

2048 Deep Reinforcement Learning

Deep Learning Final Project
Tej Shah, Gloria Liu, Max Calero

Outline

- I. 2048 Game Introduction
- II. Simulation Results: Random, DDQN, and PPO Agents
- III. Background I: Value Estimation, DQN, DDQN
- IV. Background II: Policy Gradient, TRPO, PPO, GAEs
- V. State Representation to Model
- VI. 2048 DDQN Implementation Details
- VII. 2048 PPO Implementation Details
- VIII. Conclusion & Discussion of Results

I. Game Introduction

2048 Game Introduction

Motivation:

- Test Double Deep Q-Network (DDQN) and Proximal Policy Optimization (PPO)
- Game that is simple to understand, yet has a large state space & element of randomness

2048 Game Introduction

Move tiles left, right, up, and down

Current Score: 11284.0

Max Number: 1024.0

1024.0	256.0	32.0	16.0
8.0	8.0	64.0	0.0
2.0	2.0	0.0	0.0
0.0	0.0	4.0	0.0

Last Move: up



Current Score: 11304.0

Max Number: 1024.0

1024.0	256.0	32.0	16.0
0.0	0.0	16.0	64.0
0.0	2.0	0.0	4.0
0.0	0.0	0.0	4.0

Last Move: right

2048 Game Introduction

Current Score: 668.0

Max Number: 64.0

16.0	2.0	8.0	2.0
8.0	16.0	64.0	16.0
4.0	32.0	8.0	4.0
2.0	4.0	16.0	2.0

Last Move: right

Terminal State: Board is full, no action

Current Score: 20260.0

Max Number: 2048.0

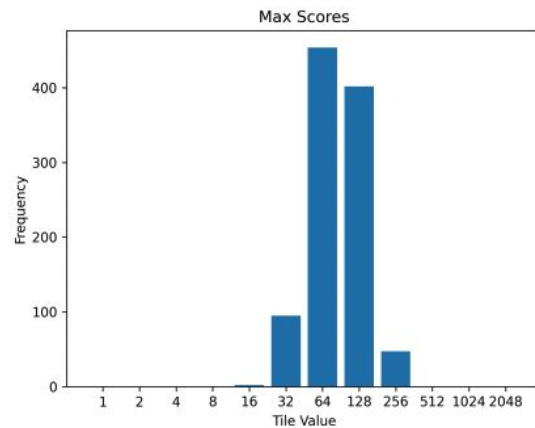
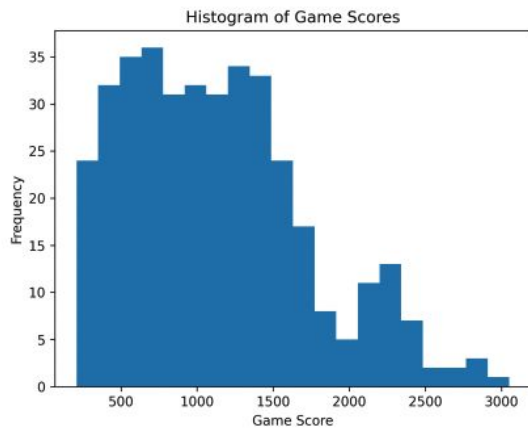
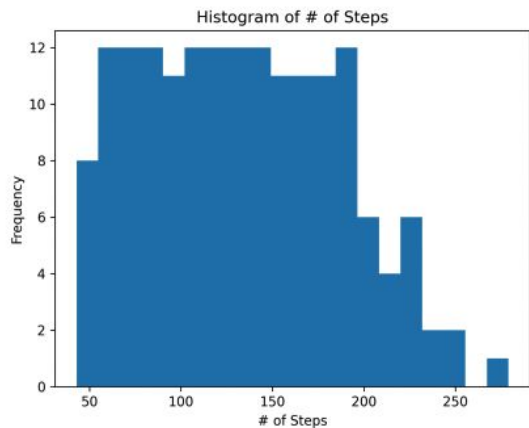
2048	0	0	2
8	16	2	0
2	4	2	0
2	0	0	0

