# Deep Q Learning with OpenAI

Feng Jiang, Yizhe Hu, Brandon Klein, Zijie(Jay) Wang

Background

Q Learning vs. Traditional NN with CartPole

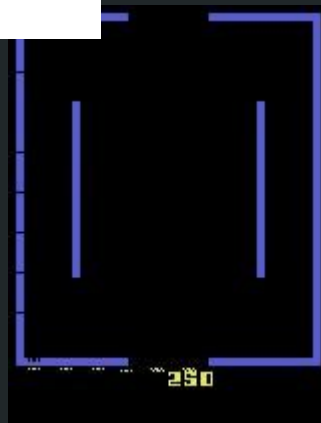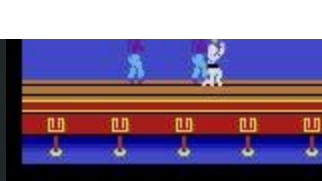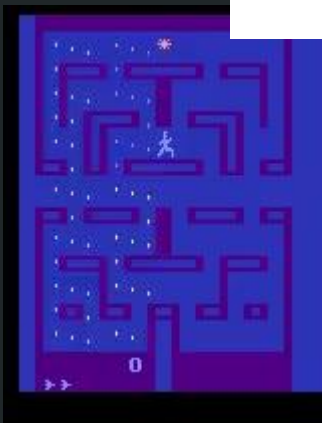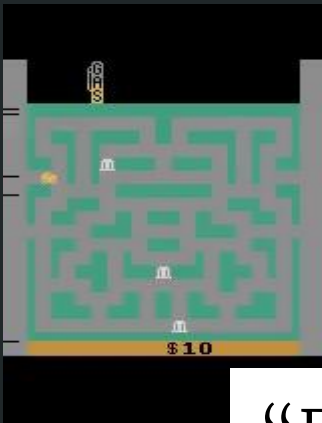Policy Based Gradient with CartPole

Deep Q Learning with Breakout

Future Work

OpenAI Gym

"Environments"

| observation | reward | done | action_space |
|---|---|---|---|
| Environment-specific object representing the observation of the environment.

Most likely the raw pixel data from the game. | Amount of reward achieved by the previous action. | Indicates whether a given episode has terminated. | The actions that can be performed at a given time. |

# CartPole

**observation**

Cart position, cart velocity, pole angle, and pole velocity at tip.

**reward**

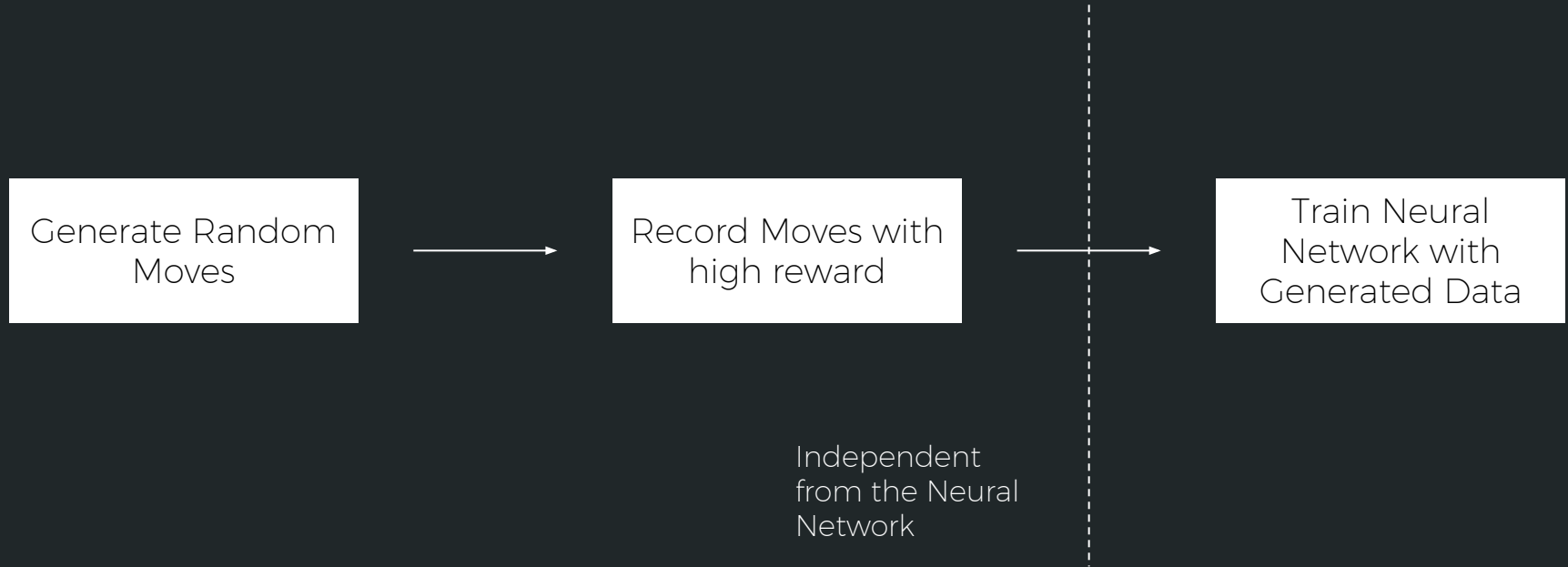Number of frames from the start episode to game termination

