

# Crossy Road Reinforcement Learning

Neel Shah, Anish Bagri, Nathan Sankar, Akul Gokaram, Aditya Sharda

May 7, 2024

## 1 Introduction

[Crossy Road](#) is an arcade video game built around the age-old joke “Why did the chicken cross the road?” In the game, the chicken (controlled by the player) has to cross the road without getting hit by vehicles.

We want to develop a reinforcement learning (RL) game agent capable of playing Crossy Road. The game’s endless and random nature makes it a great candidate for RL.

The agent will learn to maximize its score by getting the chicken to cross the road and avoid obstacles in its path, with the ultimate goal of crossing the road as many times as possible without collisions. Once the agent is capable of successfully getting the chicken to cross the road and reach the goal position, another goal could be to minimize the time it takes for the chicken to cross the road.

## 2 Configuration

There are several things we need to configure before we can begin.

### 2.1 Importing Libraries

We’ll start by importing the necessary [Python](#) libraries.

```
[507]: %pip install "gymnasium[atari, accept-rom-license, other]"
```

```
/Users/neelshah/Documents/Projects/Crossy-Road-Reinforcement-  
Learning/.venv/bin/python: No module named pip
```

Note: you may need to restart the kernel to use updated packages.

```
[508]: import copy  
import math  
import random  
from collections import deque, namedtuple  
from datetime import datetime  
from pathlib import Path  
  
import cv2  
import gymnasium as gym  
import matplotlib.pyplot as plt
```

```

import numpy as np
import PIL.Image
import torch
import torch.nn as nn
import torch.optim as optim
from gymnasium.wrappers import TransformObservation
from IPython import display
from matplotlib.ticker import MaxNLocator
from torcheval.metrics import Mean
from torchvision.transforms import v2
from tqdm import trange

# code taken from: https://github.com/googlecolab/colabtools/blob/main/google/
# ↪ colab/patches/__init__.py
def cv2_imshow(a, convert=True):
    """A replacement for cv2.imshow() for use in Jupyter notebooks.

    Args:
        a: np.ndarray. shape (N, M) or (N, M, 1) is an NxM grayscale image. For
        example, a shape of (N, M, 3) is an NxM BGR color image, and a shape of
        (N, M, 4) is an NxM BGRA color image.
        convert: boolean.
    """
    a = a.clip(0, 255).astype("uint8")
    # cv2 stores colors as BGR; convert to RGB
    if convert:
        if a.ndim == 3:
            if a.shape[2] == 4:
                a = cv2.cvtColor(a, cv2.COLOR_BGRA2RGBA)
            else:
                a = cv2.cvtColor(a, cv2.COLOR_BGR2RGB)
    display.display(PIL.Image.fromarray(a))

```

## 2.2 Controlling Randomness

Next, we'll make our results more reproducible (deterministic) by controlling sources of randomness, following the [suggestions outlined in the PyTorch docs](#).

```

[509]: random.seed(0)
       np.random.seed(0)
       torch.manual_seed(0)

```

```

[509]: <torch._C.Generator at 0x14919c990>

```

## 2.3 Setting Device

We'll also select a [device](#) to store our tensors on.

```
[510]: device = (  
    "cuda"  
    if torch.cuda.is_available()  
    else "mps"  
    if torch.backends.mps.is_available()  
    else "cpu"  
)  
print(f"Using {device} device")
```

Using mps device

## 3 Crossy Road Environment

Our first step is to implement a Crossy Road environment, which will encapsulate our representation of the reinforcement learning problem that the game poses.

For this, we will utilize the [Gymnasium](#) library (a fork of the [OpenAI Gym](#) library), which provides a standard API for RL and various reference environments.

Specifically, we will use the [Freeway](#) environment, which models an [Atari](#) game that closely resembles Crossy Road. This gives us a Pythonic interface to work with, which we can later use to develop RL models and create an agent that can play Crossy Road successfully.

### 3.1 Initializing the Environment

We will start off by initializing the environment using Gymnasium.

We pass in the following arguments (documented [here](#)) to specify the environment:

#### Environment Flavor:

The environment `id`, `mode`, and `difficulty` combine to specify the specific flavor of the environment:

- `id="ALE/Freeway-v5"`: simulates the Atari game Freeway via the [Arcade Learning Environment \(ALE\)](#) through the [Stella](#) emulator
- `mode=0`: selects [Game 1 \(Lake Shore Drive, Chicago, 3 A.M.\)](#) as the map to use
- `difficulty=0`: selects the default difficulty setting

#### Stochasticity:

As stated in the documentation:

As the Atari games are entirely deterministic, agents can achieve state-of-the-art performance by simply memorizing an optimal sequence of actions while completely ignoring observations from the environment.

To combat this, we use `frameskip` and `repeat_action_probability`:

- `frameskip=4`: enables frame skipping (sets the number of frames to skip on each skip to 4)
- `repeat_action_probability=0.25`: enables sticky actions (sets the probability of repeating the previous action instead of executing the current action to 25%)

## Simulation:

The parameters `full_action_space` and `render_mode` are used to specify how the environment is simulated:

- `full_action_space=False`: limits the action space to the 3 legal actions we will actually use instead of all 18 possible actions that can be performed on an Atari 2600 console
- `render_mode="rgb_array"`: specifies that the game should be rendered as an RGB frame

```
[511]: env = gym.make(  
        id="ALE/Freeway-v5",  
        mode=0,  
        difficulty=0,  
        obs_type="rgb",  
        frameskip=4,  
        repeat_action_probability=0.25,  
        full_action_space=False,  
        render_mode="rgb_array",  
    )
```

We will also modify the metadata to set `render_fps` to 30, meaning the game will run at 30 frames per second.

```
[512]: env.metadata["render_fps"] = 30
```

Now, we are ready to learn a little more about how our environment is implemented.

## 3.2 Observations

Let's start with the observation space.

```
[513]: env.observation_space
```

```
[513]: Box(0, 255, (210, 160, 3), uint8)
```

This observation space represents the RGB image that is displayed to a human player.

