

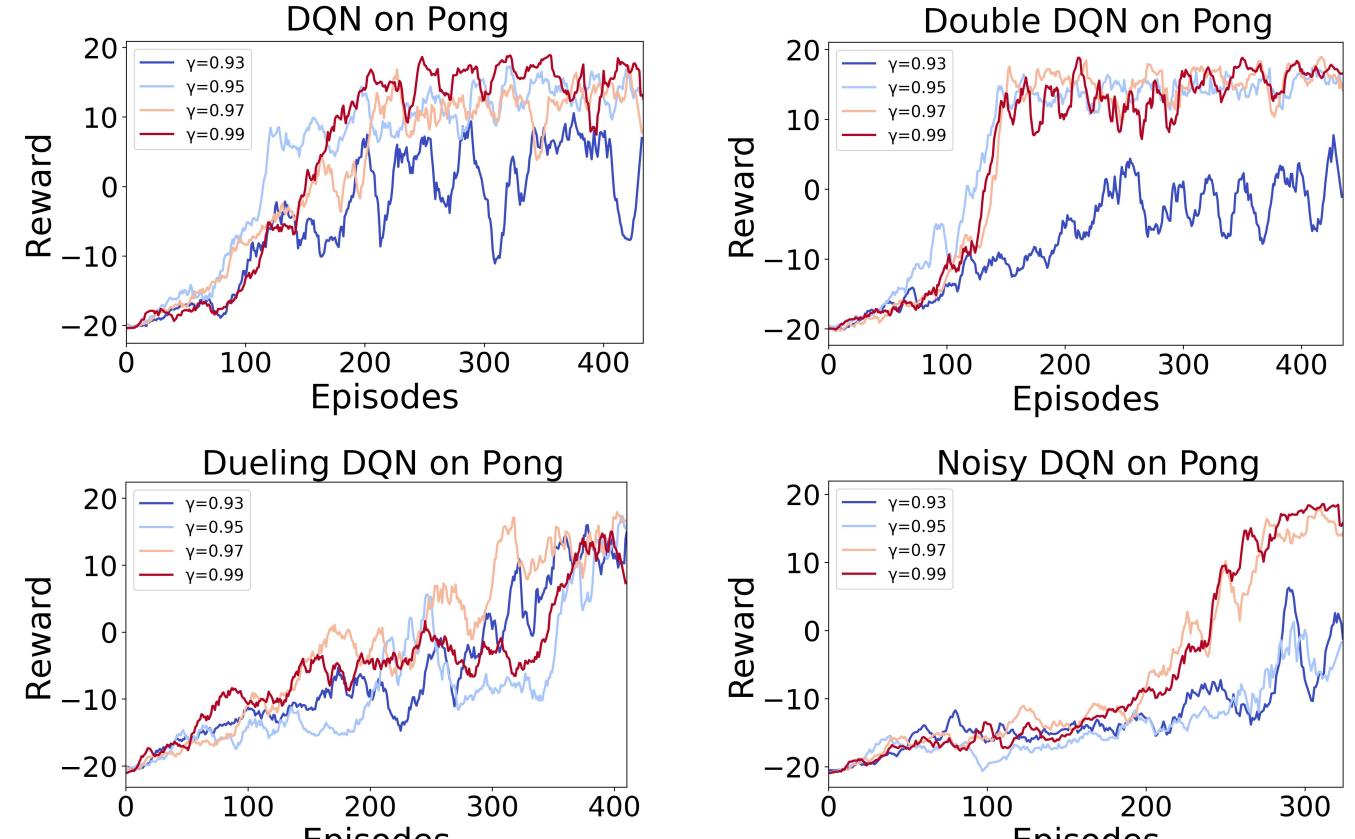


Aditya Golatkar<sup>1</sup> (UID: 505221372), Albert Zhao<sup>1</sup> (UID: 605231634), Cagatay Isil<sup>2</sup> (UID: 705430690), Xiaoran Zhang<sup>2</sup> (UID: 605429606)

<sup>1</sup> Department of Computer Science, <sup>2</sup> Department of Electrical and Computer Engineering

