# PACMAN + DQNs

Elijah Tai

# Contents

Topic + Background

Hypothesis + Importance

Progress from Midterm Report

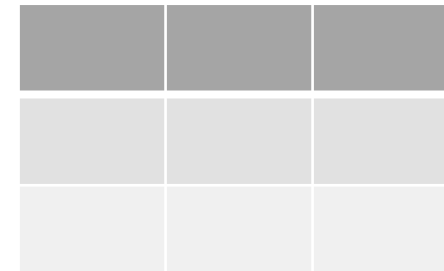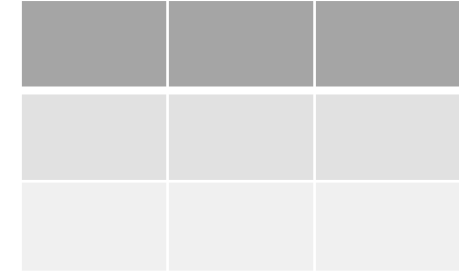Results + Code

Issues + Solutions

Next Steps

# Contents

$$Q^{new}(s_t, a_t) \leftarrow \underbrace{Q(s_t, a_t)}_{\text{current value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \left( \underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \overbrace{\underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}}}^{\text{temporal difference}} - \underbrace{Q(s_t, a_t)}_{\text{current value}} \right)$$

$$\underbrace{\phantom{r_t + \gamma \cdot \max_a Q(s_{t+1}, a)}}_{\text{new value (temporal difference target)}}$$

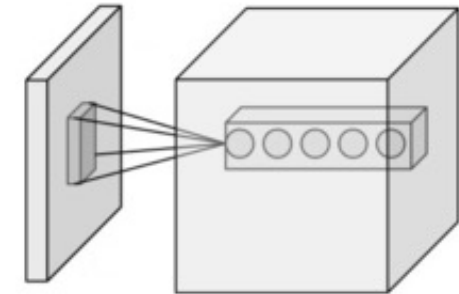Watkins, C.J.C.H., Dayan, P. *Q*–learning. *Mach Learn* **8**, 279–292 (1992). https://doi.org/10.1007/BF00992698

$$Q^{new}(s_t, a_t) \leftarrow \underbrace{Q(s_t, a_t)}_{\text{current value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \bigg( \overbrace{\underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}} - \underbrace{Q(s_t, a_t)}_{\text{current value}}}^{\text{temporal difference}} \bigg)$$