## 3.3 Actions

Next, let's move on to the action space.

```
[514]: env.action_space
```

```
[514]: Discrete(3)
```

This action space represents the actions that the chicken can take in each step:

```
[515]: action_meaning = {  
        0: "NOOP",  
        1: "UP",  
        2: "DOWN",  
    }
```

```
print(action_meaning)
```

```
{0: 'NOOP', 1: 'UP', 2: 'DOWN'}
```

### 3.4 Rewards

Finally, let's move on to the reward range.

```
[516]: env.reward_range
```

```
[516]: (-inf, inf)
```

We can see that the reward range is  $(-\infty, \infty)$ . However, this is not very informative.

As the documentation tells us:

You receive a point for every chicken that makes it to the top of the screen after crossing all the lanes of traffic.

Our goal will be to maximize the total score in the time limit automatically set by the environment.

## 4 Agent-Environment Interaction

Our next step is to create a mechanism by which an agent can interact with the environment.

First, let's create a `log_step()` function that logs information about a certain time step.

```
[517]: def log_step(step, action, observation, reward, terminated, truncated, info, score):
        print(f"\n***** Step {step} *****\n")
        if action is not None:
            print(f"Action: {action} ({action_meaning[action]})")
        if observation is not None:
            print("Observation:", observation)
        if reward is not None:
            print("Reward:", reward)
        if terminated is not None:
            print("Terminated:", terminated)
        if truncated is not None:
            print("Truncated:", truncated)
        if info is not None:
            print("Info:", info)
        if score is not None:
            print("Score:", score)
```

Now, let's create a `simulate()` function that takes in an environment, agent, and number of episodes.

It simulates running `num_episodes` episodes in the environment `env`, where the player's actions are defined by the behavior of `agent`.

```

[518]: def simulate(
    env: gym.Env, transform: any, agent: any, num_episodes: int, step_show_freq:
    ↪ int
):
    for episode in range(num_episodes):
        print(f"##### Episode {episode} #####")
        step = 0
        action = None
        observation, info = env.reset()
        observation = (
            transform(observation).to(device).clone().detach().unsqueeze(dim=0)
        )
        reward, terminated, truncated = 0.0, False, False
        score = 0
        if step % step_show_freq == 0:
            log_step(step, action, None, reward, terminated, truncated, info, ↪
            ↪ score)
            view = env.render()
            cv2_imshow(view, convert=False)
            while not (terminated or truncated):
                step += 1
                action = agent.sample_action(observation)
                observation, reward, terminated, truncated, info = env.step(action.
                ↪ item())
                score += reward
                observation = transform(observation).clone().detach().
                ↪ unsqueeze(dim=0)
                if step % step_show_freq == 0:
                    log_step(
                        step,
                        action.item(),
                        None,
                        reward,
                        terminated,
                        truncated,
                        info,
                        score,
                    )
                view = env.render()
                cv2_imshow(view, convert=False)

```

## 5 Random Agent

To test our environment, let's create a RandomAgent.

## 5.1 Sampling an Action at Random

The agent will simply move the chicken randomly by sampling actions at random from the action space of the environment.

```
[519]: class RandomAgent:
        def __init__(self, env):
            self.env = env

        def sample_action(self, _observation):
            action = self.env.action_space.sample()
            return torch.tensor(data=[[action]], dtype=torch.long, device=device)
```

## 5.2 Simulation

Now, let's run a simulation.

To limit the amount of output we need to scroll through, we will only show every 100 steps.

```
[520]: transform = v2.Compose(
        [
            v2.ToImage(),
            v2.ToDtype(dtype=torch.float32, scale=True),
        ]
    )
    random_agent = RandomAgent(env)
    simulate(env, transform, random_agent, 1, 100)
```

```
##### Episode 0 #####
```

```
***** Step 0 *****
```

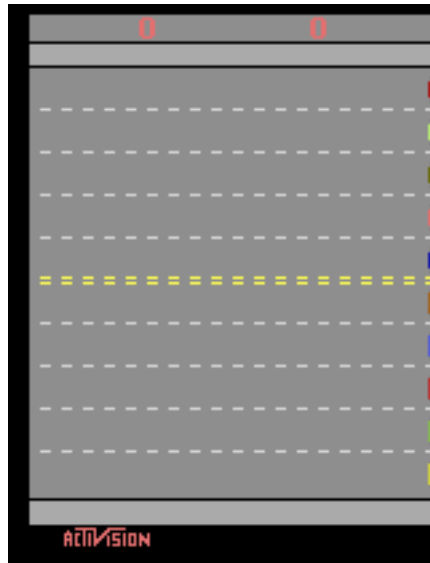
```
Reward: 0.0
```

```
Terminated: False
```

```
Truncated: False
```

```
Info: {'lives': 0, 'episode_frame_number': 0, 'frame_number': 0}
```

```
Score: 0
```



\*\*\*\*\* Step 100 \*\*\*\*\*

Action: 2 (DOWN)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 400, 'frame\_number': 400}

Score: 0.0





\*\*\*\*\* Step 200 \*\*\*\*\*

Action: 1 (UP)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 800, 'frame\_number': 800}

Score: 0.0



\*\*\*\*\* Step 300 \*\*\*\*\*

Action: 0 (NOOP)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 1200, 'frame\_number': 1200}

Score: 0.0



\*\*\*\*\* Step 400 \*\*\*\*\*

Action: 2 (DOWN)

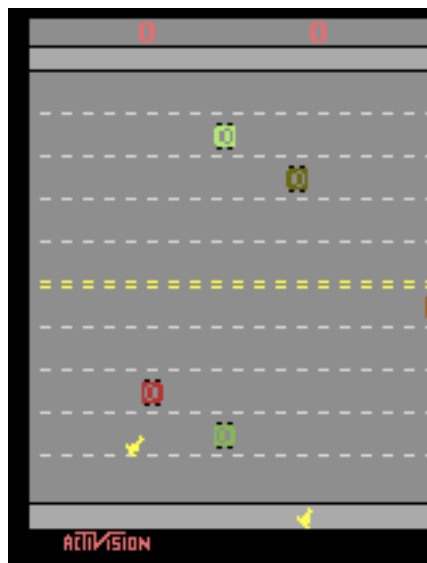
Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 1600, 'frame\_number': 1600}

Score: 0.0



\*\*\*\*\* Step 500 \*\*\*\*\*

Action: 1 (UP)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 2000, 'frame\_number': 2000}

Score: 0.0



\*\*\*\*\* Step 600 \*\*\*\*\*

Action: 2 (DOWN)