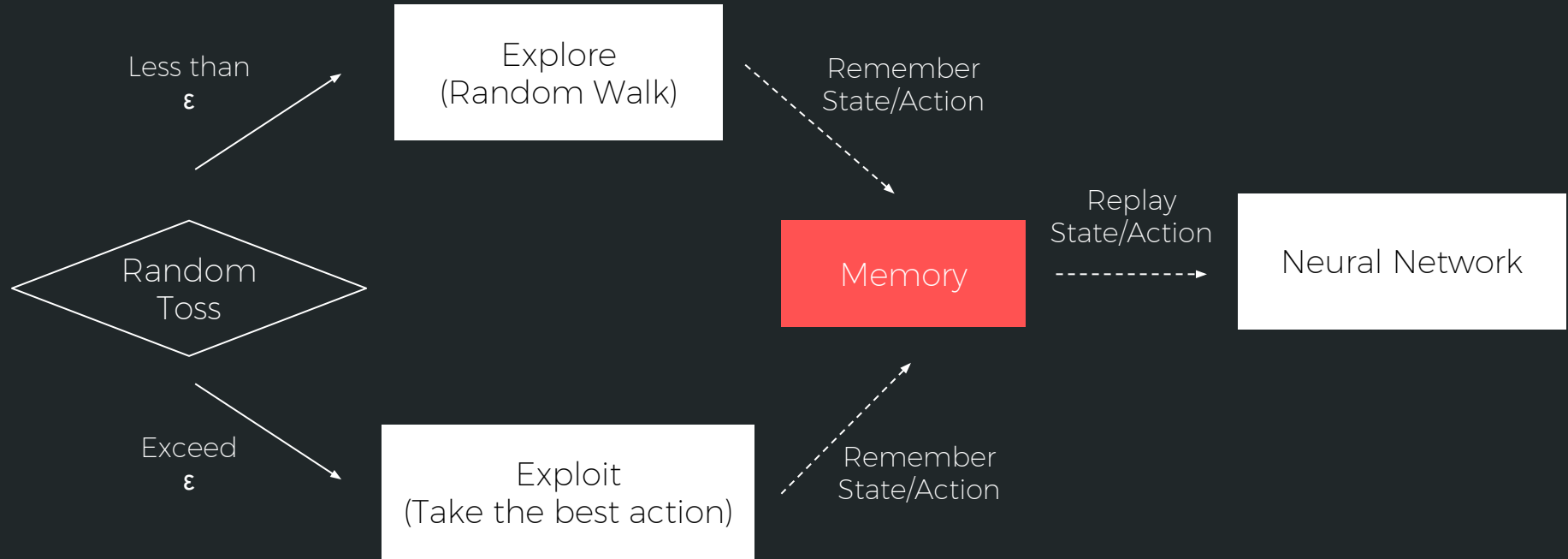**done**

The game is terminated if the pole falls off.

**action_space**

[-1, 1] indicating moving left or right.

# Traditional Neural Network

Generate Random Moves → Record Moves with high reward → Train Neural Network with Generated Data