Last Move: left

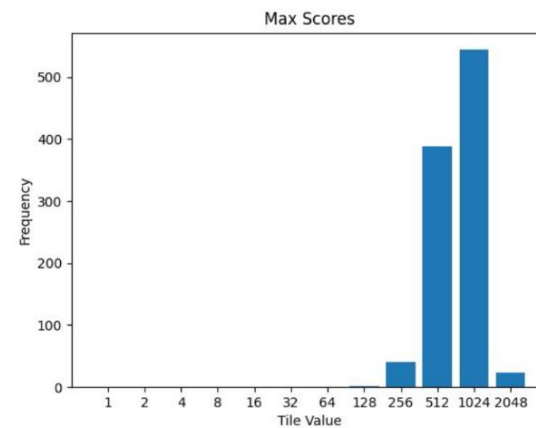
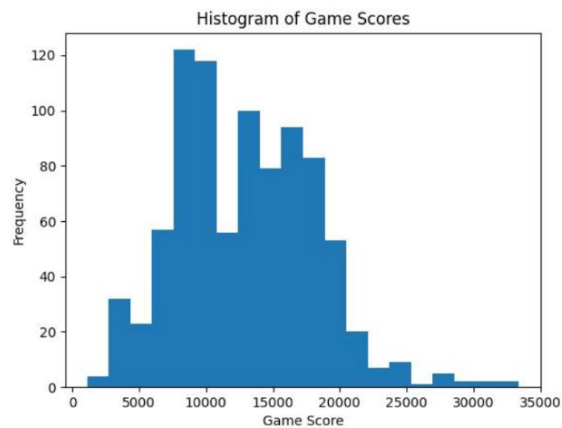
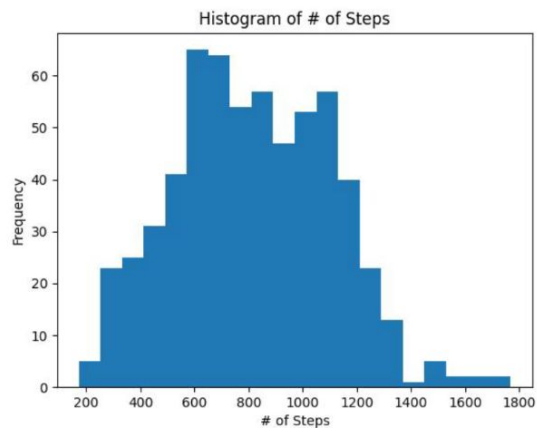
Win, can continue playing

II. Simulation Results - Random, DDQN, PPO Agents

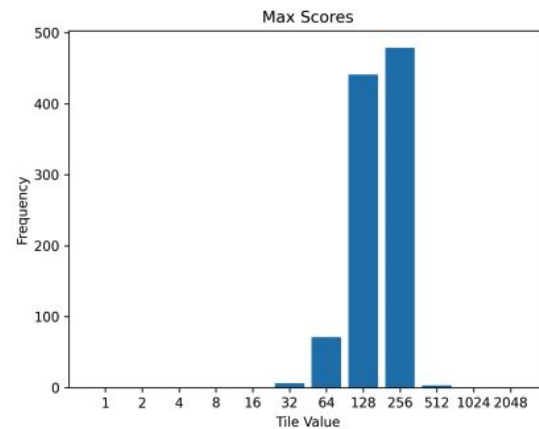
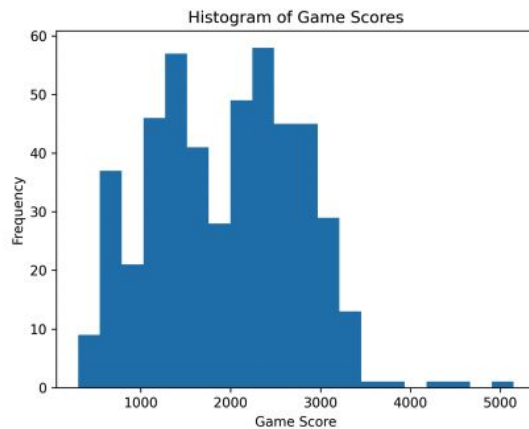
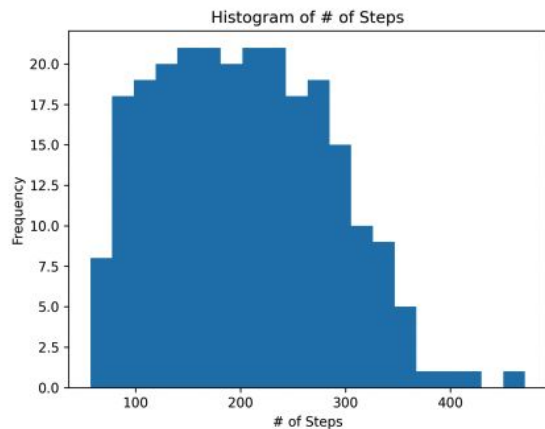
Simulation Results: Random Agent