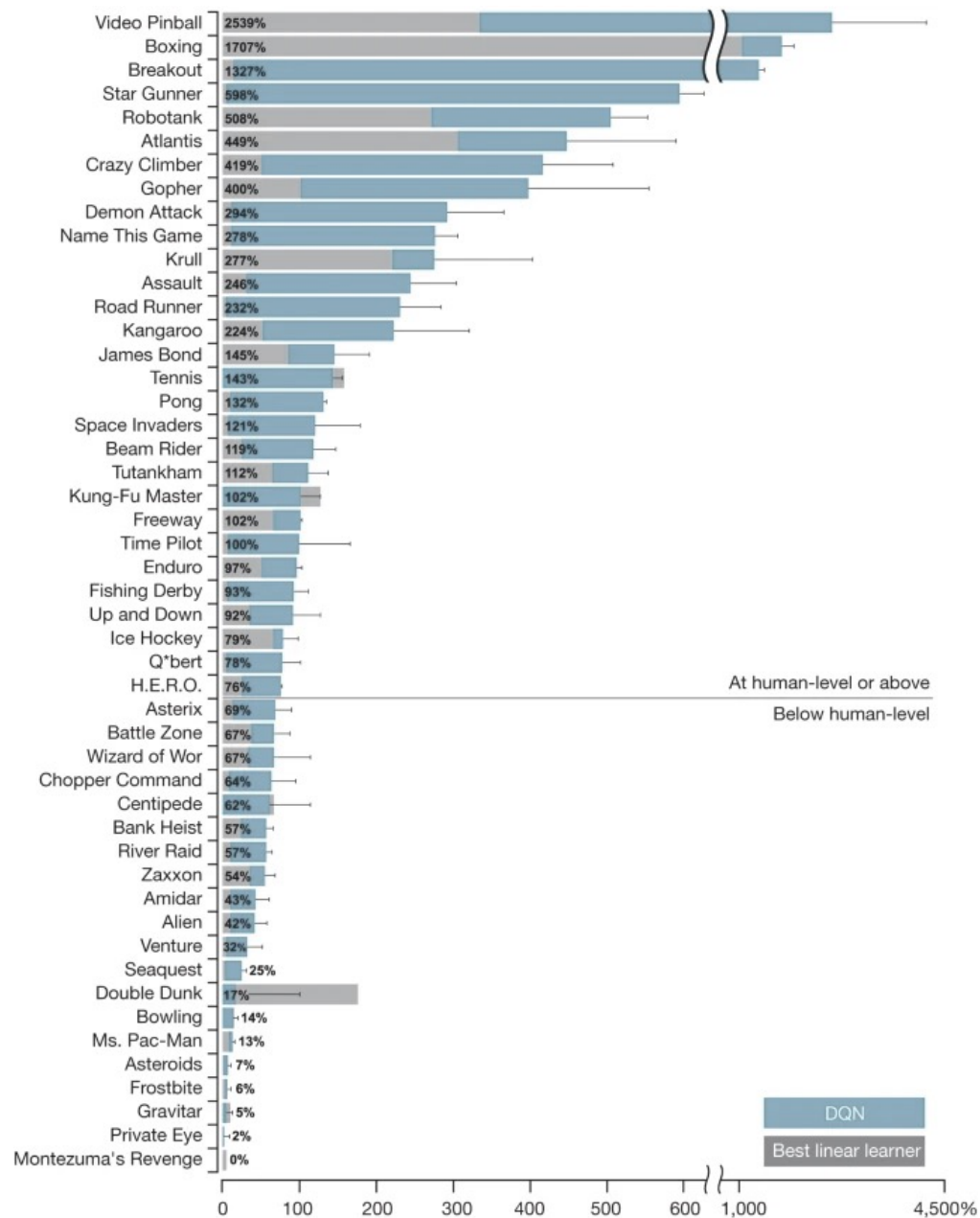$$\underbrace{\qquad\qquad\qquad\qquad\qquad}_{\text{new value (temporal difference target)}}$$

$$Q^A_{t+1}(s_t, a_t) = Q^A_t(s_t, a_t) + \alpha_t(s_t, a_t) \left( r_t + \gamma Q^B_t \left( s_{t+1}, \arg\max_a Q^A_t(s_{t+1}, a) \right) - Q^A_t(s_t, a_t) \right)$$