Independent from the Neural Network

# Q Learning

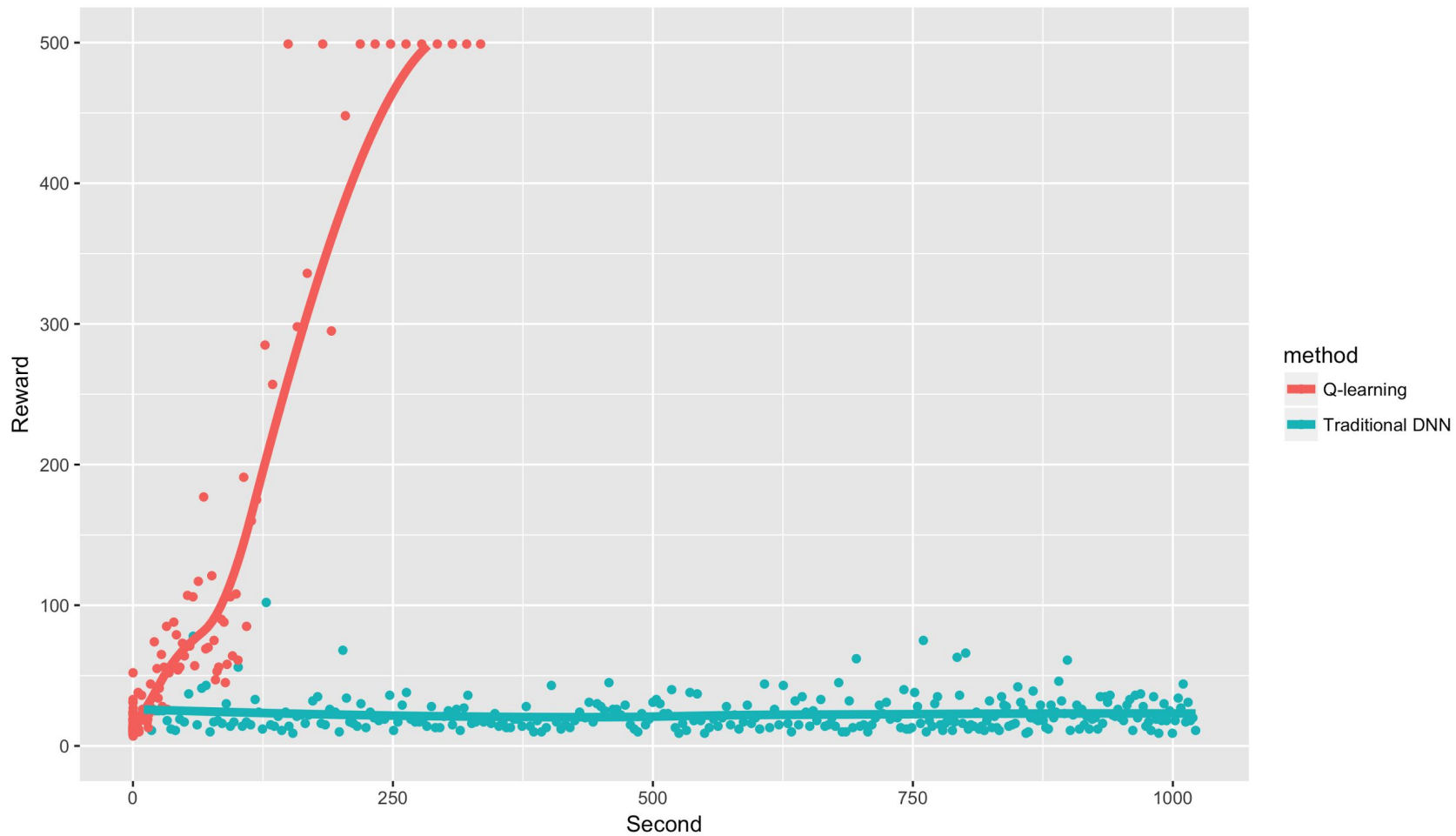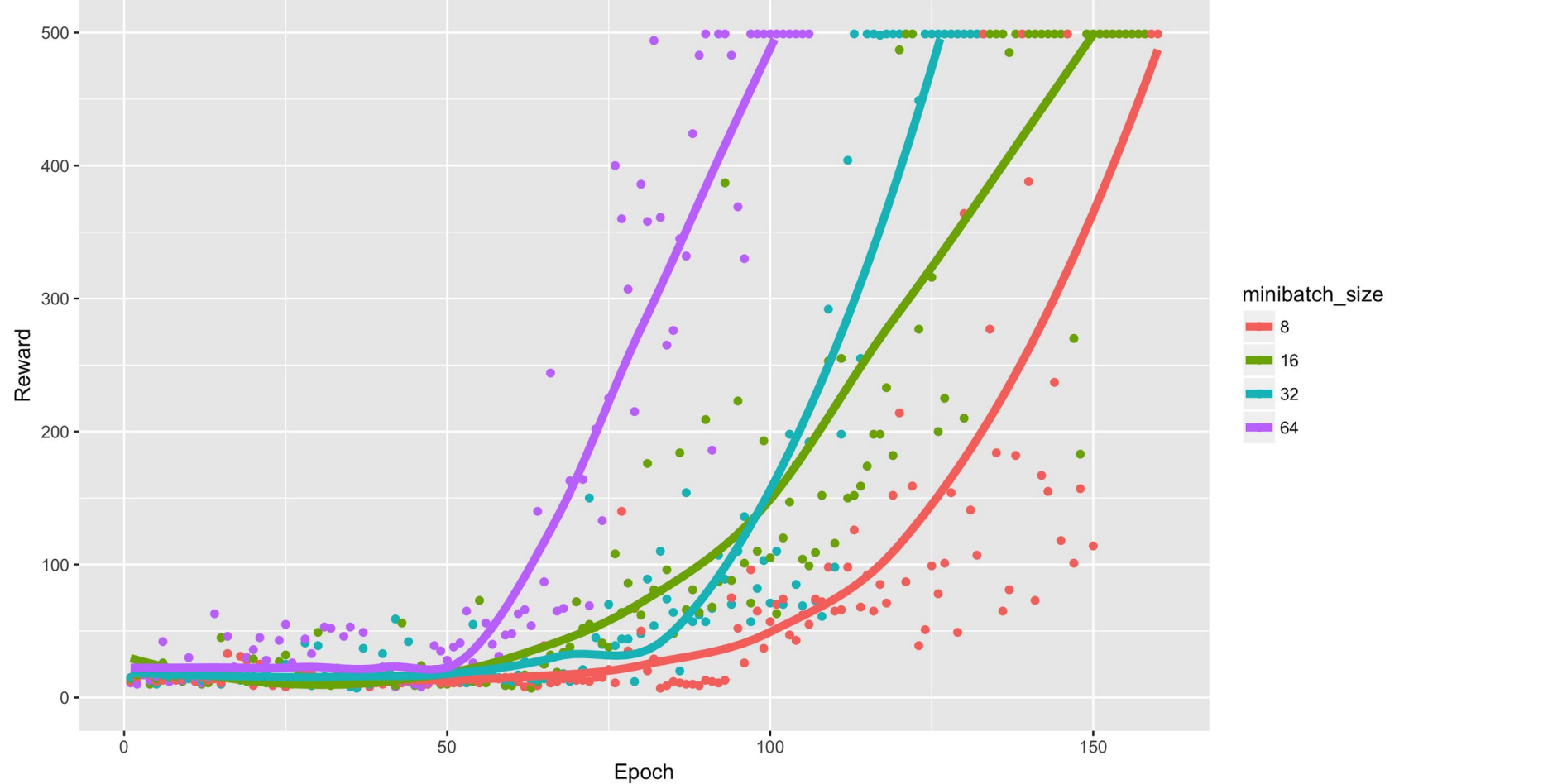| DNN | Q-Learning |
|:---:|:---:|
| 5 | 3 |
| Layers | Layers |
| 128, 256, 512 | 24 |
| Units / Layer | Units / Layer |
| Dropout Ratio: 0.2 | Dropout Ratio: 0 |
| Activation: RELU | Activation: RELU |