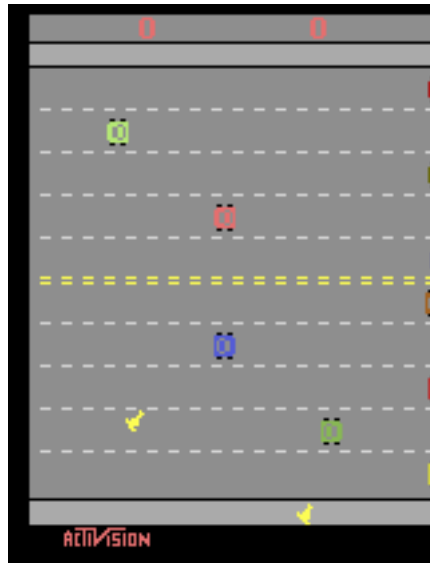
Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 2400, 'frame\_number': 2400}

Score: 0.0



\*\*\*\*\* Step 700 \*\*\*\*\*

Action: 0 (NOOP)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 2800, 'frame\_number': 2800}

Score: 0.0



\*\*\*\*\* Step 800 \*\*\*\*\*

Action: 1 (UP)

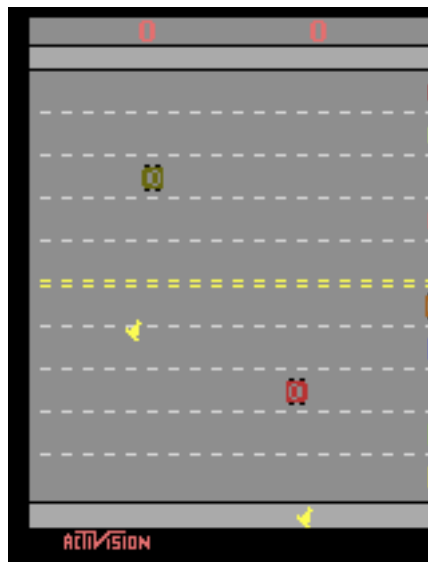
Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 3200, 'frame\_number': 3200}

Score: 0.0



\*\*\*\*\* Step 900 \*\*\*\*\*

Action: 0 (NOOP)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 3600, 'frame\_number': 3600}

Score: 0.0



\*\*\*\*\* Step 1000 \*\*\*\*\*

Action: 1 (UP)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 4000, 'frame\_number': 4000}

Score: 0.0



\*\*\*\*\* Step 1100 \*\*\*\*\*

Action: 2 (DOWN)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 4400, 'frame\_number': 4400}

Score: 0.0



\*\*\*\*\* Step 1200 \*\*\*\*\*

Action: 2 (DOWN)

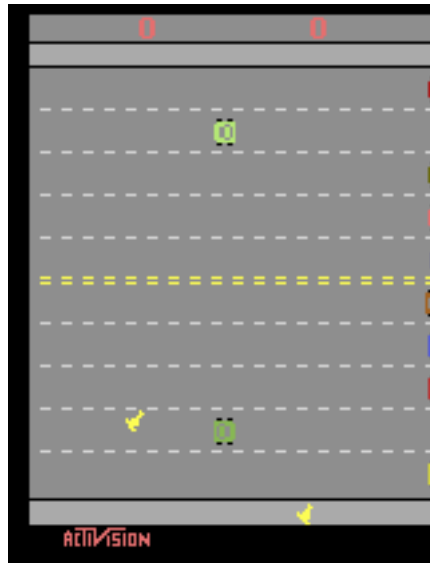
Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 4800, 'frame\_number': 4800}

Score: 0.0



\*\*\*\*\* Step 1300 \*\*\*\*\*

Action: 1 (UP)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 5200, 'frame\_number': 5200}

Score: 0.0





\*\*\*\*\* Step 1400 \*\*\*\*\*

Action: 1 (UP)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 5600, 'frame\_number': 5600}

Score: 0.0



\*\*\*\*\* Step 1500 \*\*\*\*\*

Action: 2 (DOWN)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 6000, 'frame\_number': 6000}

Score: 0.0



\*\*\*\*\* Step 1600 \*\*\*\*\*

Action: 0 (NOOP)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 6400, 'frame\_number': 6400}

Score: 0.0



\*\*\*\*\* Step 1700 \*\*\*\*\*

Action: 1 (UP)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 6800, 'frame\_number': 6800}

Score: 0.0



\*\*\*\*\* Step 1800 \*\*\*\*\*

Action: 0 (NOOP)

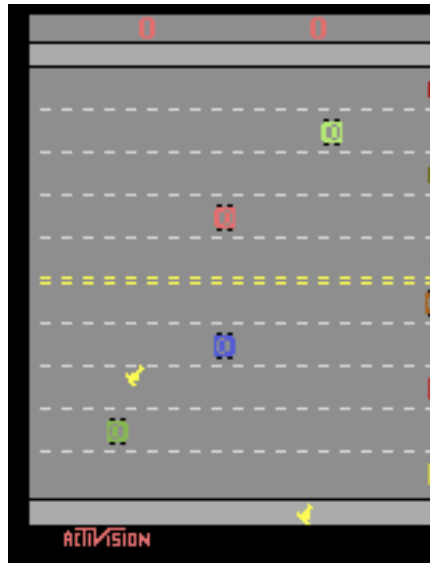
Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 7200, 'frame\_number': 7200}

Score: 0.0



\*\*\*\*\* Step 1900 \*\*\*\*\*

Action: 2 (DOWN)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 7600, 'frame\_number': 7600}

Score: 0.0



\*\*\*\*\* Step 2000 \*\*\*\*\*

Action: 2 (DOWN)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 8000, 'frame\_number': 8000}

Score: 0.0



## 6 Reinforcement Learning Agent

Now that we have a way for an agent to interact with our environment, we can go ahead and develop a RL agent capable of playing the game Crossy Road.

This is a fairly complicated game, so [ordinary Q-learning](#) won't necessarily work.

Instead, we will use a variation called [deep Q-learning](#) that utilizes deep Q-networks (DQNs).

It works similarly to ordinary Q-learning, except it approximates the Q-values for each possible action at a state using a deep neural network instead of calculating the exact values.

We will follow PyTorch's official [Reinforcement Learning \(DQN\) Tutorial](#) as a guide, and adapt it to suit our needs.