Simulation Results: DDQN



Simulation Results: PPO



III. Background I: Value Estimation, DQN, DDQN

Bellman Equation

$$V^*(s) = \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')],$$

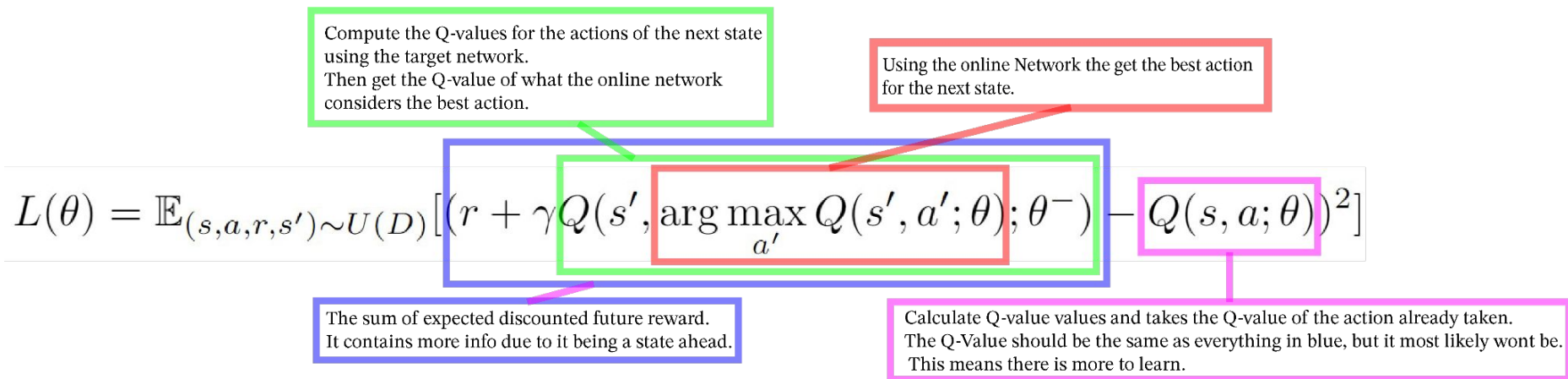
$$Q^*(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma \max_{a'} Q(s', a')],$$

Temporal difference

Y^{DDQN}

$$\delta_t = r + \gamma Q(s', \arg \max_{a'} Q(s', a'; \theta); \theta^-) - Q(s, a; \theta)$$

DDQN Loss



DDQN key concepts

- Q-Values are an estimation of expected reward for each action in a given state
- Temporal difference determines how much to learn from the sampled mini batch

IV. Background II: Policy Gradient, TRPO, PPO, GAEs

Policy Methods

$$\pi : \mathcal{S} \rightarrow \mathcal{A}$$

Main Goal: Update the parametrized probability distribution π over actions based on interaction with the environment.

Vanilla Policy Gradients

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) A^{\pi_{\theta}}(s_t, a_t) \right]$$

The diagram illustrates the components of the Vanilla Policy Gradient equation. The equation is $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) A^{\pi_{\theta}}(s_t, a_t) \right]$. The terms within the summation are highlighted in colored boxes with corresponding callouts:

- Red box:** ∇_{θ} is highlighted. Callout: "Update the policy over the trajectory samples in expectation".
- Blue box:** $\log \pi_{\theta}(a_t | s_t)$ is highlighted. Callout: "The log probability of action selection from policy."
- Yellow box:** $A^{\pi_{\theta}}(s_t, a_t)$ is highlighted. Callout: "The advantage of a particular action over the average."

Trust Region Policy Optimization (TRPO)

Relative “importance” ratio

maximize $\mathbb{E}_t \left[\frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)} A^{\pi_{\theta_{\text{old}}}}(s_t, a_t) \right]$

The advantage of a particular action over the average.

subject to

“Trust region”. You don’t want to update your new policy distribution too much from the old distribution.

$$\mathbb{E}_t [KL(\pi_{\theta_{\text{old}}}(s_t) || \pi_\theta(s_t))] \leq \delta$$

User specified bound for the trust region

Key Idea: “Constrain” policy updates to a “trust region.” However, it is an expensive computation practically.