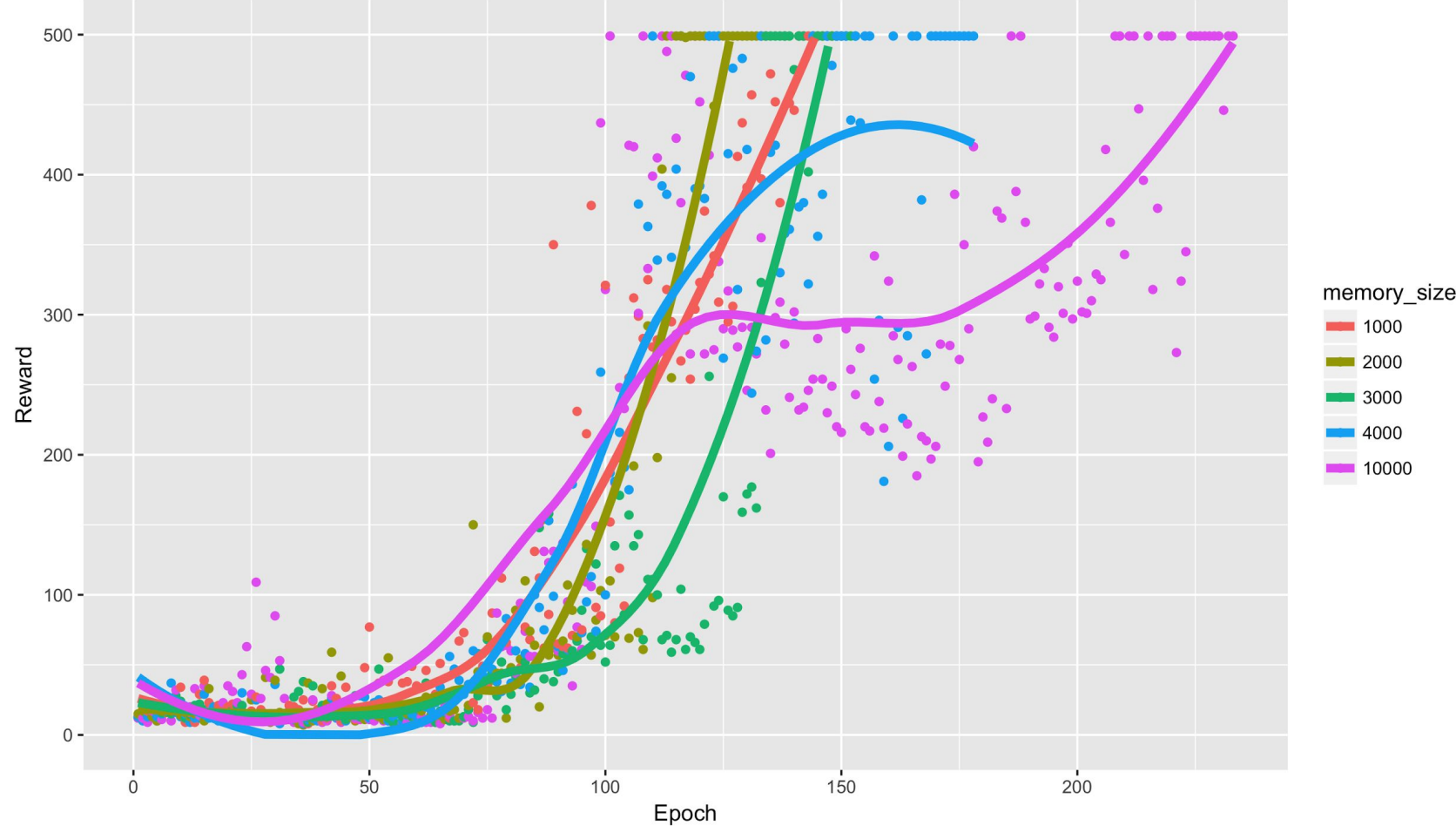Traditional DNN v.s. Q-learning

# Experiments
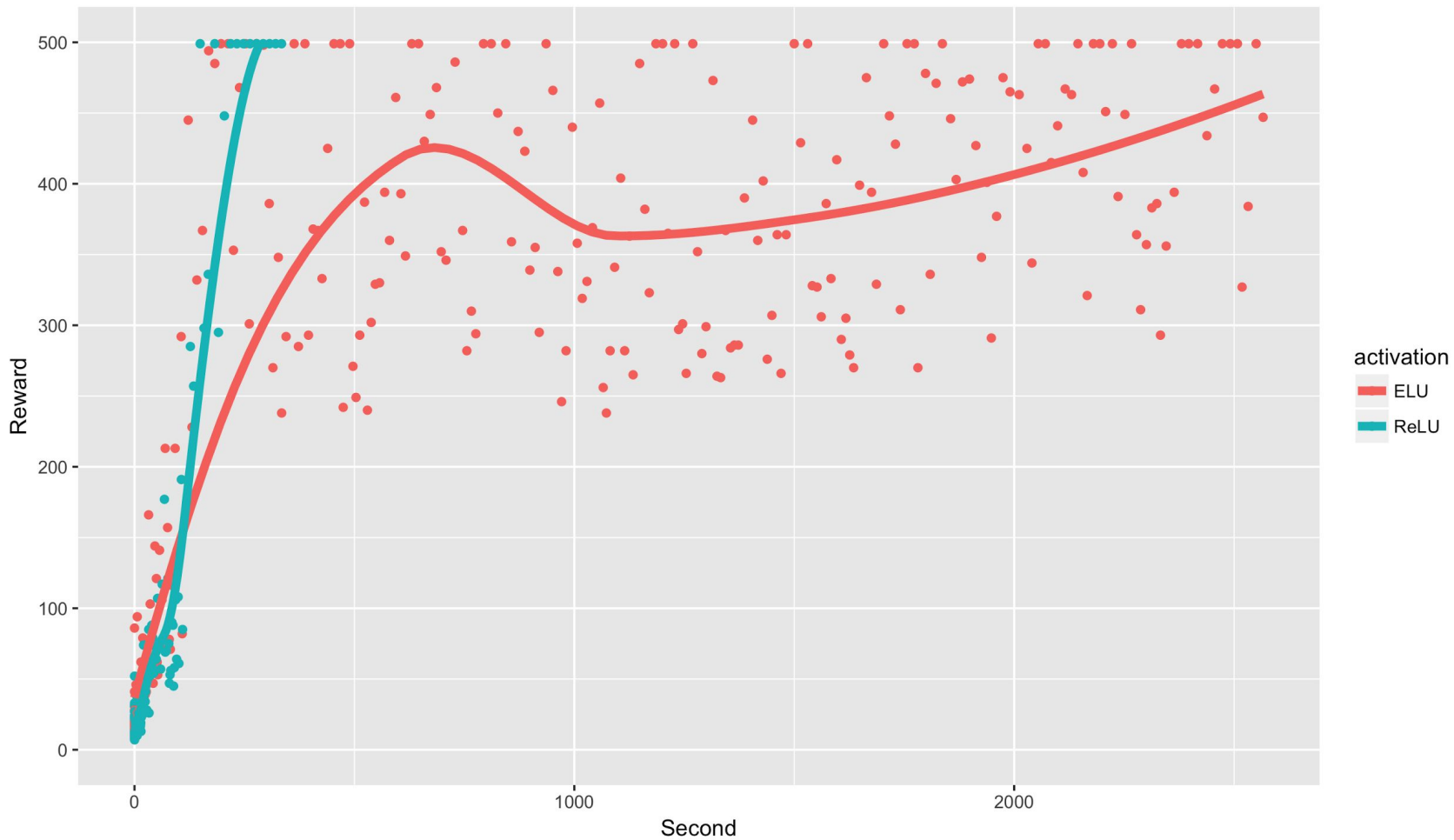
Reward v.s. Epoch with Different Minibatch Size