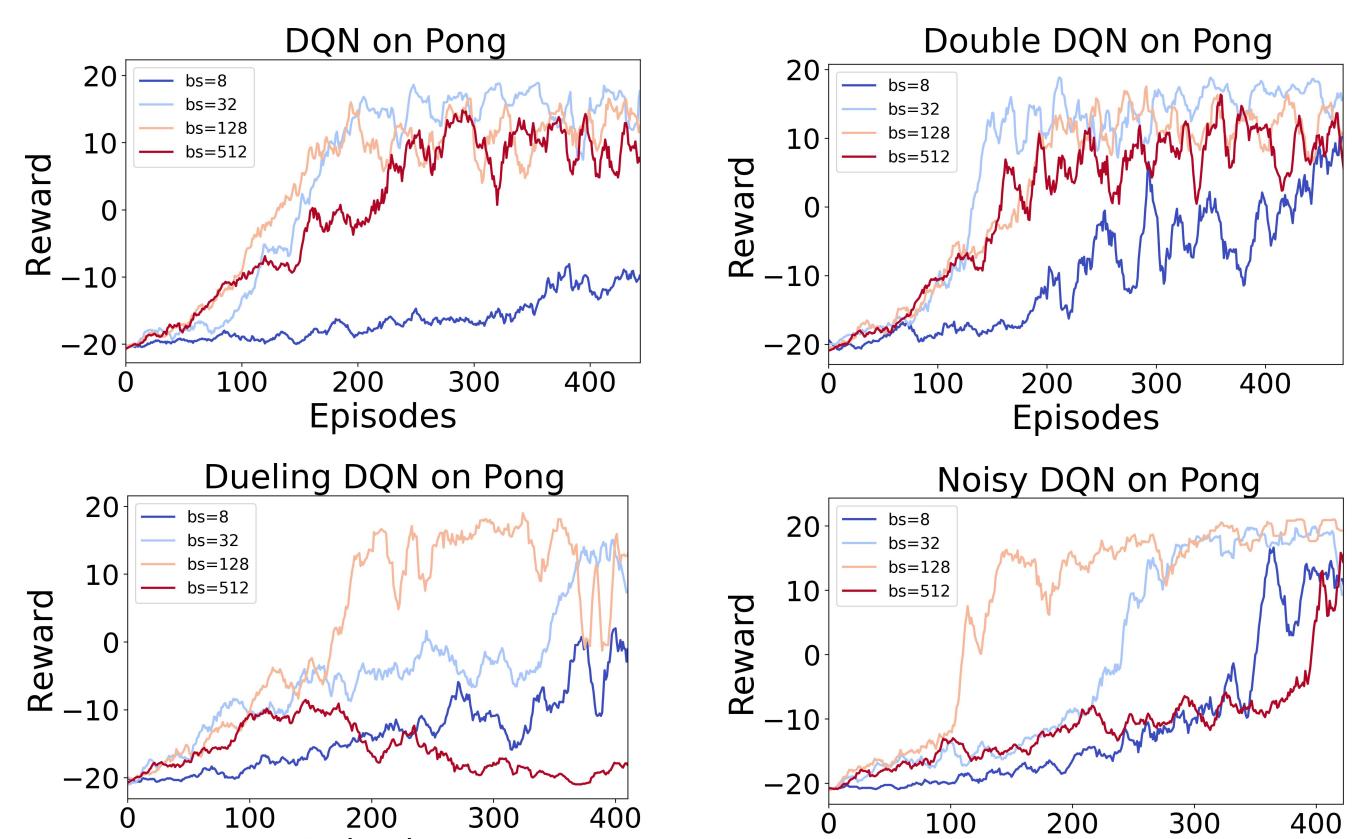
## Effective Horizon

Agents with wider effective horizon ( $=1/(1-\gamma)$ ) learn better policies.