Proximal Policy Optimization (PPO) KL Penalty

$$L^{\text{PEN}}(\theta) = \mathbb{E}_t [r_t(\theta) A_t - \beta_t KL(\pi_{\theta_{\text{old}}}(s_t) || \pi_{\theta}(s_t))]]$$

$$r_t(\theta) = \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)}$$

Key Idea:

Reformulate TRPO objective in a dual form and include KL penalty in the objective function instead of looking at a constrained optimization problem

Proximal Policy Optimization (PPO) Clipped

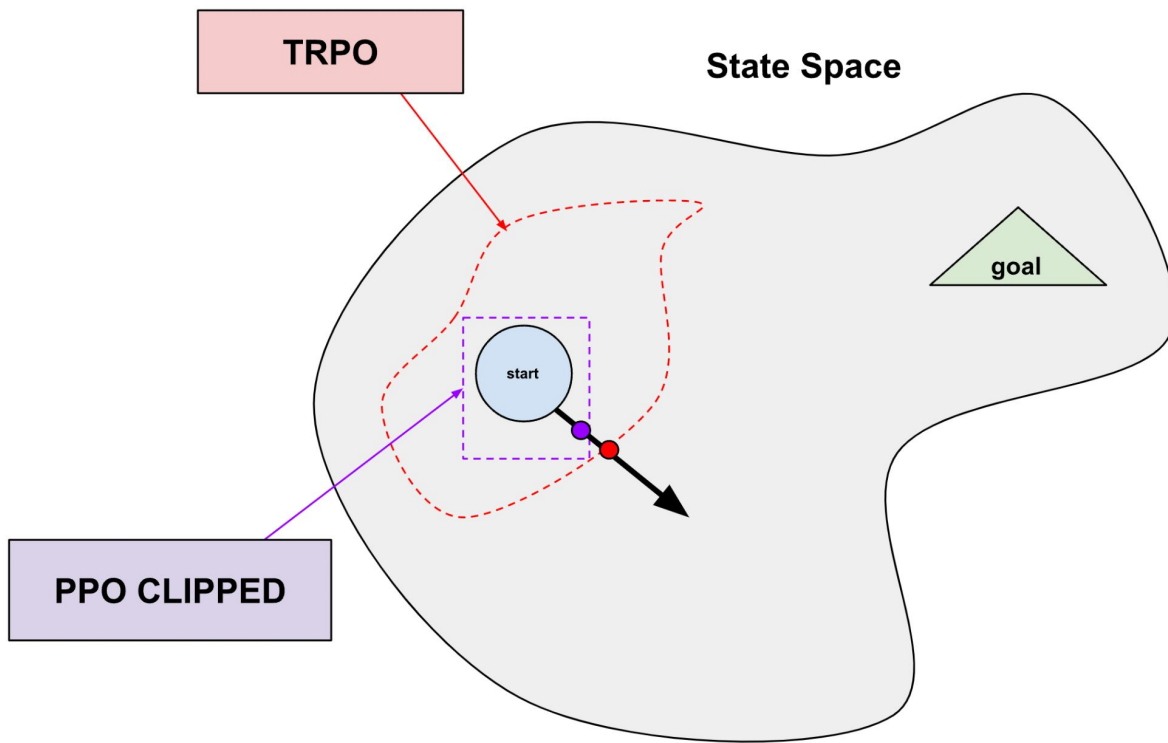
$$L^{\text{CLIP}}(\theta) = \mathbb{E}_t [\min (r_t(\theta) A_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) A_t)]$$

$$r_t(\theta) = \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)}$$

Key Idea:

Approximate TRPO and PPO Penalty behavior with pseudo surrogate objective that is easier to compute. We will use PPO clipped for 2048.

A Conceptual Picture: Policy Gradients, TRPO, PPO

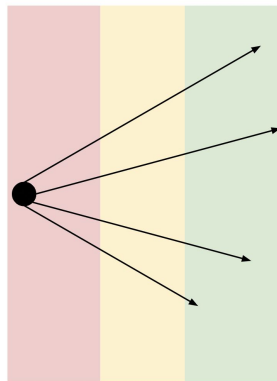


Generalized Advantage Estimation

- We want to score actions as advantageous based on some n-step lookahead, where earlier actions are prioritized more (lambda based exponential average)
- Further, we discount (gamma) future rewards for actions back to present time t

$$\hat{A}_t^{\text{GAE}(\gamma, \lambda)} = \sum_{l=0}^{\infty} (\gamma \lambda)^l \delta_{t+l}^V$$