Reward v.s. Epoch with Different Memory Szie

# ELU Activation Function

$$f(x) = \begin{cases} x & \text{if } x \geq 0 \\ \alpha \left( \exp(x) - 1 \right) & \text{if } x < 0 \end{cases}$$
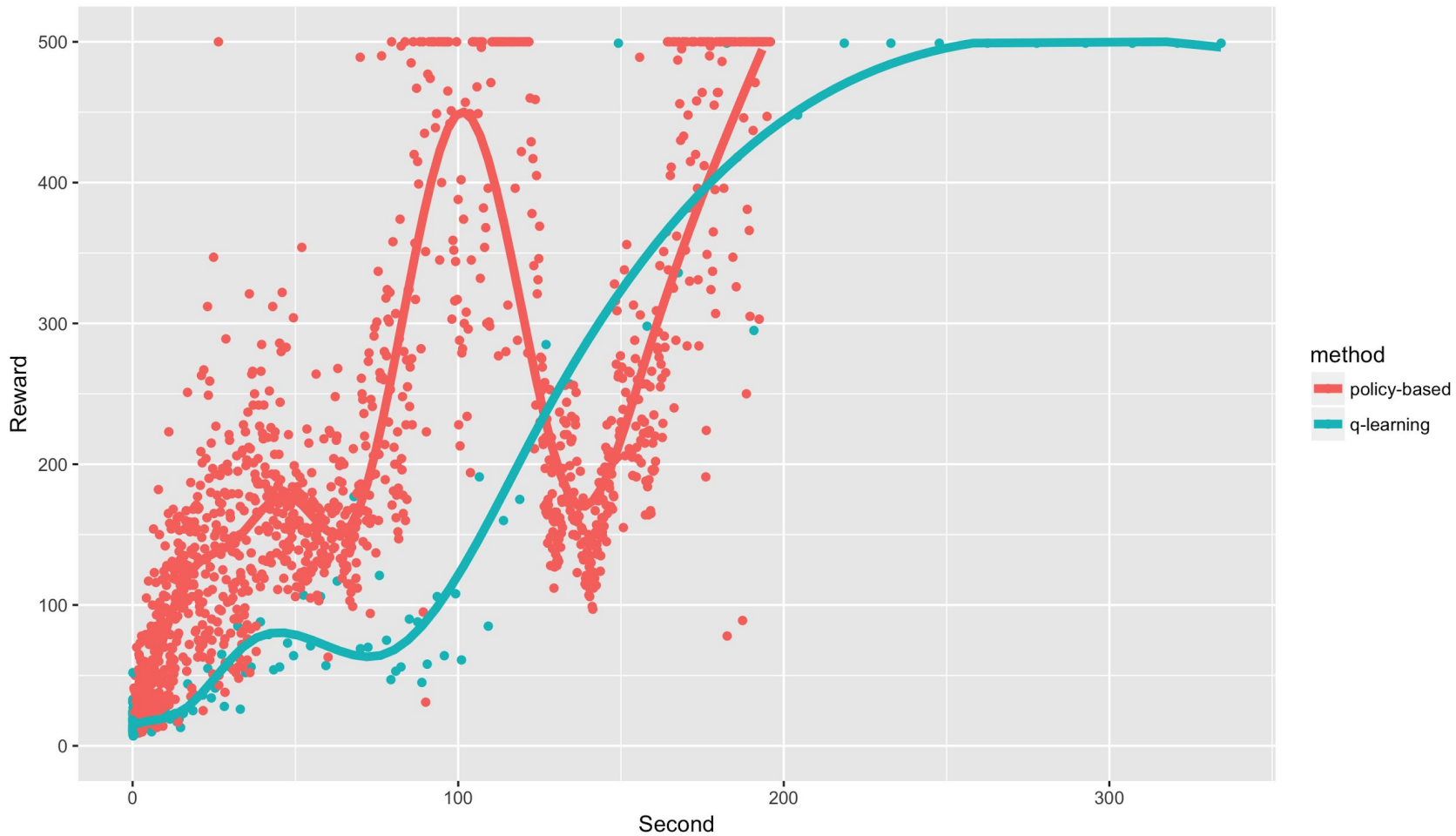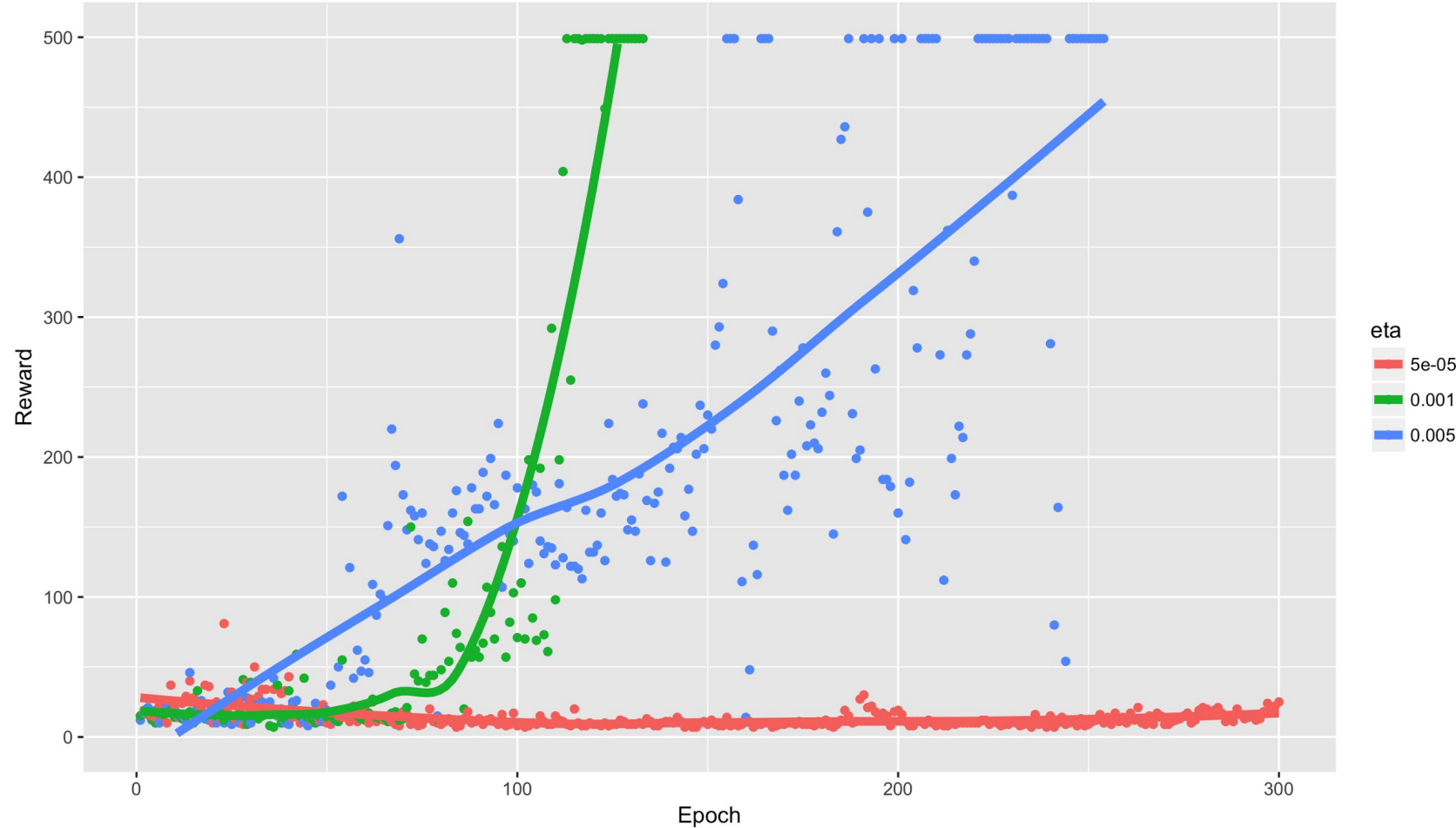
ReLU v.s. ELU