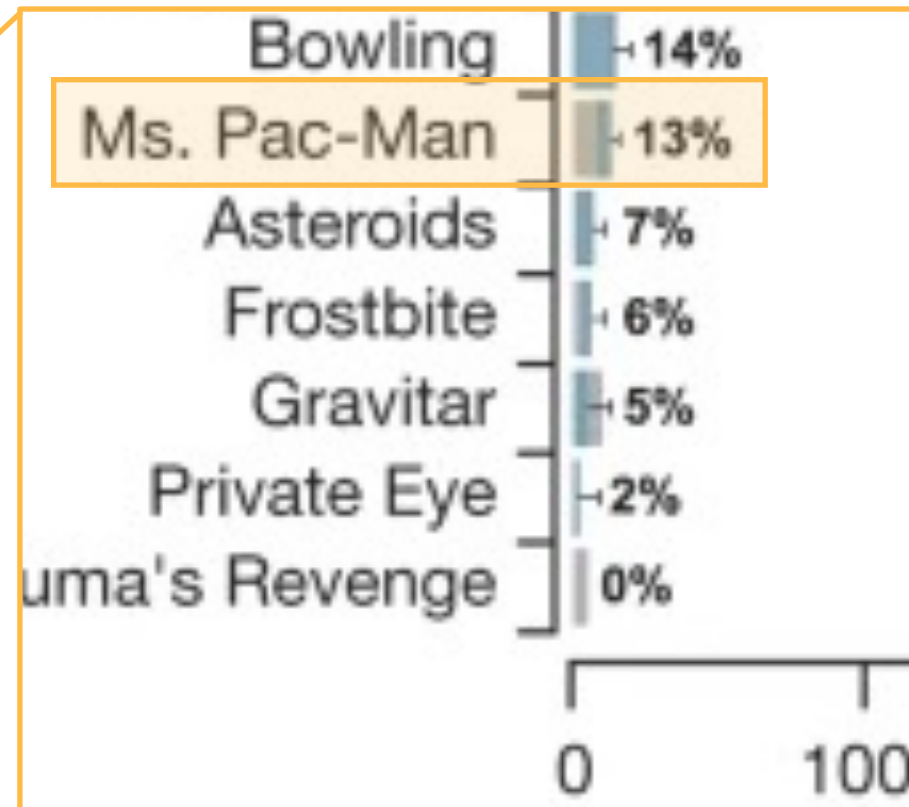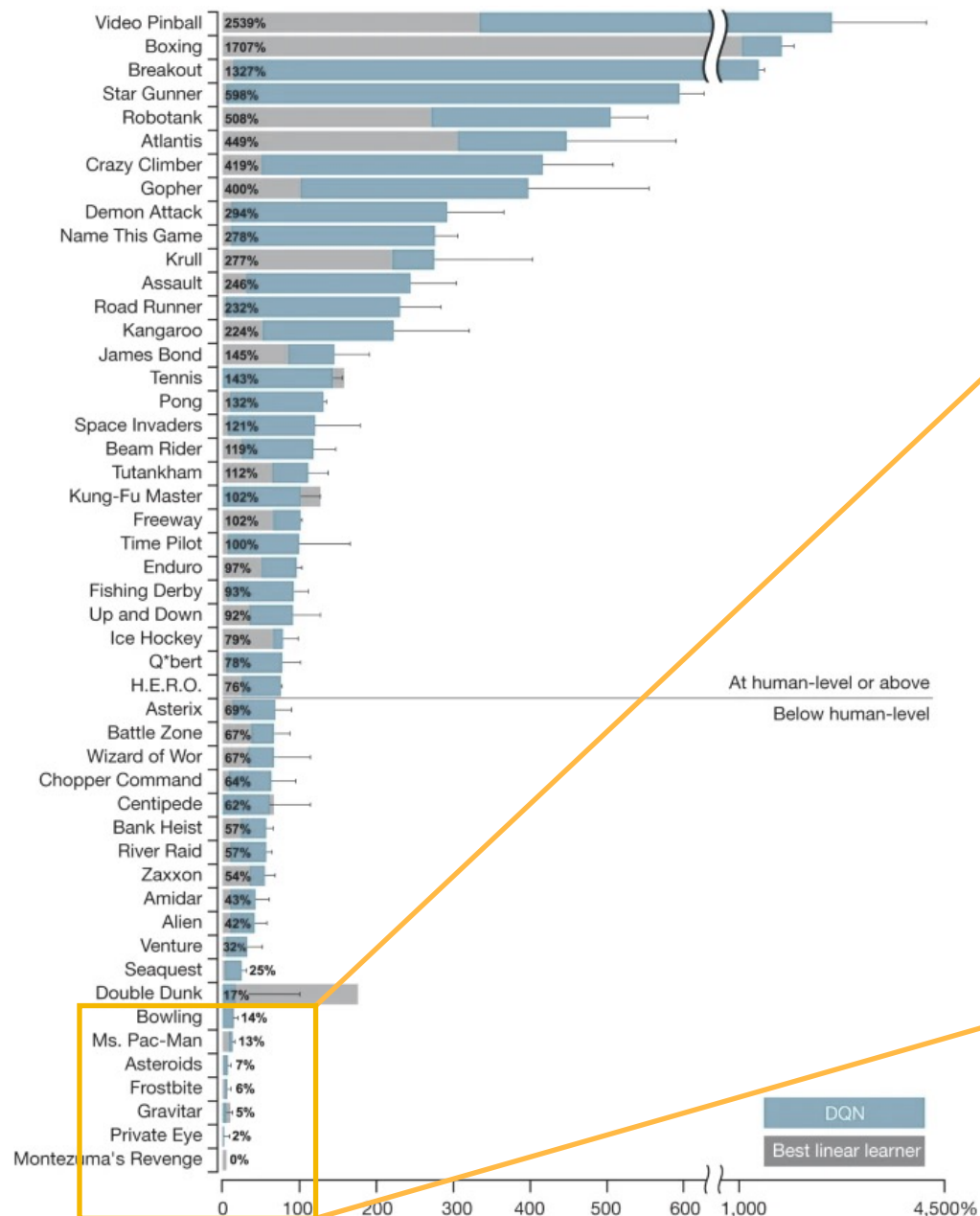Watkins, C.J.C.H., Dayan, P. *Q*–learning. *Mach Learn* **8**, 279–292 (1992). https://doi.org/10.1007/BF00992698