$$\delta_t^V = r_t + \gamma V(s_{t+1}) - V(s_t)$$

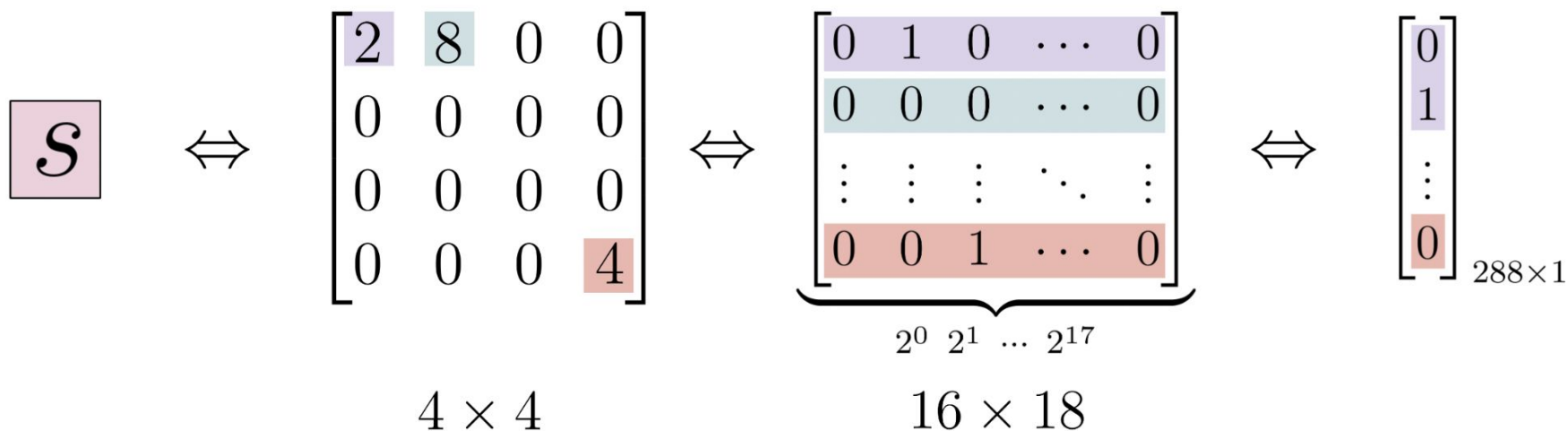


Increasing lambda:

- High bias
- Bias/variance tradeoff
- High variance

V. State Representation to Model

2048: State Representation (DDQN & PPO)

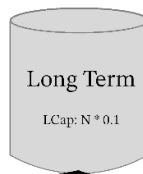
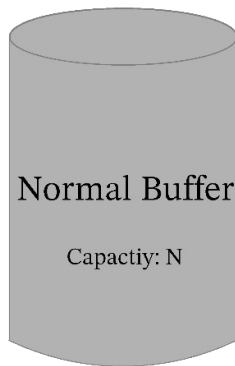


VI. 2048 DDQN Implementation Details

2048 DDQN: Implementation Details

Hyperparameters	Value
Replay Capacity	50000
Batch size	128
Epsilon start	0.9
Epsilon end	0.01
Epsilon decay	10000
Gamma	0.99
Tau	0.001
Learning Rate	1×10^{-5}
Number of Shared Layers	1
Shared Hidden Layer Size	256
Activation Function	nn.ReLU()
Gradient Clip Value	1000
Episodes	75000

U = DDQN Buffer



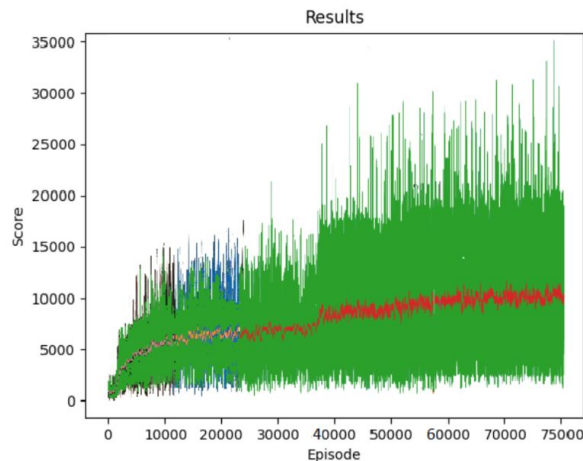
Stores tuples

$u_n = (\text{state}, \text{action}, \text{reward}, \text{next state}, \text{done})$