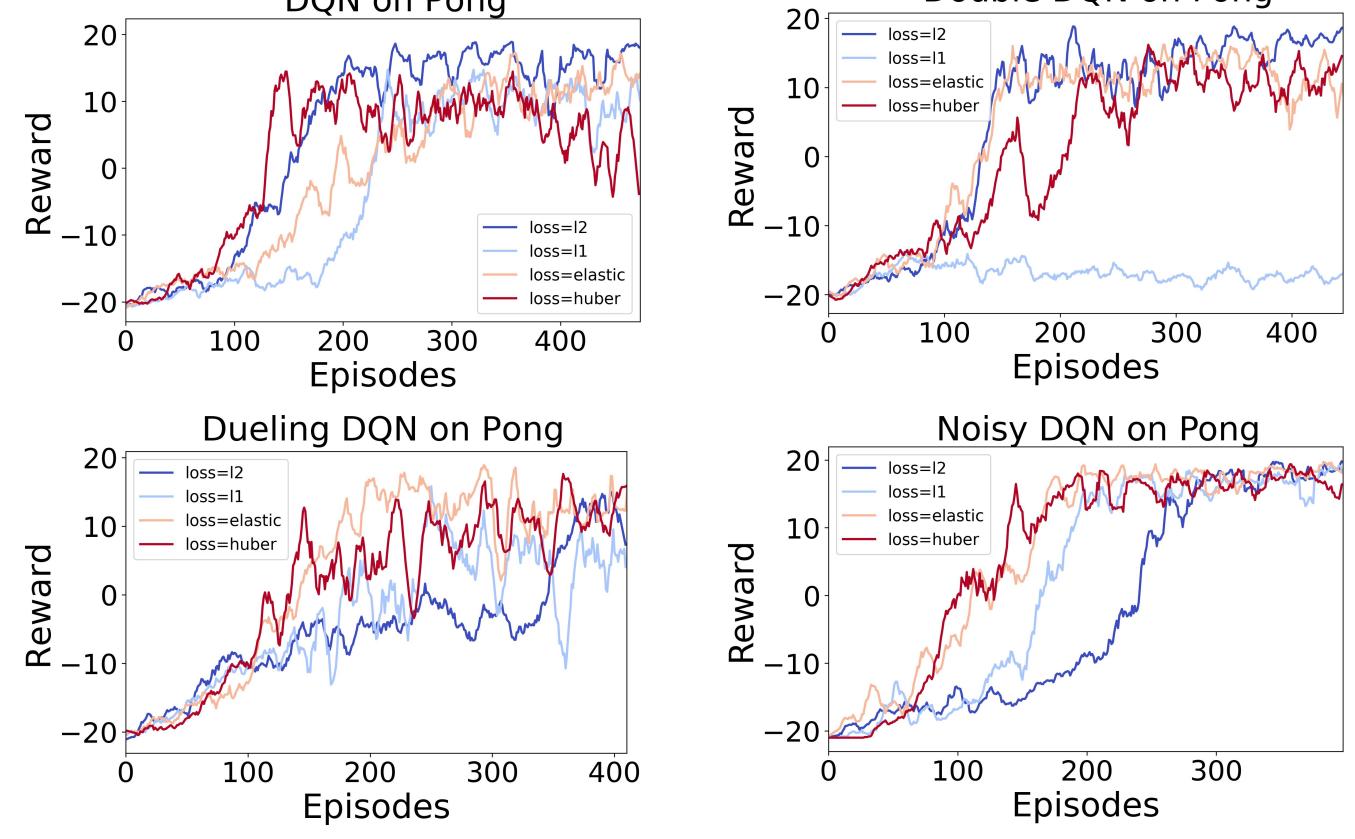


## Batch-Size

Low batch-size (or high gradient noise) results in slower convergence.



## Loss Function



## Deep Q-Networks and Variants

### DQN[1]

$$Y_t^{DQN} = R_{t+1} + \gamma \max_a Q(S_{t+1}, a; (\theta_t^-, w_t^-))$$

### Double DQN[2]

$$Y_t^{DoubleDQN} = R_{t+1} + \gamma Q(S_{t+1}, \text{argmax}_a Q(S_{t+1}, a; (\theta_t, w_t)); (\theta_t^-, w_t^-))$$