Volodymyr Mnih et al. "Playing atari with deep reinforcement learning". *arXiv preprint arXiv:1312.5602* (2013).

| Game | DQN (%) |
|---|---|
| Video Pinball | 2539% |
| Boxing | 1707% |
| Breakout | 1327% |
| Star Gunner | 598% |
| Robotank | 508% |
| Atlantis | 449% |
| Crazy Climber | 419% |
| Gopher | 400% |
| Demon Attack | 294% |
| Name This Game | 278% |
| Krull | 277% |
| Assault | 246% |
| Road Runner | 232% |
| Kangaroo | 224% |
| James Bond | 145% |
| Tennis | 143% |
| Pong | 132% |
| Space Invaders | 121% |
| Beam Rider | 119% |
| Tutankham | 112% |
| Kung-Fu Master | 102% |
| Freeway | 102% |
| Time Pilot | 100% |
| Enduro | 97% |
| Fishing Derby | 93% |
| Up and Down | 92% |
| Ice Hockey | 79% |
| Q*bert | 78% |
| H.E.R.O. | 76% |
| Asterix | 69% |
| Battle Zone | 67% |
| Wizard of Wor | 67% |
| Chopper Command | 64% |
| Centipede | 62% |
| Bank Heist | 57% |
| River Raid | 57% |
| Zaxxon | 54% |
| Amidar | 43% |
| Alien | 42% |
| Venture | 32% |
| Seaquest | 25% |
| Double Dunk | 17% |
| Bowling | 14% |
| Ms. Pac-Man | 13% |
| Asteroids | 7% |
| Frostbite | 6% |
| Gravitar | 5% |
| Private Eye | 2% |
| Montezuma's Revenge | 0% |

At human-level or above / Below human-level

Legend: DQN / Best linear learner

Mnih, V., Kavukcuoglu, K., Silver, D. *et al.* Human-level control through deep reinforcement learning. *Nature* **518**, 529–533 (2015). https://doi.org/10.1038/nature14236

# Contents

# Hypothesis

- DQN algorithm struggles in MS PACMAN because of the power-up pellets. Suddenly, the best strategy is to chase ghosts instead of fleeing. This change confuses the network.

| | | | |
|---|---|---|---|
| Removing the reward from eating ghosts during training paradoxically increases rewards. | Initializing powered-up MS PACMAN to do the opposite movement increases rewards. | Discretely context switching between two different networks will increase rewards. | Valid combinations of the other optimizations will further increase rewards. |

# Importance

- While the optimizations are specific to MS PACMAN, this generally demonstrates tactics to improve performance, especially within settings of low computational resources when the neural network may not have a chance to stably converge.

# Contents

# Progress

- Previously, I got an existing DQN running from: [github.com/bourbonut/dqn-pacman](github.com/bourbonut/dqn-pacman)

- Learned about DQNs and fixed bugs in the original code.
- Wrote a script to determine the invincibility frames of ms pacman after eating a power-pellet
- Coded different behaviors triggered by certain score rewards are triggers.
- Ran 10k frames ~65 epochs on three of the four specific hypotheses
- Generated plots for visualization

# Contents

Topic + Background

Hypothesis + Importance

Progress from Midterm Report

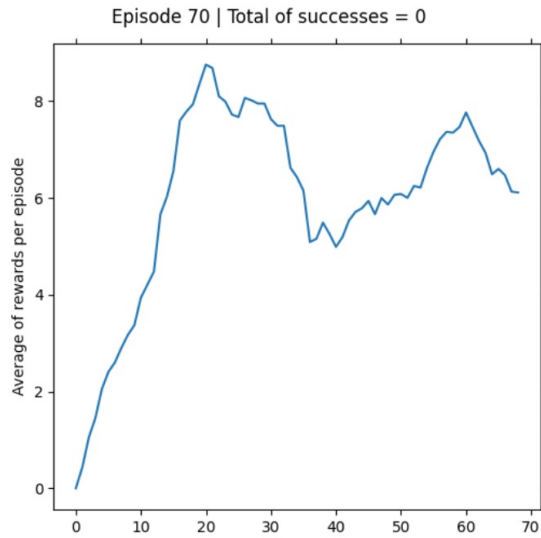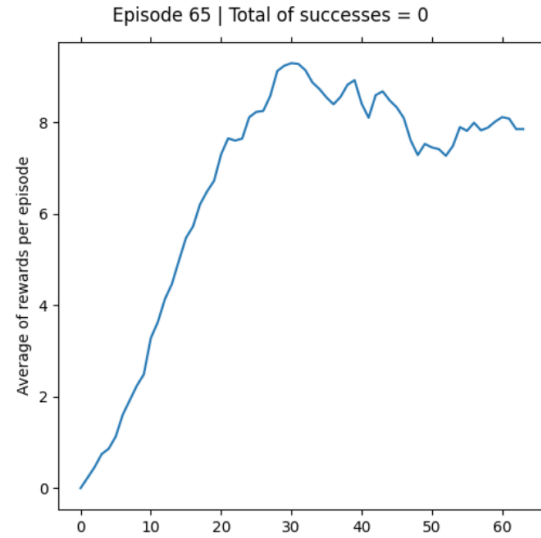**Results + Code**