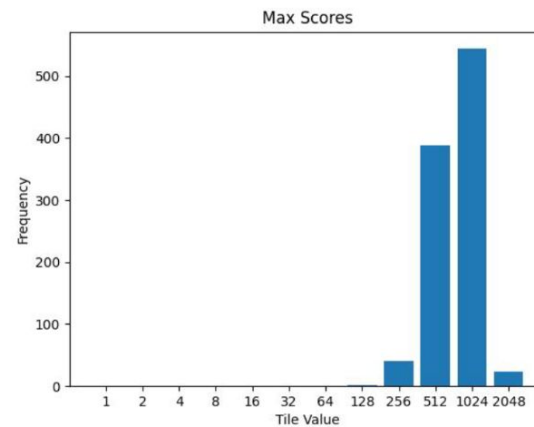
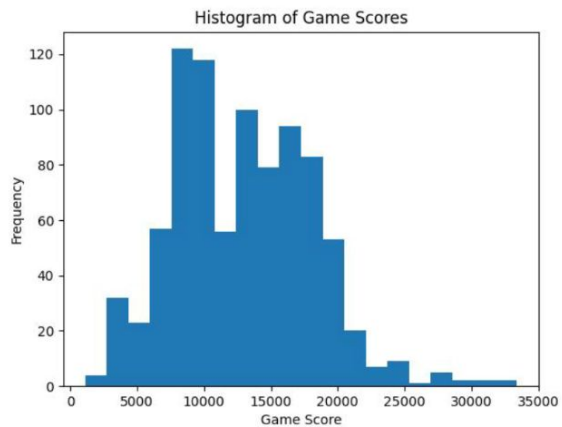
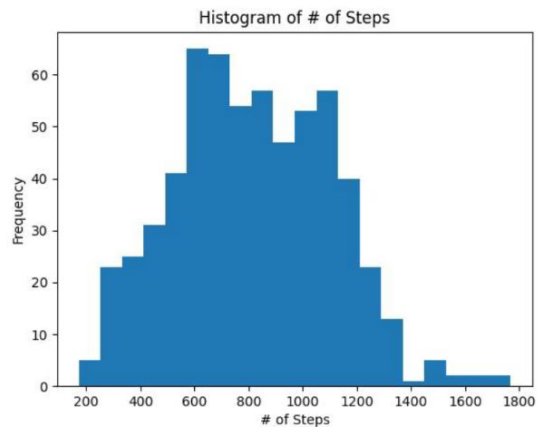
Normal_Buffer = $\{u_1, u_2, u_3, \dots, u_{LCap * 0.1}, \dots, u_{N/2}, u_{N/2 + 1}, \dots, u_{N/2 + LCap * 0.1}, \dots, u_N\}$

2048 DDQN: Implementation Results

Hyperparameters	Value
Replay Capacity	50000
Batch size	128
Epsilon start	0.9
Epsilon end	0.01
Epsilon decay	10000
Gamma	0.99
Tau	0.001
Learning Rate	1×10^{-5}
Number of Shared Layers	1
Shared Hidden Layer Size	256
Activation Function	nn.ReLU()
Gradient Clip Value	1000
Episodes	75000



2048 DDQN: Results



2048 DDQN Behavior

Current Score: 11476.0

Max Number: 1024.0

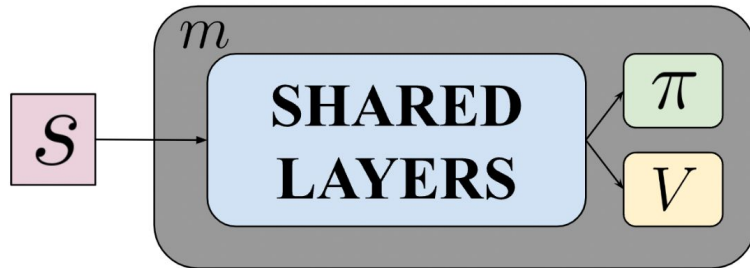
1024.0	256.0	32.0	16.0
32.0	64.0	8.0	4.0
8.0	16.0	2.0	2.0
4.0	4.0	0.0	0.0