### Dueling DQN[3]

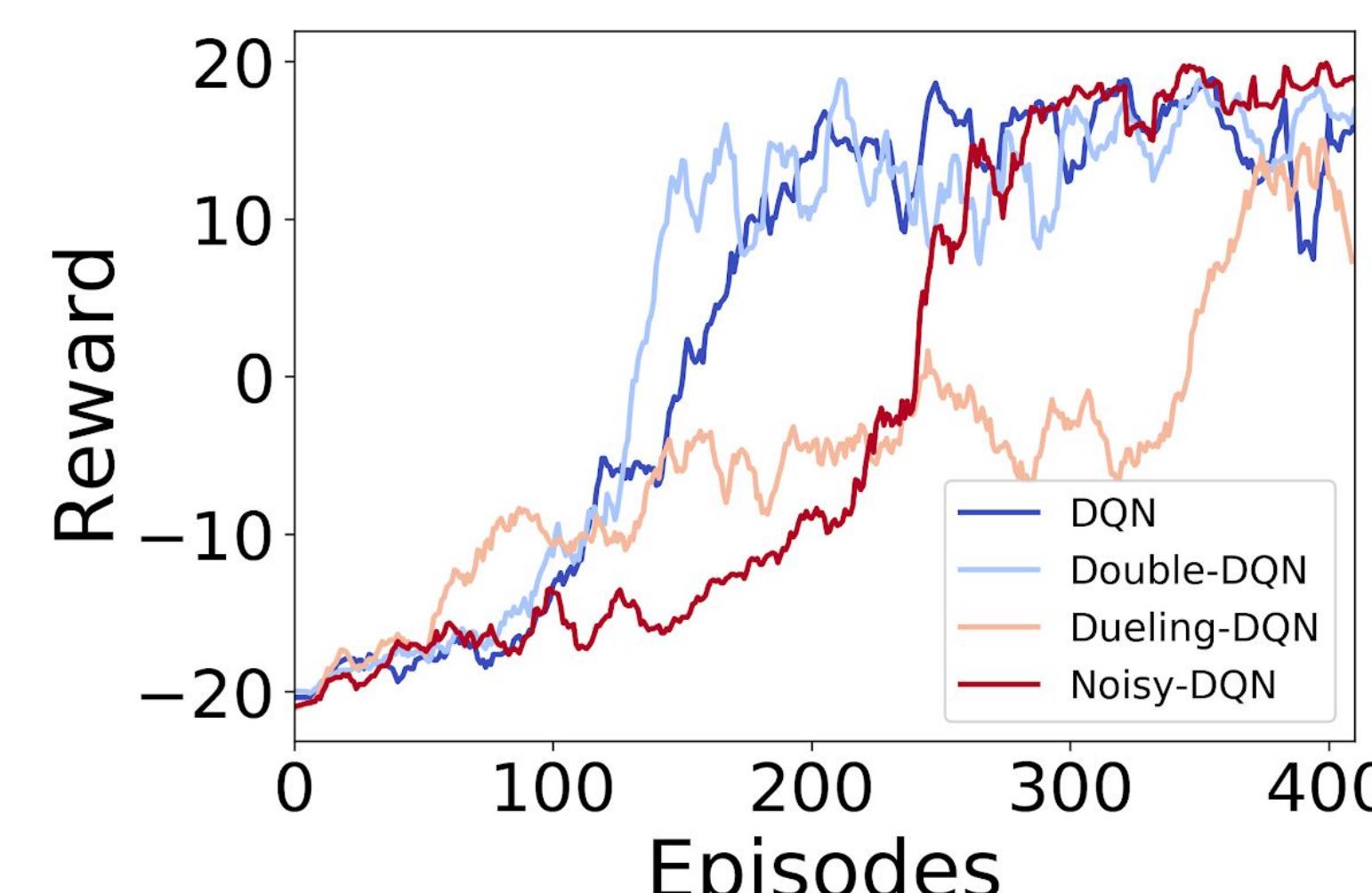
$$Y_t^{DuelDQN} = R_{t+1} + \gamma \max[V(S_{t+1}; (\theta_t^-, \beta_t^-) + \frac{1}{|\mathcal{A}|} \sum_a A(S_{t+1}, a; (\theta_t^-, \alpha_t^-))]$$

### Noisy DQN[4]

$$Y_t^{NoisyDQN} = R_{t+1} + \gamma \max_a Q(S_{t+1}, a; (\theta_t^-, \mu^{w_t^-} + \sigma^{w_t^-} \odot \epsilon)) \quad \epsilon \sim \mathcal{N}(0, I)$$

Mean of FC layers  
Standard deviation of FC layers  
Noise with normal distribution