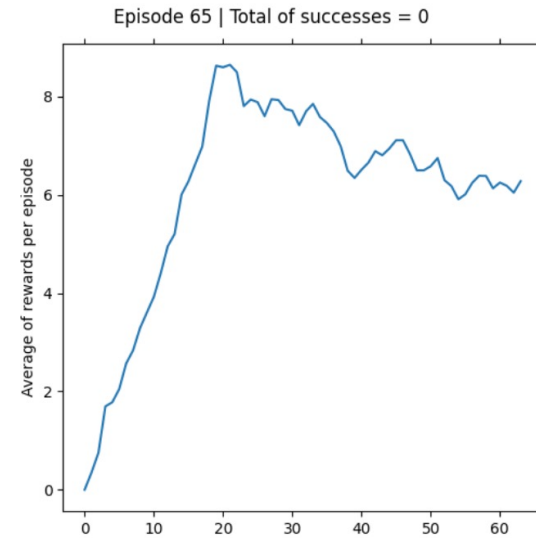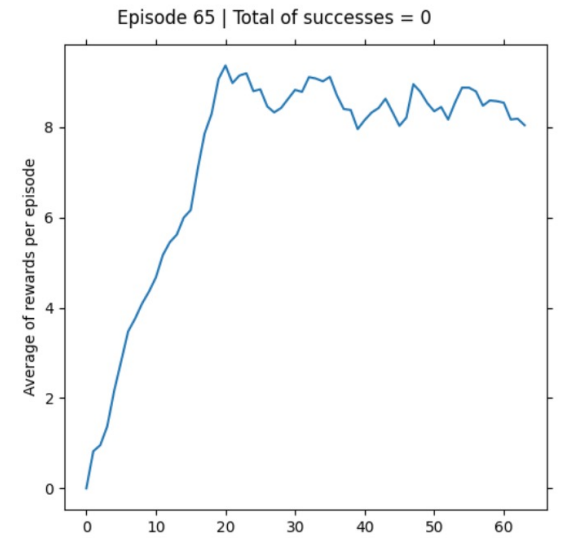Issues + Solutions