Last Move: left

VII. 2048 PPO Implementation Details

2048 PPO: Implementation Details

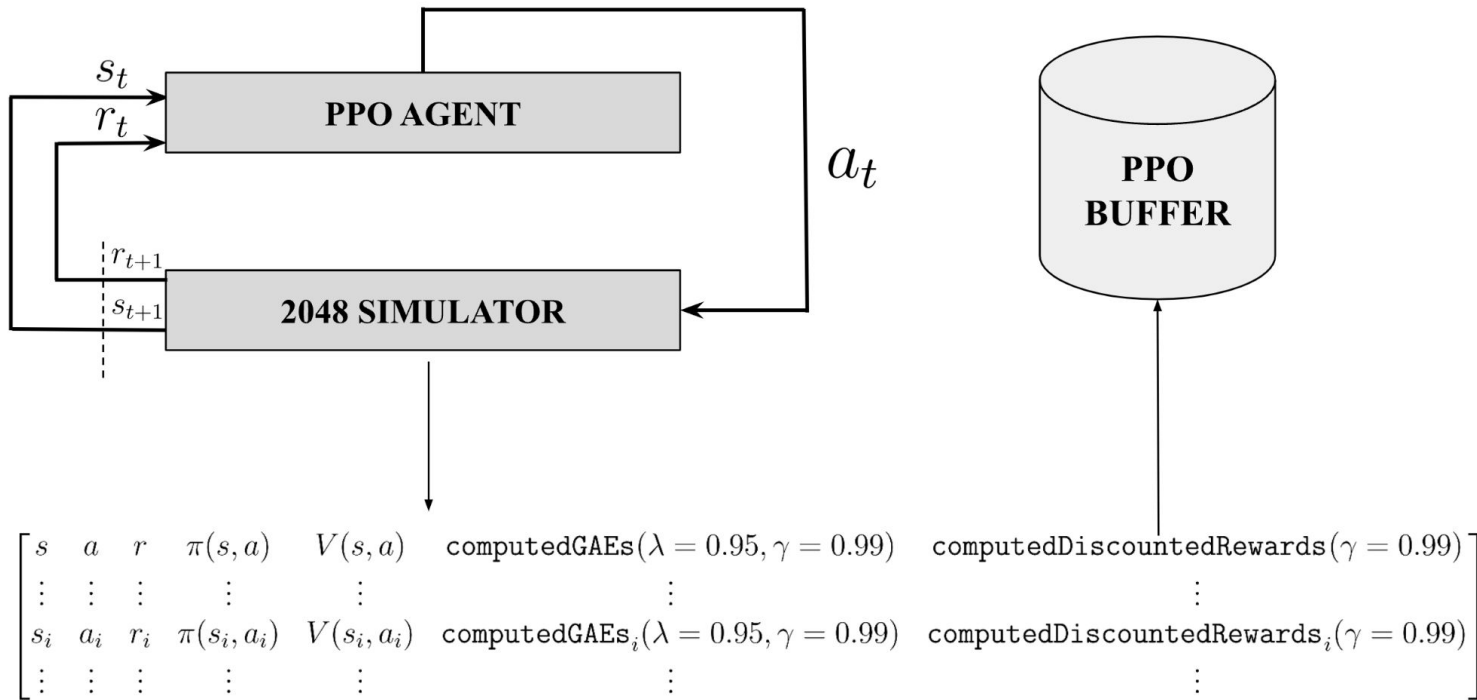
Hyperparameters	Value
Shared Hidden Layer Size	256
Number of Shared Layers	1
Activation Function	nn.Tanh()
PPO Clip Value	0.10
PPO Policy Learning Rate	1×10^{-5}
PPO Value Learning Rate	1×10^{-4}
PPO Epochs	60
Value Epochs	60
KL Target	0.02
Number of Rollouts	8/16/64



$$L^{\text{CLIP}} = \mathbb{E}_t [\min(r_t(\theta)A_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)A_t)]$$

$$L^{\text{Value}} = \frac{1}{N} \sum_{i=1}^N (r_i - V(s_i))^2$$

2048 PPO: Implementation Details



2048 PPO: Results

Hyperparameters	Value
Shared Hidden Layer Size	256
Number of Shared Layers	1
Activation Function	nn.Tanh()
PPO Clip Value	0.10
PPO Policy Learning Rate	1×10^{-5}
PPO Value Learning Rate	1×10^{-4}
PPO Epochs	60
Value Epochs	60
KL Target	0.02
Number of Rollouts	8/16/64