### 6.1 Defining Hyperparameters

Let's start off by defining the hyperparameters we will use.

```
[521]: hparams = {
    "memory_capacity": 10000,
    "batch_size": 128,
    "num_episodes": 4,
    "learning_rate": 1e-3, # alpha
    "exploration_rate": { # epsilon
        "start": 0.99,
        "end": 0.05,
        "decay": 1000,
    },
    "discount_factor": 0.99, # gamma
    "update_rate": 0.005, # tau
    "gradient_clip_value": 100,
}
```

## 6.2 Preprocessing Data

Before we can perform deep Q-learning, we'll need to preprocess our data to make it suitable for training.

We can specify a [transformation function](#) to transform or augment our data to fit our needs.

Here, we will:

- Turn the RGB values into an image
- Make the image grayscale
- Change the datatype to `float32`
- Flatten the image into a single dimension

```
[522]: transform = v2.Compose(
    [
        v2.ToImage(),
        v2.Lambda(lambda x: x.to(device)),
        # v2.Grayscale(),
        v2.ToDtype(dtype=torch.float32, scale=True),
        # v2.Lambda(lambda x: torch.flatten(input=x)),
    ]
)
```

## 6.3 Defining Models

Next, we can define our models, i.e. build our neural networks for use in the deep Q-learning algorithm. We will be using a deep Q-network (DQN).

Our model will be a [feed-forward neural network](#) that takes in the difference between the current and previous screen patches in order to predict the Q-values.

Essentially, the network is used to predict the expected return of taking each action given the current state.

```
[523]: class DQN(nn.Sequential):
    def __init__(self, observation_size, action_size):
        super().__init__(
            nn.Conv2d(3, 32, kernel_size=8, stride=4),
            nn.ReLU(),
            nn.Conv2d(32, 64, kernel_size=4, stride=2),
            nn.ReLU(),
            nn.Conv2d(64, 64, kernel_size=3, stride=1),
            nn.ReLU(),
            nn.Flatten(),
            nn.Linear(64 * 352, 512),
            nn.ReLU(),
            nn.Linear(512, action_size),
        )
```

## 6.4 Defining Loss Function and Optimizer

With our network defined, we can now define our loss function and optimizer.

### 6.4.1 Loss Function

We will use [smooth L1 loss](#), as is common for DQN.

```
[524]: loss_fn = nn.SmoothL1Loss()
```

### 6.4.2 Optimizer

We will use the [AdamW \(adaptive momentum estimation with weight decay\)](#) optimizer, since it is particularly robust.

It combines the benefits of momentum (such as in [SGD \(stochastic gradient descent\)](#) with momentum) and adaptive learning rates (such as in [RMSprop \(root mean square propagation\)](#)).

```
[525]: optimizer = optim.AdamW
```

## 6.5 Defining Experience Replay Memory

To perform deep Q-learning, we'll rely on a technique known as [experience replay memory](#).

### 6.5.1 Representing a Transition

First, we'll represent a transition using a tuple of the form (observation, action, next\_observation, reward), which maps a (observation, action) pair to the corresponding result (next\_observation, reward).

```
[526]: Transition = namedtuple(
    "Transition", ("observation", "action", "next_observation", "reward")
)
```

### 6.5.2 Representing Replay Memory

Next, we'll represent replay memory using a deque (double-ended queue) of fixed capacity. It will store the transitions observed most recently.

We'll also have a method for selecting a random batch of transitions.

```
[527]: class ReplayMemory:
    def __init__(self, capacity):
        self.memory = deque([], maxlen=capacity)

    def push(self, *args):
        self.memory.append(Transition(*args))

    def sample(self, batch_size):
        return random.sample(self.memory, batch_size)

    def __len__(self):
        return len(self.memory)
```

## 6.6 Deep Q-Network Agent

Now, we are ready to build our DQN agent.

### 6.6.1 Initializing the Agent

In the `__init__()` method, we do the following:

- Define the hyperparameters using the given values:
  - `hparams`: the dictionary of hyperparameters
- Preprocess the data to form a suitable training environment `env`:
  - Wrap the given environment `env` in a `TransformObservation` wrapper that transforms the observation using the given transformation function `transform`.
  - Calculate the relevant sizes of the environment:
    - \* `observation_size`: the size of an observation in the observation space
    - \* `action_size`: the size of an action in the action space
- Define the models (neural networks):
  - `policy_net`: used for action selection (to select the best action to take at a given state, i.e. the action with the highest Q-value)
  - `target_net`: used for action evaluation (to calculate the target Q-value of taking that action at the given state)
    - \* Note that `target_net` is just a copy of `policy_net` that's updated less frequently.
- Define the loss function and optimizer using the given values:
  - `loss_fn`: the loss function
  - `optimizer`: the optimizer
- Define the experience replay memory:



- memory: the memory
- Define variables used for training:
  - checkpoint\_save\_interval: how often to save a checkpoint
  - total\_steps: the total number of steps (the number of times an action was sampled)
  - episode\_scores: the total score for each episode (the number of times the chicken made it across the road)
  - episode\_losses: the loss for each episode
  - loss\_metric: the metric used to calculate the loss for each episode

```
[528]: class DQNAgent:
    def __init__(self, hparams, env, transform, loss_fn, optimizer):
        # defining hyperparameters
        self.hparams = copy.deepcopy(hparams)

        # preprocessing data
        self.env = TransformObservation(env, transform)
        observation, _ = self.env.reset()
        self.observation_size = len(observation)
        self.action_size = self.env.action_space.n

        # defining models
        self.policy_net = DQN(self.observation_size, self.action_size).
        ↪to(device)
        self.target_net = DQN(self.observation_size, self.action_size).
        ↪to(device)
        self.target_net.load_state_dict(self.policy_net.state_dict())

        # defining loss function and optimizer
        self.loss_fn = loss_fn
        self.optimizer = optimizer(
            params=self.policy_net.parameters(),
            lr=self.hparams["learning_rate"],
            amsgrad=True,
        )

        # defining experience replay memory
        self.memory = ReplayMemory(self.hparams["memory_capacity"])

        # defining variables for training
        self.checkpoint_save_interval = 2
        self.total_steps = 0
        self.episode_scores = []
        self.episode_losses = []
        self.loss_metric = Mean()
```

### 6.6.2 Sampling an Action from the Policy

In the `sample_action()` method, the agent will sample an action using an epsilon-greedy policy, which means balancing the amount of exploration vs. exploitation as per the hyperparameter for exploration rate (epsilon,  $\epsilon$ ).