Next steps

# Results: 10k steps

Neither

No Ghost
Reward

Opposite
Initialization

Both



Episode 70 | Total of successes = 0



Episode 65 | Total of successes = 0



Episode 65 | Total of successes = 0



Episode 65 | Total of successes = 0

6.0

7.8

6.5

8.1

# Results: 10k steps
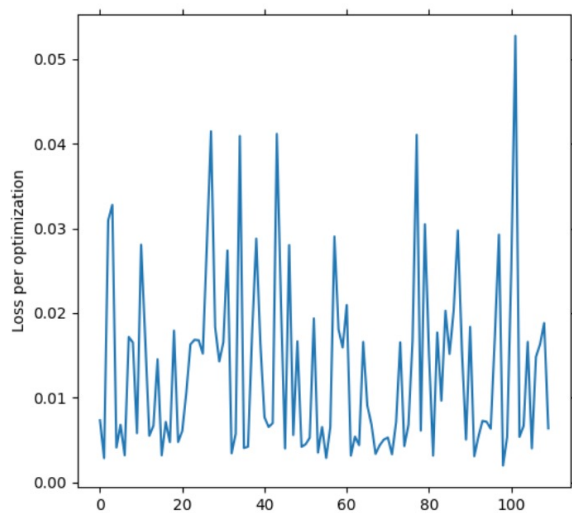
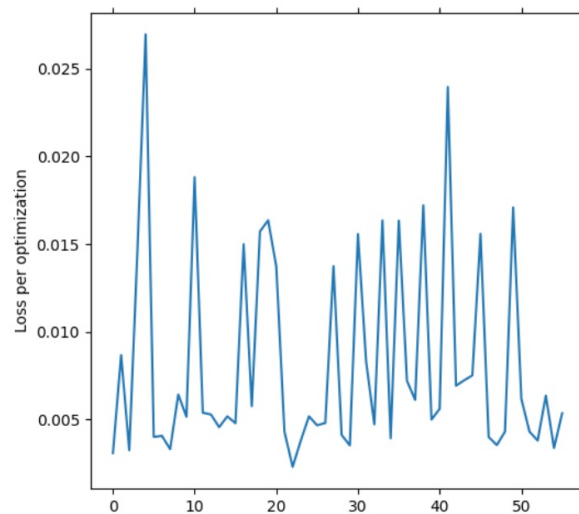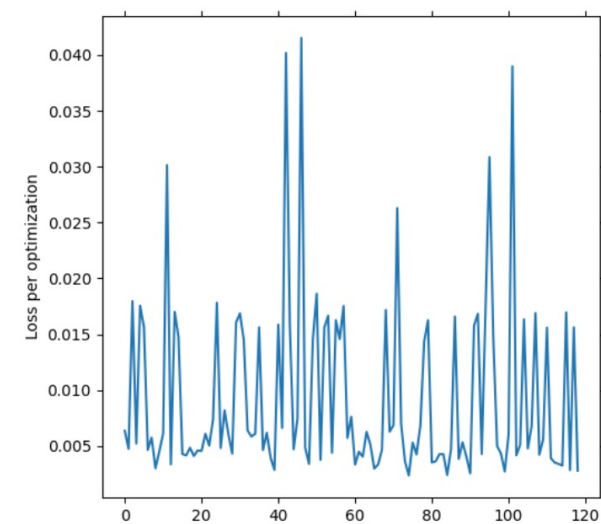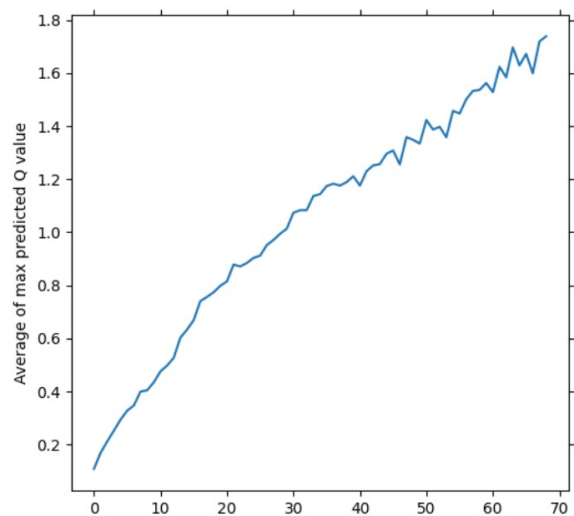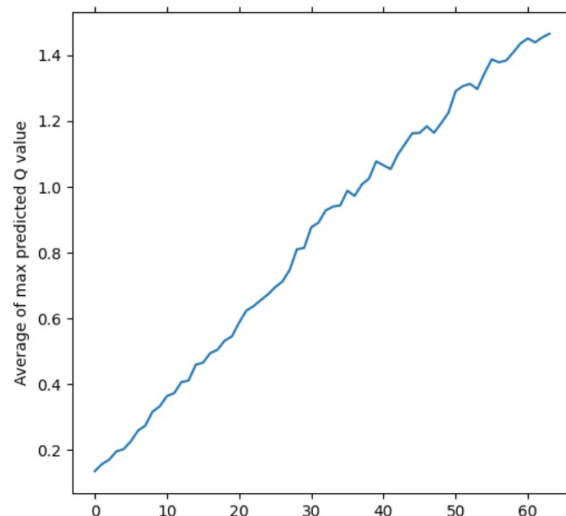| Neither | No Ghost Reward | Opposite Initialization | Both |
|---------|-----------------|-------------------------|------|
|  |  |  |  |
| 0.054 | 0.027 | 0.040 | 0.042 |

