cc:296] systemMemory: 48.00 GB
2024-12-10 11:25:15.863038: I metal_plugin/src/device/metal_device.cc:313] maxCacheSize: 18.00 GB
2024-12-10 11:25:15.863050: I tensorflow/core/common_runtime/pluggable_device/pluggable_device_factory.cc:305] Could not identify NUMA node of platform GPU ID 0, defaulting to 0. Your kernel
1 may not have been built with NUMA support.
2024-12-10 11:25:15.863059: I tensorflow/core/common_runtime/pluggable_device/pluggable_device_factory.cc:271] Created TensorFlow device (/job:localhost/replica:0/task:0/device:GPU:0 with 0
MB memory) -> physical PluggableDevice (device: 0, name: METAL, pci bus id: <undefined>)
0 | 50 | 00:00<7, ?it/s|2024-12-10 11:25:17.753816: I tensorflow/core/grappler/optimizers/custom_graph_optimizer_registry.cc:117] Plugin optimizer for device_type GPU is enabl
ed.
100% ██████████ 50/50 [20:17<00:00, 24.35s/it]
Episode: 50, Avg Score (last 50 eps): 0.47, Epsilon: 0.07

```

training scores to smooth out noise (

All these visualizations rely only on scores and test_scores collected during training and testing, without requiring any modifications to previous code cells.

2019年10月10日



Test Episode 1/5: Score = -5.0
Test Episode 2/5: Score = -5.0
Test Episode 3/5: Score = -5.0
Test Episode 4/5: Score = -4.0
Test Episode 5/5: Score = -5.0
Average Test Score: -4.8