Note that we are using epsilon decay, so the exploration rate will decrease exponentially over time.

```
[529]: class DQNAgent(DQNAgent):
        def sample_action(self, observation):
            # epsilon
            exploration_rate = self.hparams["exploration_rate"]["end"] + (
                self.hparams["exploration_rate"]["start"]
                - self.hparams["exploration_rate"]["end"]
            ) * math.exp(
                -1.0 * self.total_steps / self.hparams["exploration_rate"]["decay"]
            )
            self.total_steps += 1

            if np.random.uniform(0, 1) < exploration_rate: # explore
                action = self.env.action_space.sample()
            else: # exploit
                action = self.policy_net(observation).argmax(dim=1).item()

            # convert action to tensor
            return torch.tensor(
                data=[[action]],
                dtype=torch.long,
                device=device,
            )
```

### 6.6.3 Plotting Values

To visualize our training progress, we will create some plots.

**Episode Scores** First, we will create a `plot_episode_scores()` method that plots the total score for each episode (the number of times the chicken made it across the road).

```
[530]: class DQNAgent(DQNAgent):
        def plot_episode_scores(self):
            # convert episode scores to tensor
            episode_scores_tensor = torch.tensor(self.episode_scores, dtype=torch.
            ↪float32)

            # set up figure
            plt.figure()
            plt.title("Episode Scores")
            plt.xlabel("Episode")
            plt.ylabel("Score")
```

```

        # make x-axis use integers
        plt.gca().xaxis.set_major_locator(MaxNLocator(integer=True,
↪min_n_ticks=0))

        # plot episode scores
        plt.plot(episode_scores_tensor.numpy())

```

**Episode Losses** Next, we will create a `plot_episode_losses()` method that plots the loss for each episode.

```

[531]: class DQNAgent(DQNAgent):
        def plot_episode_losses(self):
            # convert episode losses to tensor
            episode_losses_tensor = torch.tensor(self.episode_losses, dtype=torch.
↪float32)

            # set up figure
            plt.figure()
            plt.title("Episode Losses")
            plt.xlabel("Episode")
            plt.ylabel("Loss")

            # make x-axis use integers
            plt.gca().xaxis.set_major_locator(MaxNLocator(integer=True,
↪min_n_ticks=0))

            # plot episode losses
            plt.plot(episode_losses_tensor.numpy())

```

#### 6.6.4 Optimizing the Policy Network

In the `optimize_policy_net()` method, we update the policy network's weights.

We do the following:

- Sample a batch of **Transitions** from the experience replay memory.
- Convert the batch-array of **Transitions** into a **Transition** of batch-arrays.
- Concatenate all of the tensors, accounting for terminal observations as necessary.
- Compute the Q-values (utility for observation-action pairs).
- Compute the V-values (utility for next observations) and the corresponding expected Q-values.
- Compute the loss and optimize the policy network's weights.

```

[532]: class DQNAgent(DQNAgent):
        def optimize_policy_net(self):
            # sample batch of transitions from memory
            if len(self.memory) < self.hparams["batch_size"]:

```

```

        return
    transitions = self.memory.sample(self.hparams["batch_size"])

    # convert batch-array of Transitions into Transition of batch-arrays
    batch = Transition(*zip(*transitions, strict=True))

    # concatenate tensors
    # observation, action, reward
    observation_batch = torch.cat(batch.observation)
    action_batch = torch.cat(batch.action)
    reward_batch = torch.cat(batch.reward)
    # nonterminal next observation
    nonterminal_next_observation_mask = torch.tensor(
        tuple(map(lambda obs: obs is not None, batch.next_observation)),
        dtype=torch.bool,
        device=device,
    )
    nonterminal_next_observation_batch = torch.cat(
        [obs for obs in batch.next_observation if obs is not None]
    )

    # compute Q-values
    observation_action_values = self.policy_net(observation_batch).gather(
        dim=1, index=action_batch
    )

    # compute V-values and corresponding expected Q-values
    next_observation_values = torch.zeros(
        size=(self.hparams["batch_size"],), device=device
    )
    with torch.no_grad():
        next_observation_values[nonterminal_next_observation_mask] = (
            self.target_net(nonterminal_next_observation_batch).max(dim=1).
↪ values
        )
    expected_observation_action_values = (
        next_observation_values * self.hparams["discount_factor"]
    ) + reward_batch

    # compute loss and optimize policy net
    loss = self.loss_fn(
        observation_action_values,
        expected_observation_action_values.unsqueeze(dim=1),
    )
    self.loss_metric.update(
        input=loss.detach().cpu(), weight=self.hparams["batch_size"]
    )

```

```

self.optimizer.zero_grad()
loss.backward()
nn.utils.clip_grad_value_(
    parameters=self.policy_net.parameters(), clip_value=100
)
self.optimizer.step()

```

### 6.6.5 Optimizing the Target Network

In the `optimize_target_net()` method, we update the target network’s weights. It’s a soft update based on the hyperparameter for update rate ( $\tau$ ,  $\tau$ ).

```

[533]: class DQNAgent(DQNAgent):
        def optimize_target_net(self):
            # get state dictionaries
            policy_net_state_dict = self.policy_net.state_dict()
            target_net_state_dict = self.target_net.state_dict()

            # update target net state dictionary
            for key in policy_net_state_dict:
                old_weight = (1 - self.hparams["update_rate"]) *
↪target_net_state_dict[key]
                new_weight = self.hparams["update_rate"] *
↪policy_net_state_dict[key]
                target_net_state_dict[key] = old_weight + new_weight

            # load target net state dictionary into target net
            self.target_net.load_state_dict(target_net_state_dict)

```

### 6.6.6 Training the Agent

In the `train()` method, we train our agent in order to make its policy “better” at solving the environment.