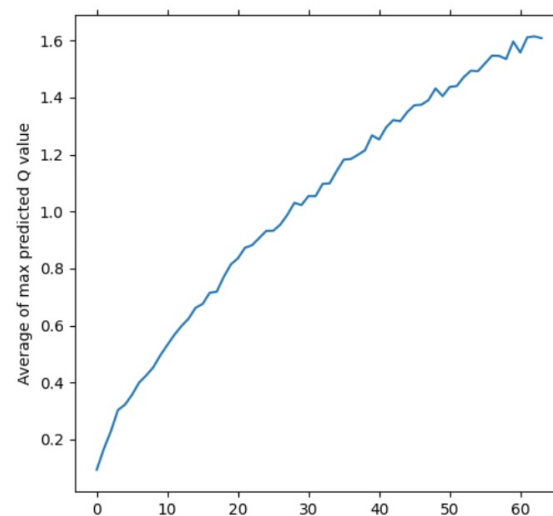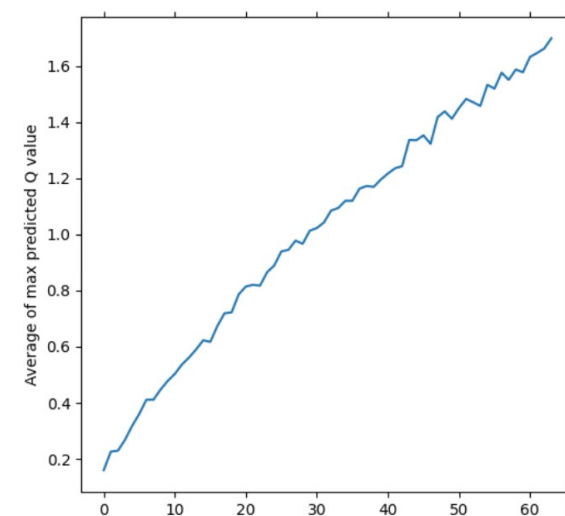# Results: 10k steps

Neither

No Ghost
Reward

Opposite
Initialization

Both



1.72

1.47

1.62

1.70

# Main Code Alterations

## Invincibility Framing

```python
1   import gym
2   import numpy as np
3   import matplotlib.pyplot as plt
4   import torch
5   import random as r
6
7   env = gym.make("MsPacman-v4", render_mode='human')
8   env.reset()
9
10  stepmode = False
11  phase = 0
12  counter = 0
13  for _ in range(10000):
14      # env.action_space.n |=> 9
15      # env.get_action_meanings() |=> ['NOOP', 'UP', 'RIGHT', 'LEFT', 'D'
16
17      if phase == 0:
18          next_state, reward, terminated, truncated, info = env.step(4)
19          counter += 1
20          if counter == 120:
21              phase += 1
22              counter = 0
23      if phase == 1:
24          next_state, reward, terminated, truncated, info = env.step(3)
25          counter += 1
26          if counter >= 80:
27              phase += 1
28              counter = 0
29      if phase == 2:
30          counter += 1
31          next_state, reward, terminated, truncated, info = env.step(1)
32          if reward == 50:
33              phase += 1
34              counter = 0
35      if phase == 3:
36          counter += 1
37          input(f'continue {counter}')
38
39          next_state, reward, terminated, truncated, info = env.step(0)
40
41      env.render()
42      print(phase)
43
44  env.close()
```

## No Ghost Reward

```python
16  parser.add_argument(
17      "--ghostbonus",
18      action="store_true",
19      dest="ghostbonus",
20      help="no reward for eating a ghost or a strawberry",
21  )
```

```python
86          # Elijah: Change the reward to ignore eating a ghost
87          display_reward = reward_
88          if args.ghostbonus and (reward_ == 200 or reward_ % 400 == 0):
89              display_reward = 0
```

```python
113             # Elijah: We want to display the true reward, but we w
114             if args.ghostbonus:
115                 display.data.rewards.append(display_reward)
116             else:
117                 display.data.rewards.append(reward)
118             reward = torch.tensor([reward], device=device)
119
120             old_action = action_
121             if reward != 0:
122                 dmaker.old_action = action.item()
123
124             next_state = preprocess_observation(observations, obs)
```

```python
119     def _save(self, image=False):
120         """Save data in `pickle` file and an image"""
121         if image:
122             PATH_PLOTS = self.data.path / '..' / 'plots'
123             self.update_axis()
124             self.fig.tight_layout()
125             plt.savefig(PATH_PLOTS / f"episode-{self.data.ep}.png")
126             print(f"Figure {self.data.ep} saved.")
127             for axis in self.axis:
128                 axis.cla()
129         self.data.save()
```