# Policy–Based Gradient

$$\theta_{t+1} = \theta_t + \alpha R_{t+1} \nabla_\theta \log \pi_{\theta_t}(A_t|S_t)$$

- ❖ Monte-Carlo Policy Gradient
- ❖ Some environments have easy policy and complex score
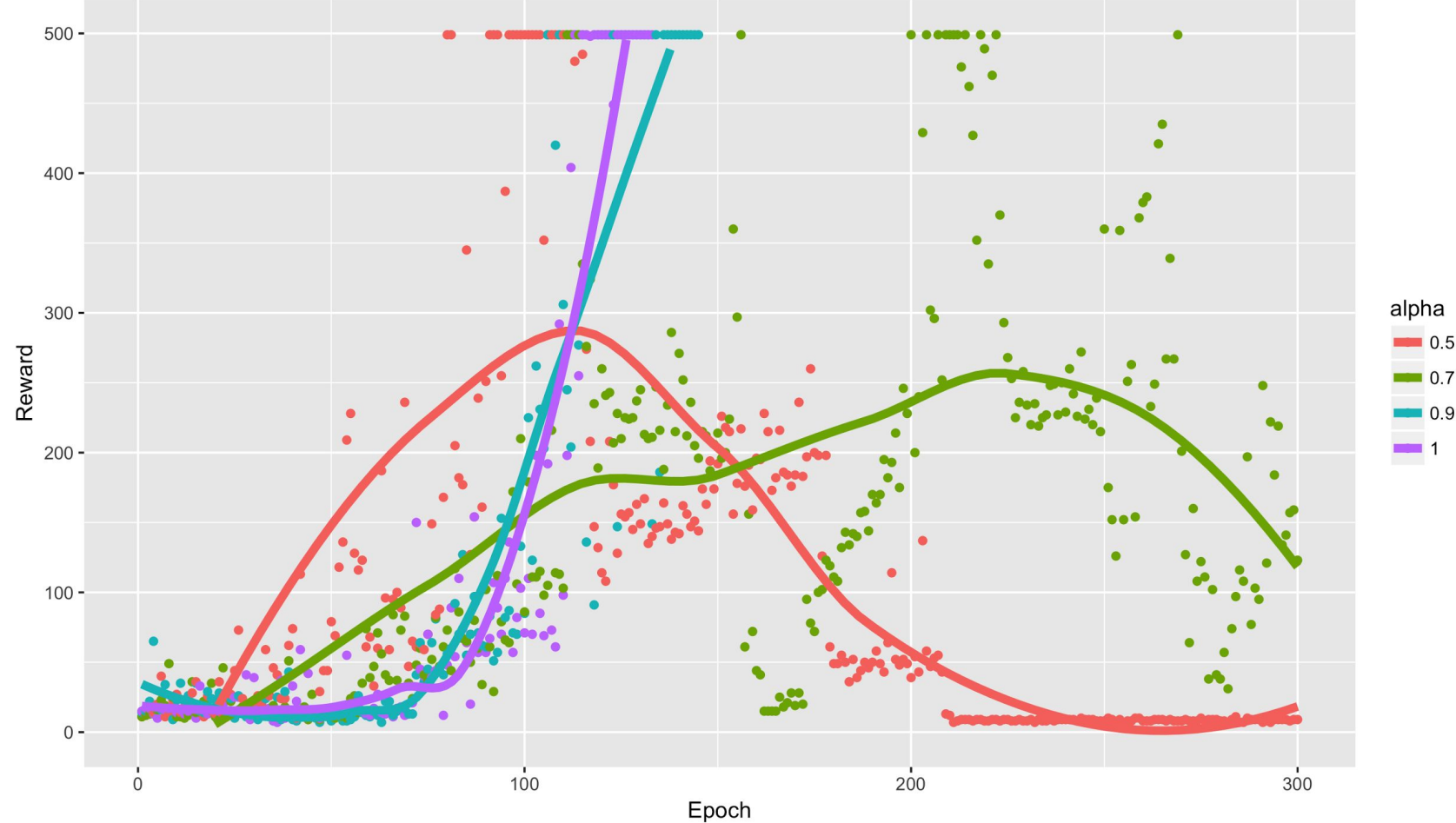- ❖ High variance