DQN variants on Pong

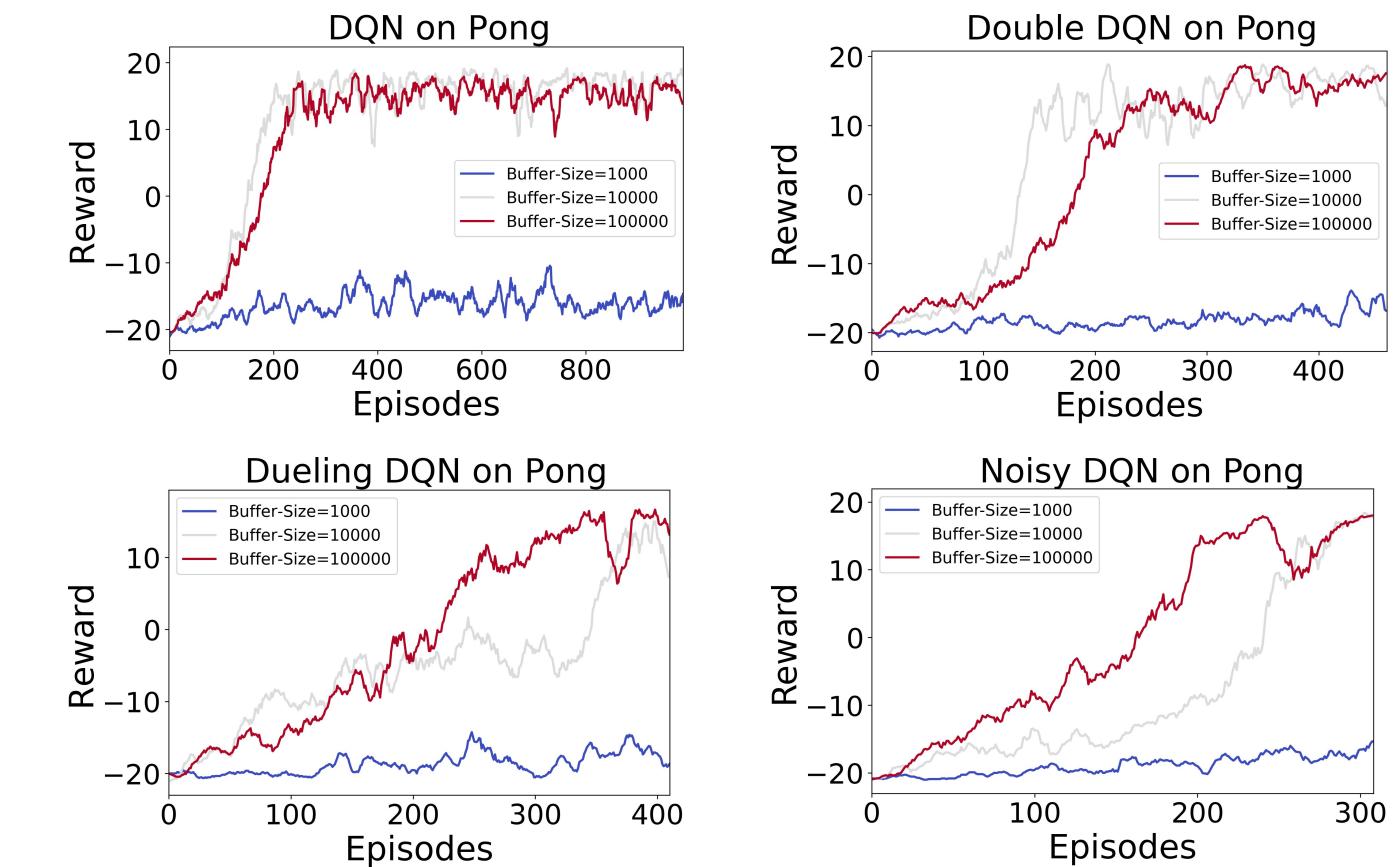


## References

- [1] Volodymyr Mnih et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- [2] Hado Van Hasselt et al. Deep reinforcement learning with double q-learning. In Thirtieth AAAI Conference on Artificial Intelligence, 2016.
- [3] Ziyu Wang et al. Dueling network architectures for deep reinforcement learning. *arXiv preprint arXiv:1511.06581*, 2015.
- [4] Meire Fortunato et al. Noisy networks for exploration. In International Conference on Learning Representations, 2018.
- [5] Tom Schaul et al. Prioritized experience replay. In International Conference on Learning Representations, Puerto Rico, 2016.