## Opposite Initialization

```python
22  parser.add_argument(
23      "--ghostrev",
24      action="store_true",
25      dest="ghostrev",
26      help="train a single NN, reversing actions under power pellet"
```

```python
77          # Elijah: Reverse the action that we do
78          if args.ghostrev and rev_counter > 0:
79              action = [0, 4, 3, 2, 1, 9, 7, 6, 5][action]
80              rev_counter -= 1
81
82          action_ = ACTIONS[old_action][action.item()] # Elijah: This sort of softens the new action.
83
84          obs, reward_, done, info = env.step(action_)
```

```python
20  REVERSED = {0: 1, 1: 0, 2: 3, 3: 2}
21  isreversed = (
22      lambda last_action, action: "default" if REVERSED[action] - last_action else "reverse"
23
24
25  ACTIONS = {
26      1: [1, 4, 6, 5],
27      2: [5, 7, 3, 2],
28      3: [6, 8, 3, 2],
29      4: [1, 4, 8, 7],
30      5: [1, 4, 3, 2],
31      6: [1, 4, 3, 2],
32      7: [1, 4, 3, 2],
33      8: [1, 4, 3, 2],
34  }
```

```python
22  TARGET_UPDATE = 400  # here, Elijah changed from 8_
23  REPLAY_MEMORY_SIZE = 3 * 1200 #here, Elijah changed
24
25  # Environment constants
26  N_ACTIONS = 4
27  AVOIDED_STEPS = 80  # At the beginning, there is a
28  DEAD_STEPS = 36  # frames to avoid when the agent d
29  K_FRAME = 2
30
31  # Optimizer parameters
32  LEARNING_RATE = 2.5e-4
33  # DECAY_RATE = 0.99
34  MOMENTUM = 0.95
35
36  # Algorithm constant
37  MAX_FRAMES = 20_000 #Elijah changed from 2,000,000
38  SAVE_MODEL = 5 # Elijah changed from 20
```

### File tree (right panel)

- ghostbonusDQN
  - models
    - policy-model-5.pt
    - policy-model-10.pt
    - policy-model-15.pt
    - policy-model-20.pt
    - policy-model-25.pt
    - policy-model-30.pt
    - policy-model-35.pt
    - policy-model-40.pt
    - policy-model-45.pt
    - policy-model-50.pt
    - policy-model-55.pt
    - policy-model-60.pt
    - policy-model-65.pt
    - policy-model-final.pt
    - target-model-5.pt
    - target-model-10.pt
    - target-model-15.pt
    - target-model-20.pt
    - target-model-25.pt
    - target-model-30.pt
    - target-model-35.pt
    - target-model-40.pt
    - target-model-45.pt
    - target-model-50.pt
    - target-model-55.pt
    - target-model-60.pt
    - target-model-65.pt
    - target-model-final.pt

- plots
  - episode-5.png
  - episode-10.png
  - episode-15.png
  - episode-20.png
  - episode-25.png
  - episode-30.png
  - episode-35.png
  - episode-40.png
  - episode-45.png
  - episode-50.png
  - episode-55.png
  - episode-60.png
  - episode-65.png
- recorded-data
  - episode-5.pkl
  - episode-10.pkl
  - episode-15.pkl
  - episode-20.pkl
  - episode-25.pkl
  - episode-30.pkl
  - episode-35.pkl
  - episode-40.pkl
  - episode-45.pkl
  - episode-50.pkl
  - episode-55.pkl
  - episode-60.pkl
  - episode-65.pkl

# Demo

# Contents

# Issues + Solutions

- Need variance metric for comparing the learning curves
  - Need to run the code multiple times to make better comparison

- Only have results for short time frames
  - Plan to let it run for ~1000 episodes

- Still need to add in the context switching (double) double DQN
  - Requires a new class and code structure changes

- This will add more quantitative rigor

# Contents

Topic + Background

Hypothesis + Importance

Progress from Midterm Report

Results + Code

Issues + Solutions

**Next Steps**

# Planning for Next Steps

- Running more times for longer

- Implementing the Double-Double DQN (the third hypothesis)

- Generating more quantitative metrics for learning curve comparison

- Extra: Implementing Dueling DQN or soft steps (Polyak averaging)

- Extra: Comparing greyscale versus color to understand ghost personalities

Thank you!