We also handle saving and loading of the models and plots as necessary.

```

[534]: class DQNAgent(DQNAgent):
        def train(self):
            start_episode = 0
            end_episode = self.hparams["num_episodes"]

            # load checkpoint if available
            if Path("checkpoint.tar").is_file():
                checkpoint = torch.load("checkpoint.tar")
                start_episode = checkpoint["start_episode"]
                self.policy_net.load_state_dict(checkpoint["policy_net"])
                self.target_net.load_state_dict(checkpoint["target_net"])
                self.optimizer.load_state_dict(checkpoint["optimizer"])

```

```

self.memory.memory = checkpoint["memory"]
self.total_steps = checkpoint["total_steps"]
self.episode_scores = checkpoint["episode_scores"]
self.episode_losses = checkpoint["episode_losses"]

for episode in trange(
    start_episode,
    end_episode,
    desc="Episodes",
    initial=start_episode,
    total=end_episode,
):
    # reset environment
    display.clear_output(wait=True)
    step = 0
    observation, _ = self.env.reset()
    observation = observation.clone().detach().unsqueeze(dim=0)
    reward, terminated, truncated = 0.0, False, False
    score = 0
    self.loss_metric.reset()

    while not (terminated or truncated):
        step += 1

        # sample action from policy
        action = self.sample_action(observation)
        action = action.clone().detach()
        next_observation, reward, terminated, truncated, _ = self.env.
↪step(
            action.item()
        )
        score += reward
        next_observation = (
            None
            if terminated
            else next_observation.clone().detach().unsqueeze(0)
        )
        reward = torch.tensor(data=[reward], dtype=torch.float32, ↵
↪device=device)
        reward = reward.clone().detach()

        # store transition in memory
        self.memory.push(observation, action, next_observation, reward)

        # move to next observation
        observation = next_observation

```

```

        # optimize model
        self.optimize_policy_net()
        self.optimize_target_net()

    # log episode score and update plots
    self.episode_scores.append(score)
    self.plot_episode_scores()
    plt.show()

    # log episode loss and update plots
    self.episode_losses.append(self.loss_metric.compute().item())
    self.plot_episode_losses()
    plt.show()

    # save checkpoint on interval or at the end
    if (episode + 1) % self.checkpoint_save_interval == 0 or (
        episode + 1
    ) == self.hparams["num_episodes"]:
        # save checkpoint
        Path("checkpoints/").mkdir(parents=True, exist_ok=True)
        checkpoint = {
            "timestamp": datetime.now().isoformat(timespec="minutes"),
            "start_episode": episode + 1,
            "policy_net": self.policy_net.state_dict(),
            "target_net": self.target_net.state_dict(),
            "optimizer": self.optimizer.state_dict(),
            "memory": self.memory.memory,
            "total_steps": self.total_steps,
            "episode_scores": self.episode_scores,
            "episode_losses": self.episode_losses,
        }
        torch.save(checkpoint, "checkpoint.tar")
        torch.save(
            checkpoint,
            f"checkpoints/
↪{checkpoint['timestamp']}-{checkpoint['start_episode']}.tar",
        )

    # clear output
    display.clear_output(wait=True)

    # plot episode scores
    self.plot_episode_scores()
    plt.show()

    # plot episode losses
    self.plot_episode_losses()

```

```
plt.show()
```

### 6.6.7 Inference Mode

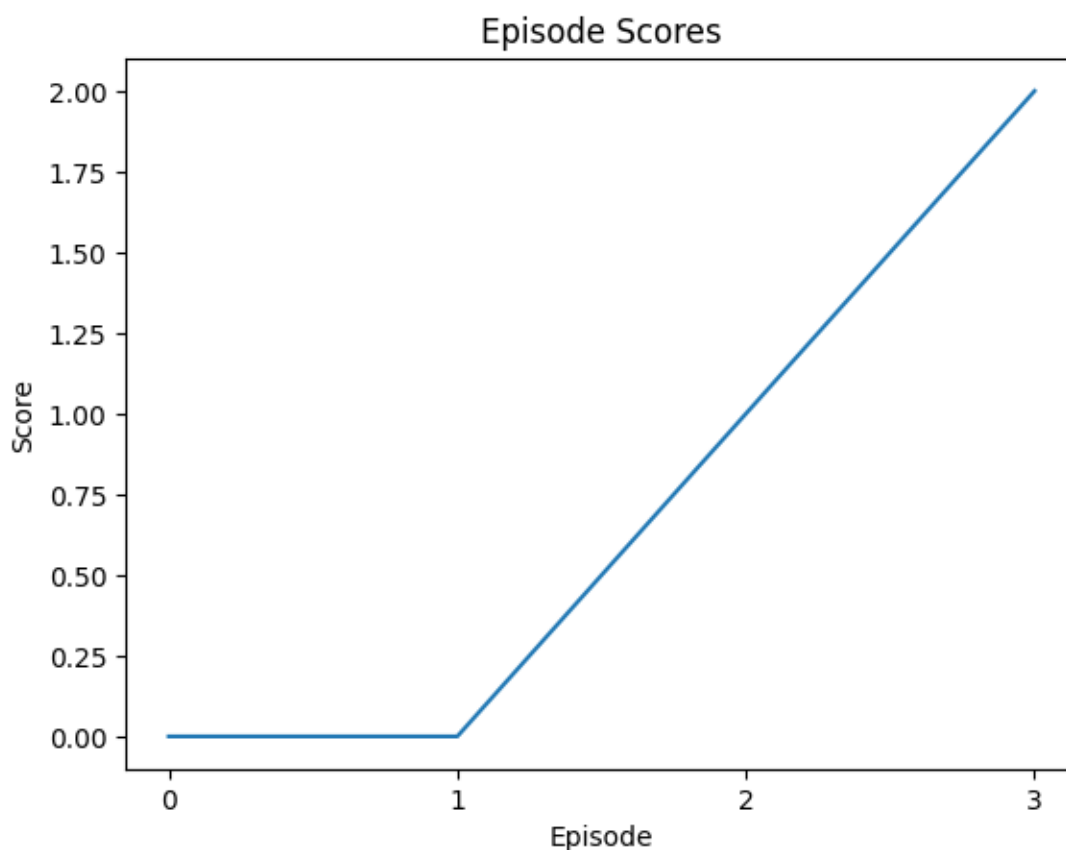
In the `inference_mode()` method, we set the hyperparameters for exploration rate (epsilon,  $\epsilon$ ) to 0, so we don't explore anymore and only exploit.

```
[535]: class DQNAgent(DQNAgent):  
        def inference_mode(self):  
            self.hparams["exploration_rate"]["start"] = 0.0  
            self.hparams["exploration_rate"]["end"] = 0.0
```