Q-learning v.s. policy-based Gradient

Reward v.s. Epoch with Different Learning Rate

# Use the Old Q Value...

$$y = (1 - \alpha)\hat{Q}(s_0, a) + \alpha \left( r + \gamma \max_{a'} \hat{Q}(s_1, a') \right)$$

Reward v.s. Epoch with Different Alpha

# Breakout

## observation

Raw pixel value of the game frame.
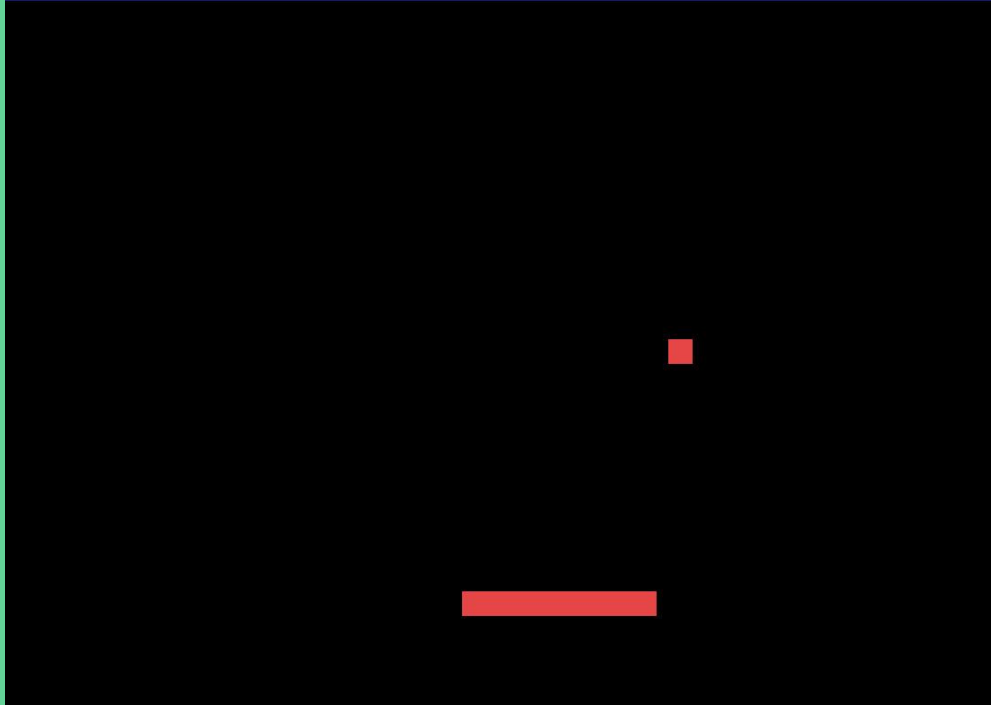
## reward

Scores gained by hitting the bricks.

## done

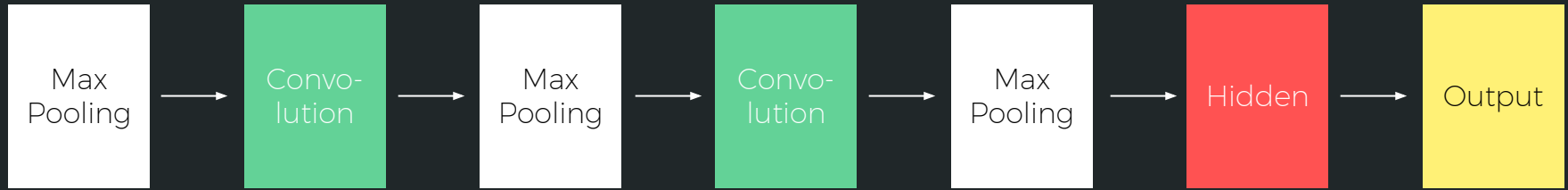The game is terminated if the ball falls off or all bricks are broken.

## action_space

[fire, fire left, fire right, left, right, noop]

# Q Learning with CNN

| Max Pooling | → | Convo-lution | → | Max Pooling | → | Convo-lution | → | Max Pooling | → | Hidden | → | Output |

Running on Google Cloud (with GPU) right now...

# Future Work

Generalize and Apply to other games

LSTM

Asynchronous Method