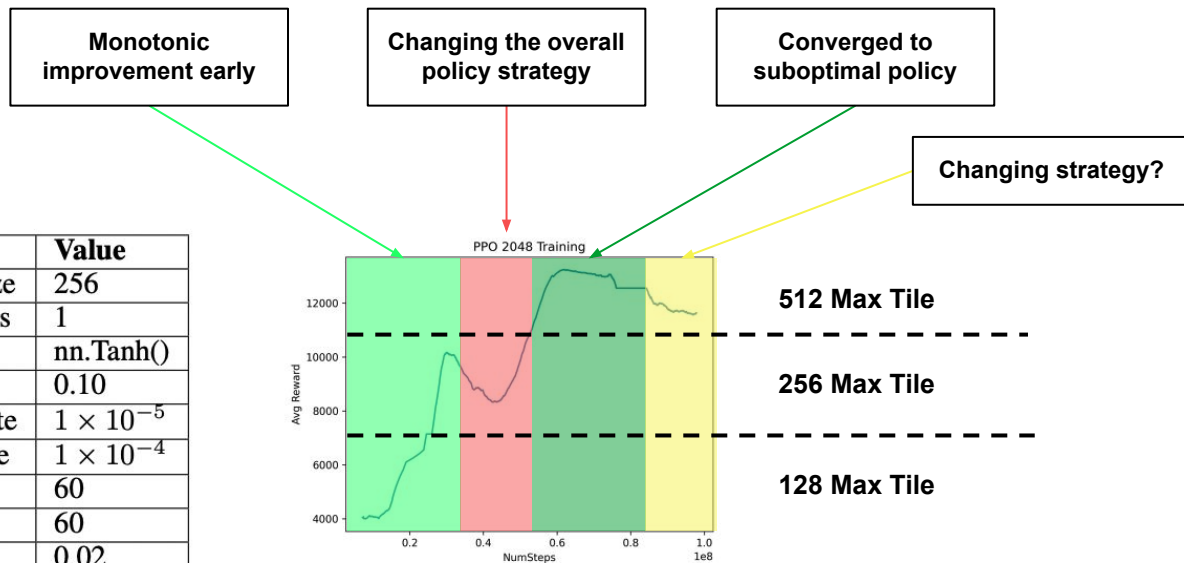
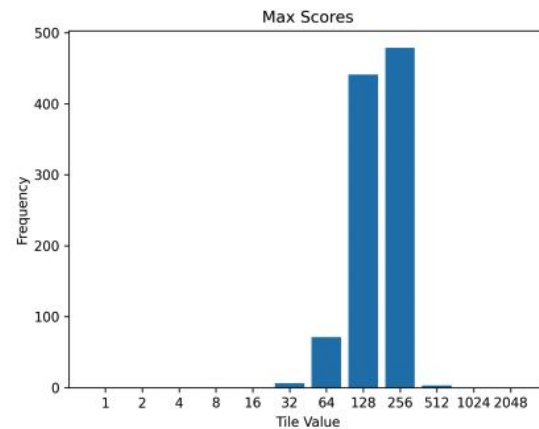
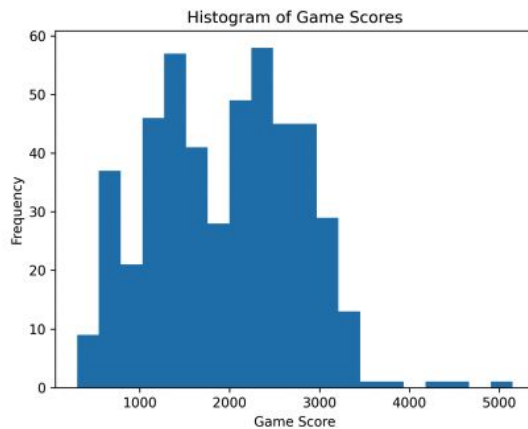
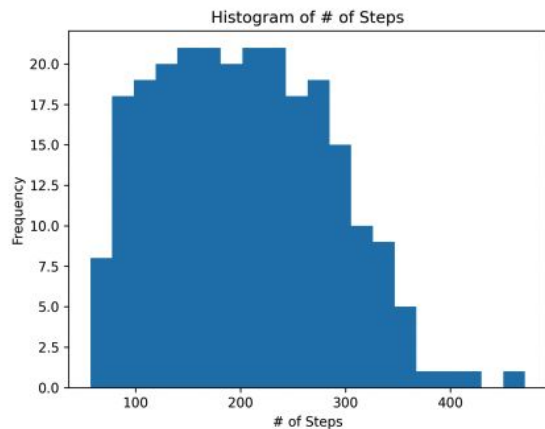


Figure 13: $r_{\text{final}} = \log_{\text{MaxTile}} \frac{1}{t} \cdot \text{deltaScore}$

2048 PPO: Results



2048 PPO Behavior

Current Score: 0.0

Max Number: 2.0

0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
2.0	0.0	0.0	0.0
0.0	2.0	0.0	0.0

Last Move: no move

VIII. Conclusion & Results Discussion

Conclusion

Results

- Both DDQN and PPO better than random - further study and experimentation with these methods
- DDQN more sample efficient than PPO for 2048

Takeaways

- Importance of reward functions and state representation
- Converging to suboptimal results