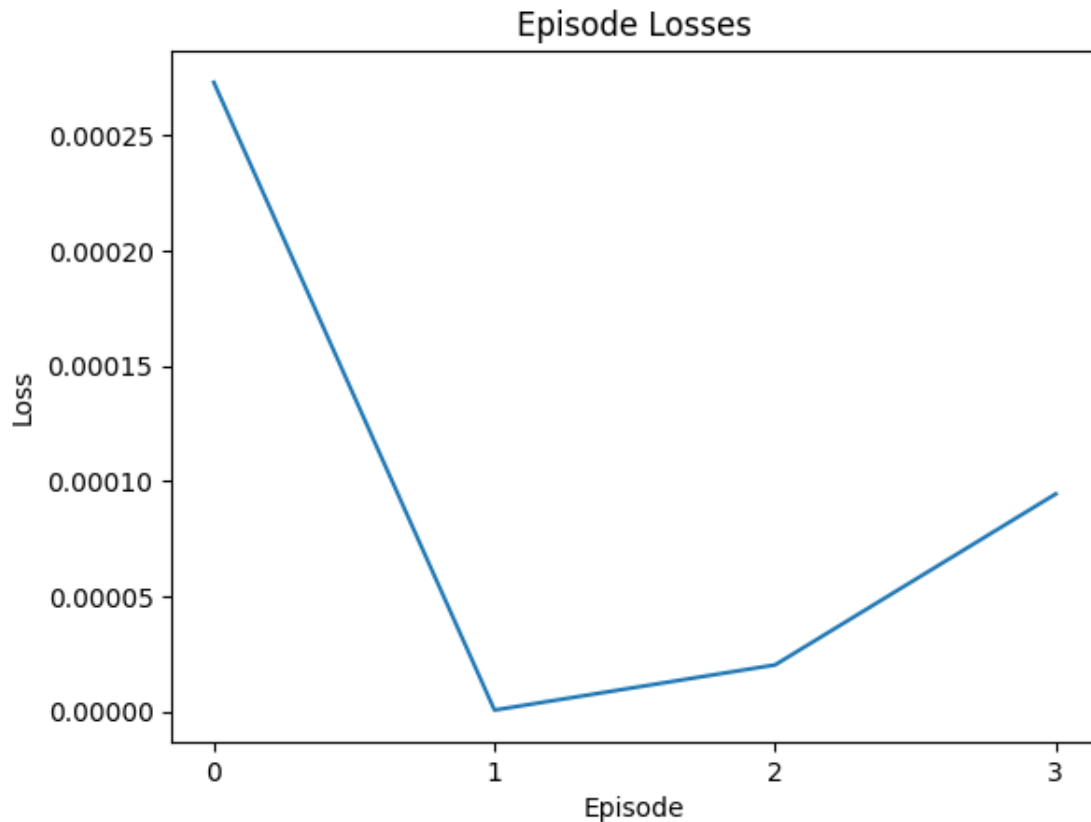
## 6.7 Training the Agent

Now, we're equipped with everything we need to train our DQN agent.

```
[536]: dqn_agent = DQNAgent(hparams, env, transform, loss_fn, optimizer)  
dqn_agent.train()
```







## 6.8 Simulation

Now, let's run a simulation.

To limit the amount of output we need to scroll through, we will only show every 100 steps.

```
[537]: dqn_agent.inference_mode()
        simulate(env, transform, dqn_agent, 1, 100)
```

```
##### Episode 0 #####
```

```
***** Step 0 *****
```

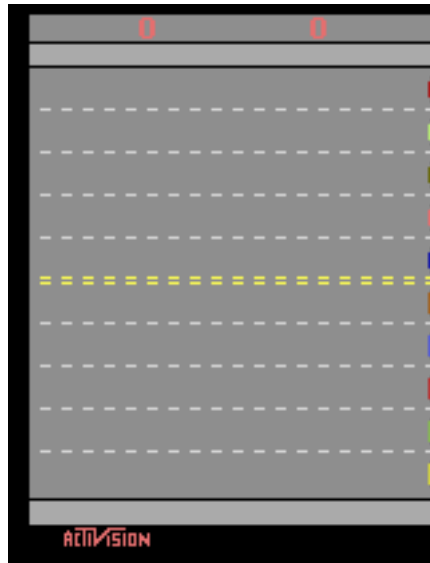
```
Reward: 0.0
```

```
Terminated: False
```

```
Truncated: False
```

```
Info: {'lives': 0, 'episode_frame_number': 0, 'frame_number': 8192}
```

```
Score: 0
```



\*\*\*\*\* Step 100 \*\*\*\*\*

Action: 2 (DOWN)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 400, 'frame\_number': 8592}

Score: 0.0



\*\*\*\*\* Step 200 \*\*\*\*\*

Action: 0 (NOOP)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 800, 'frame\_number': 8992}

Score: 0.0



\*\*\*\*\* Step 300 \*\*\*\*\*

Action: 0 (NOOP)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 1200, 'frame\_number': 9392}

Score: 0.0



\*\*\*\*\* Step 400 \*\*\*\*\*

Action: 2 (DOWN)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 1600, 'frame\_number': 9792}

Score: 0.0



\*\*\*\*\* Step 500 \*\*\*\*\*

Action: 2 (DOWN)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 2000, 'frame\_number': 10192}

Score: 0.0



\*\*\*\*\* Step 600 \*\*\*\*\*

Action: 0 (NOOP)

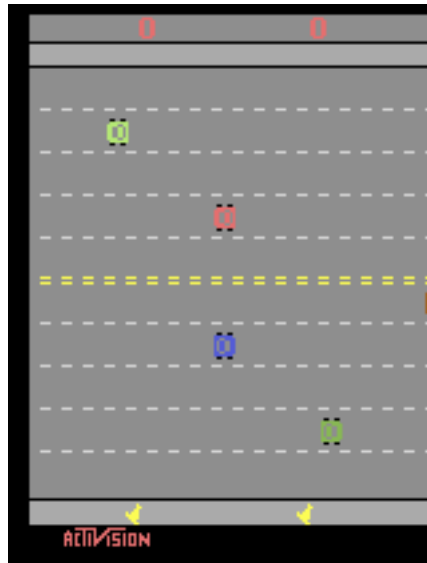
Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 2400, 'frame\_number': 10592}