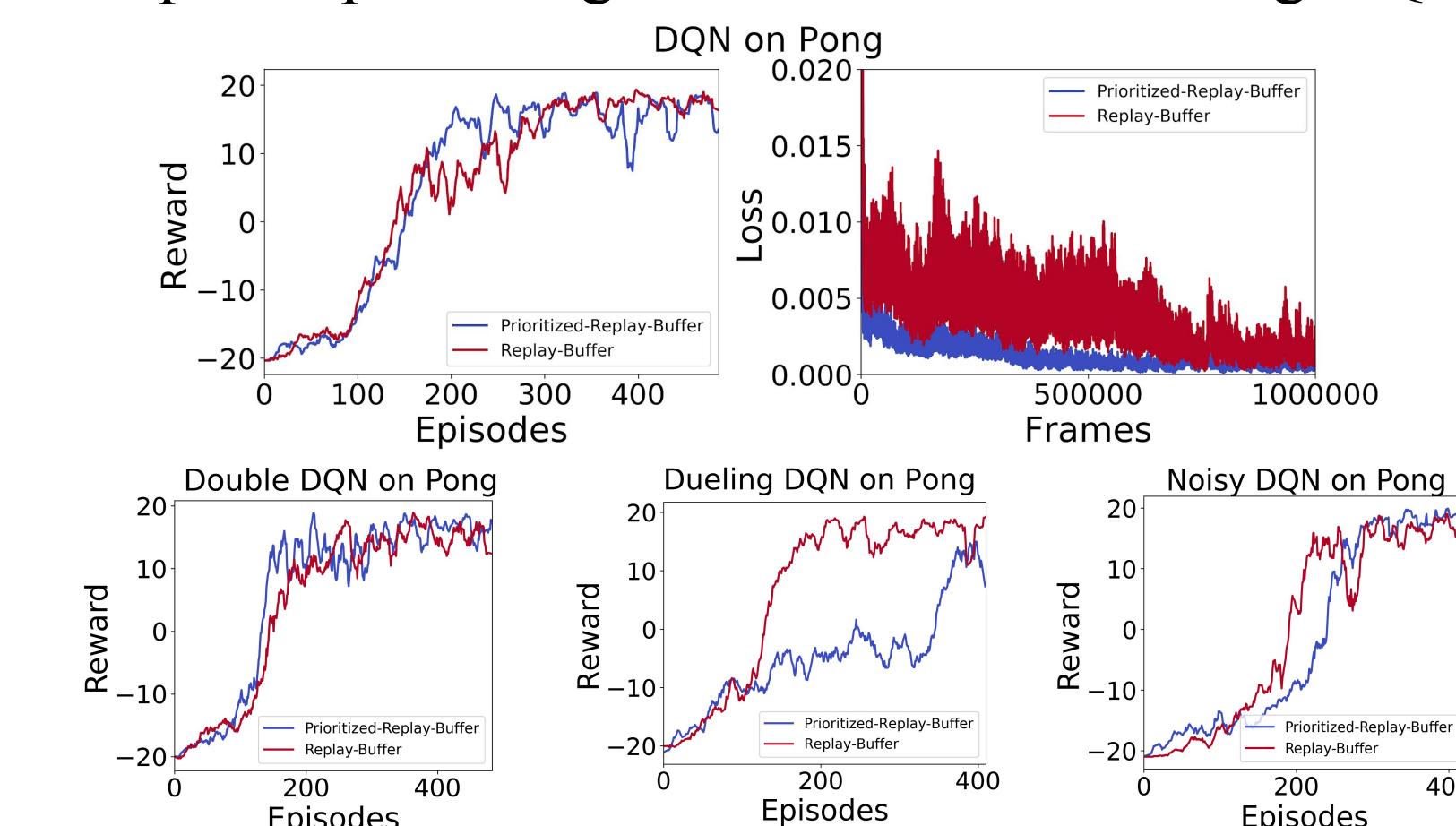
## Replay Buffer Size

Low replay buffer size leads to slow convergence due to inability to effectively reuse data samples.



## Role of Priority

Prioritization<sup>[5]</sup> reduces loss variance but does not speed up convergence & hurts for Dueling DQN.



## Exploration

Over-exploration results in slower convergence.