Score: 0.0



\*\*\*\*\* Step 700 \*\*\*\*\*

Action: 0 (NOOP)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 2800, 'frame\_number': 10992}

Score: 0.0



\*\*\*\*\* Step 800 \*\*\*\*\*

Action: 2 (DOWN)

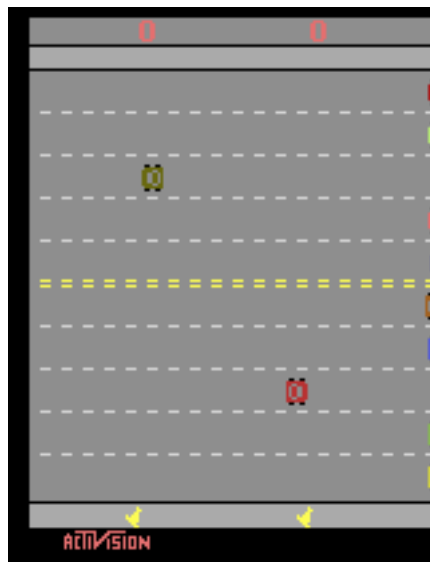
Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 3200, 'frame\_number': 11392}

Score: 0.0



\*\*\*\*\* Step 900 \*\*\*\*\*

Action: 0 (NOOP)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 3600, 'frame\_number': 11792}

Score: 0.0



\*\*\*\*\* Step 1000 \*\*\*\*\*

Action: 0 (NOOP)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 4000, 'frame\_number': 12192}

Score: 0.0





\*\*\*\*\* Step 1100 \*\*\*\*\*

Action: 0 (NOOP)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 4400, 'frame\_number': 12592}

Score: 0.0



\*\*\*\*\* Step 1200 \*\*\*\*\*

Action: 0 (NOOP)

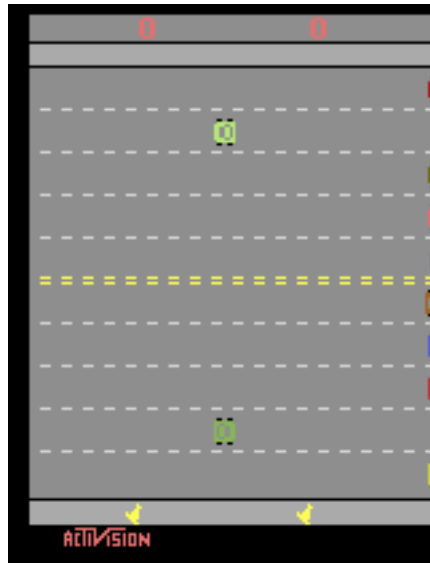
Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 4800, 'frame\_number': 12992}

Score: 0.0



\*\*\*\*\* Step 1300 \*\*\*\*\*

Action: 2 (DOWN)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 5200, 'frame\_number': 13392}

Score: 0.0



\*\*\*\*\* Step 1400 \*\*\*\*\*

Action: 0 (NOOP)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 5600, 'frame\_number': 13792}

Score: 0.0



\*\*\*\*\* Step 1500 \*\*\*\*\*

Action: 0 (NOOP)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 6000, 'frame\_number': 14192}

Score: 0.0



\*\*\*\*\* Step 1600 \*\*\*\*\*

Action: 2 (DOWN)

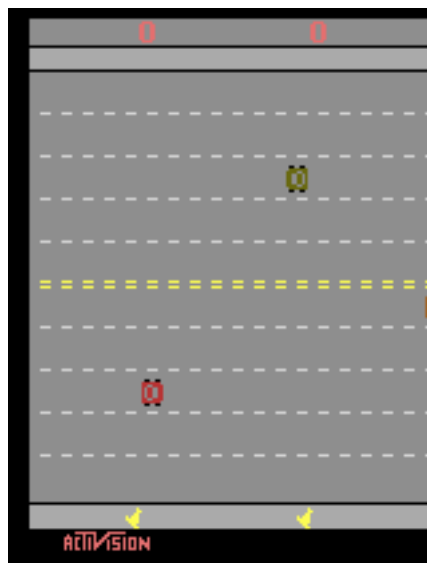
Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 6400, 'frame\_number': 14592}

Score: 0.0



\*\*\*\*\* Step 1700 \*\*\*\*\*

Action: 2 (DOWN)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 6800, 'frame\_number': 14992}

Score: 0.0



\*\*\*\*\* Step 1800 \*\*\*\*\*

Action: 0 (NOOP)

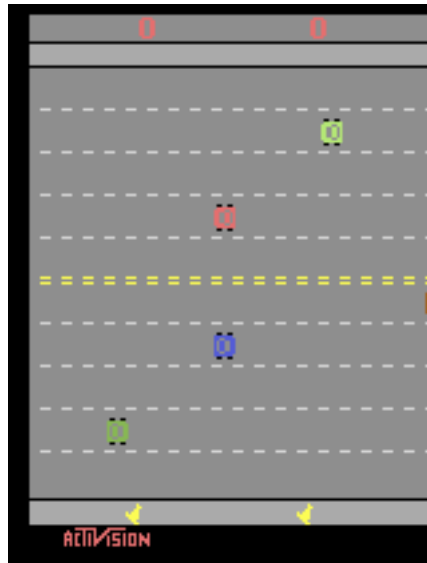
Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 7200, 'frame\_number': 15392}

Score: 0.0



\*\*\*\*\* Step 1900 \*\*\*\*\*

Action: 0 (NOOP)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 7600, 'frame\_number': 15792}

Score: 0.0



\*\*\*\*\* Step 2000 \*\*\*\*\*

Action: 2 (DOWN)

Reward: 0.0

Terminated: False

Truncated: False

Info: {'lives': 0, 'episode\_frame\_number': 8000, 'frame\_number': 16192}

Score: 0.0



## 7 Conclusion

As we can see, using the power of deep Q-learning, we were able to develop a RL game agent capable of playing Crossy Road that performs significantly better than a random agent.

While the agent is not perfect, it has impressive performance, considering the complexity of the game and the compute constraints we were under.